

An Efficient and Usable Client-Side Cross Platform Compatible Phishing Prevention Application

FINAL YEAR PROJECT REPORT

Submitted by

N. Dhanush(2016103021)

G. Santhosh(2016103057)

S. Ben Stewart(2016103513)

Guide

Dr. Angelin Gladston

Associate Professor

Department of Computer Science & Engineering

College of Engineering, Guindy



ANNA UNIVERSITY: CHENNAI 600 025

JANUARY 2020

1 PROBLEM STATEMENT

Phishing is a crime where the victim is contacted by an attacker posing as a trustworthy source and lure them into providing sensitive information like credit card details and personal identification numbers. These attacks are currently being blocked by web browsers that have a list of such phishing links. It takes several days and intense computing resources to prepare the list. Having a time lag in this process means that many victims are vulnerable at that point in time to such an attack. Inorder to make the process more efficient, the functionality required will be ported to the client side of the web browser. This makes sure that the time delay is averted and the phishing attack can be thwarted with fewer computational resources.

To solve the above mentioned problem, we will be implementing a web browser add-on that works as a background script on the client side. All the background scripts required will be made cross platform compatible to make the development easier and more efficient.

2 INTRODUCTION

2.1 THE INTERNET

The Internet had its humble beginnings as a research project to share and access resources in other computers. And in fact these “other” computers were the main frames of the time that were located in the research facilities and college laboratoires of the United States of America. This Internet is not the traditional Internet that we know as the World Wide Web. It was just a network of computers that could be operated from other nodes in the network. One main roadblock that was encountered was that the main frames were expensive and far exceeded the computational requirements of the then public. As a result of this there were only a handful of expensive mainframes with a few government and private institutions in a single country. Never would they have imagined the used for the internet unless the next era of personal computers boomed. As this phase had computers be considered almost as national assets they were sometimes compromised by enemy nations performing acts of sabotage. And

this is where the roots of performing malicious activities for the gain of nations or individuals can be traced to.

2.2 THE WORLD WIDE WEB

In the stage of the personal computers, people bought computers to perform basic operations like spreadsheets and graphics related tasks. Once the internet had been introduced to those machines, they were still restricted by the slow DSL connections that the telephone companies provided. But once this threshold was broken, we would think that the Internet as we know today would have started to thrive. But it did not. The way in which the information is accessed was not intuitive till the World Wide Web made its appearance in 1989 developed by Tim Berners-Lee at CERN. This era with the internet has its own set of unethical activities that were performed using computers or against the computers. Some elaborate sabotage attempts even involved using the internet connections of those computers.

2.3 THE INTERNET OF CONTENT

The World Wide Web started the phase known as the “Internet of Content” where the websites of different organisations could be accessed by any person on the internet to get to know the organisation and other details like the availability time, recruitment offers and so on. The impact of the World Wide Web was accelerated once the search engines were developed to index and display the billions of web pages that were loaded into servers connected to the Internet. This made possible for people to access almost anything hosted on the Internet without knowing upfront the URL related to the resource. This opened a plethora of problems related to the act of performing unethical activities using or against the World Wide Web. The most common internet based crimes were phishing and site defamation. We will look into them a bit later, once after we have explained the next era of the internet.

2.4 THE INTERNET OF E-COMMERCE

For the next era to come upon, there were a few changes that were required in the internet. They were the ability to host dynamic content based on the users and the ability for the user to interact with such content. These abilities were provided by the advent of web based programming languages with Javascript leading the way. This made possible that people could perform online cash

based transactions for the services that the internet made possible in the first place. This kind of opened the Pandora's box for the cyber related crimes of this day. The notoriety of phishing greatly increased because of the scope that you could have the banking credentials of several thousands of users.

2.5 THE INTERNET OF PEOPLE

And before deep diving into the domain of phishing let us have a short look into the other eras of the internet that have come along. We've had the "Internet of People" which was brought by the advent of social media platforms like FaceBook, Twitter and LinkedIn. People now share much more data online about themselves over the internet to the public. This has lead to even more problems like social media addiction, anxiety and attention deficit among the users. But let us not just paint a dark picture of this era and move on to what the future has in store.

2.6 THE INTERNET OF MACHINES

We are currently in this era of the "Internet of Machines" were more and more IoT devices with the capability to connect to the Internet and use it to communicate with other IoT devices and some centralised computers. This is probably exciting times as even the standards of the Internet of Machines has yet to be decided and wonder if we would be having another World Wide Web like platform available with the machines in mind. Even in this age, the unethical activities can be performed as was done in the previous eras of the Internet.

2.7 PHISHING

Now that we have an understanding of what the Internet is and why it is so important and how it is being used, let us dive into the topic of phishing. If a definition were needed, phishing is any social engineering made to trick the people to access the malicious resources that may get critical information such as passwords, bank account numbers etc. from the victims. Phishing is done not only for the monetary incentives it provides but also for the impersonation of people in social media or to compromise the networks of organisations or countries. They are targeted upon the users who have access to such details like an e-commerce customer or a company manager.

2.8 SOCIAL ENGINEERING

The most common ways to phish are to send emails to those who are related. What such emails contain are the malicious links and other content to convince the users that it is indeed legitimate. The same strategy is applied to the hosted pages that are pointed to by those links. As a result of this many unsuspecting people are tricked. And in order to prevent these instances, many methods have been devised. But the most important thing is to be constantly vigilant that phishing is possible and that you might be a target and never clicking on such links. Some such phishing sites have also been indexed and are even placed above the real site in search engines. So, care must be taken even when the links are provided by the search engine.

2.9 PHISHING PREVENTION

Let us look into the other methods to prevent phishing. Since search engines have to index all pages to be displayed for the user's query, it seems logical to use some mechanism to find such links while indexing and use the same in browsers to notify users that they are accessing a page which is probably used for phishing. This works for most cases, but fails for those dynamically created pages which are not indexed by the search engines and newly created sites which have yet to be crawled because it takes a few hours for search engines to index new pages.

Thus to prevent the problems caused by the above method, new strategies based on the content in the page and also the link that the user accesses, is taken into consideration. This brings with it the advantages its own set of disadvantages like, it having to work on the client side machines and also take the time to render the webpages. Many optimisation techniques and strategies have been developed to overcome this, which will be discussed in the next chapter.

3 RELATED WORK

3.1 AUTOMATIC PHISHING CLASSIFICATION

The work by Colin Whittaker, Brian Ryner and Marria Nazif for Google provided the base benchmark for most of the future works and so would be better if we have a better understanding of what they did. For creating the base dataset required to train a machine learning model, they used the links from the

Gmail spam filters and also those that were submitted by other users. The features used are the URL, the contents of the HTML page and also where the page is hosted. These features are then used by the model which is a logistic regression classifier to find if the site is used for phishing or not.

The training for the model is done offline using the blacklist for the last three months. This has to be done to account for the temporal resilience required from such models. This method of using a published blacklist introduces another risk, which is the risk of feedback loops, that pass down the same error to the classifier. This is because the list might have some false recognitions for the web pages that are submitted manually by other users. This means that the whole dataset has to be manually checked and such wrong classifications be removed. Though the percentage of such instances are very low the fact that this list has to be verified manually really is a black mark for this method. Though they achieve an impressive false positive rate of under 0.1%. Even though we consider that the black list is perfect, the fact that the system uses a black list to work means that there will be an inevitable time lag between the time the phishing site is up and that the page is detected to be used for phishing.

This time lag has to be reduced by decreasing the time taken to develop the model and thereby help the users to find even the most recent phishing pages. In spite of the obvious shortcomings that this work had, the main takeaway would be that the features required for the model are

1. The URL of the page
2. The HTML page contents
3. The host server details

3.2 CANTINA

This work can further taken up by Guang Xiang, Jason Hong, Carolyn P. Rose and Lorrie Cranor in their CANTINA+ which provides a feature-rich machine learning framework for detecting phishing web sites. It is split into two phases. In the first phase, the task is to find the feature values for all the records in the database. Once this is done, the second phase is to find if the site is used for phishing or not. The above feature is limited to 15 features which are used in both the phases.

This method for performance optimization and to reduce false positives uses a hash based page removal model in which the similar looking components of the website are removed, making it easier to find the differences. Once, the similar components are removed, the presence of a login form in the HTML content is searched for. If the content matched that of the login page, then the pretrained model is used.

Since the creation of the model requires an updated list of phishing sites, the list provided by PhishTank's verified blacklist is being used. And as far as the time required to find the similar looking sites is concerned, it is greatly reduced by using the SHA 1 hashing algorithm. It is noted that though this hashing algorithm can be easily broken, it is being used for the efficiency and the high accuracy with which it finds out the phishing sites.

Though this model provides a faster way of finding those sites, it definitely comes with its own set of drawbacks. The first one being that SHA 1 algorithm that defeats the whole purpose of the malicious actor never being able to circumvent the system.

Though this system has its disadvantages, the features are a few things that can be taken from them. They include the embedded domain, IP address, number of dots in URL, suspicious URL, number of sensitive words in URL, out of position top level domains (TLD), bad forms, bad action fields, non-matching URLs, out of position brand names, the age of domain, page in top search results, page rank and page results while searching for copyright company name and host name.

Many models were used and it was found out that the Bayesian models and the Random Forests performed remarkably well.

3.3 AUTO UPDATED WHITELIST

Ankit Kumar Jain and B. B. Gupta provided another approach to protect against phishing attacks at client side using auto-updated white-list. The accurate and fast detection of phishing sites in real time environment is paramount. The time constraints of the above methods are because they use a visual similarity based approach. This can be reduced by using a heuristic based approach which depends mainly on the feature set, the classifier and the training data.

The hypothesis is based on the fact that though the phishing pages look similar to the corresponding real website, they do differ steeply in the functionality they offer. But almost all but the critical phishing functionality redirect to the corresponding real website.

To provide such functionality, a whitelist is used in this method. The whitelist contains the domain name and IP address as the parameters. The whitelist provides for the faster running time and to reduce the false positive rate. The working of the whitelist is as follows.

First when the user has never visited any site, the whitelist has no records. But when the user does so and visits a site, the whitelist is checked if it has the domain along with the IP address. If the record is present, the page is said to be a safe site. Else, the second component which is almost the same as that of the previous models kicks in to find if the requested site is a phishing site or not.

Thus, because of the whitelist the model has the advantage of being language independent and is capable of finding the embedded components in the phishing website that can cause a DDOS attack.

The model was developed further into a usable application down the lane. And this paper provided the following findings. The model provided the base for using auto-updated whitelists to speed up the process of finding out if the page is used for phishing or not. This is highly advantageous because most of the sites that a person accesses will not be a phishing page, significantly reducing the average time taken to process the site.

3.4 OFF-THE-HOOK

The work by Samuel Marchal, Giovanni Armano, Tommi Grondahl, Kalle Saari, Nidhi Singh, and N. Asokan implemented the client-side phishing prevention application named Off-the-hook.

The main improvements that they provided were better privacy, realtime protection, resilience to dynamic phishing and effective warnings. It also has an emphasis on making a software application that can be easily used by the general public to protect themselves against phishing.

A brand independent, evolution resistant model to detect the phishing pages has been developed. This means that sites of all domains will be identified correctly and the temporal resistance maintained. This is that even when the malicious actor knows how the system works, and tries to figure out a work around, the system learns and adapts itself to the changing conditions.

The main upside for this project is the cross platform capabilities that it intends to provide and the user base that it can have based on the ease of use that the application provides.

The downside is as follows. All the functionality required cannot be implemented inside the browser using Javascript alone. As a result of this, machine dependence creeps into the equation. This in the long run will be a major block to the ease of use that the application says will provide.

The other main disadvantage is that the model does not take into account the static IP addresses on which the page might be hosted. This model relies on the assumption that such IP addresses that host the phishing pages will be blacklisted and removed by the host provider.

Thus based on all the above mentioned related works, we will be redeveloping the application Off-the-hook to follow the basic networking and socket connections so that the applications both running inside the browser and within the operating system of the user will remain compatible.

4 HIGH LEVEL BLOCK DIAGRAM

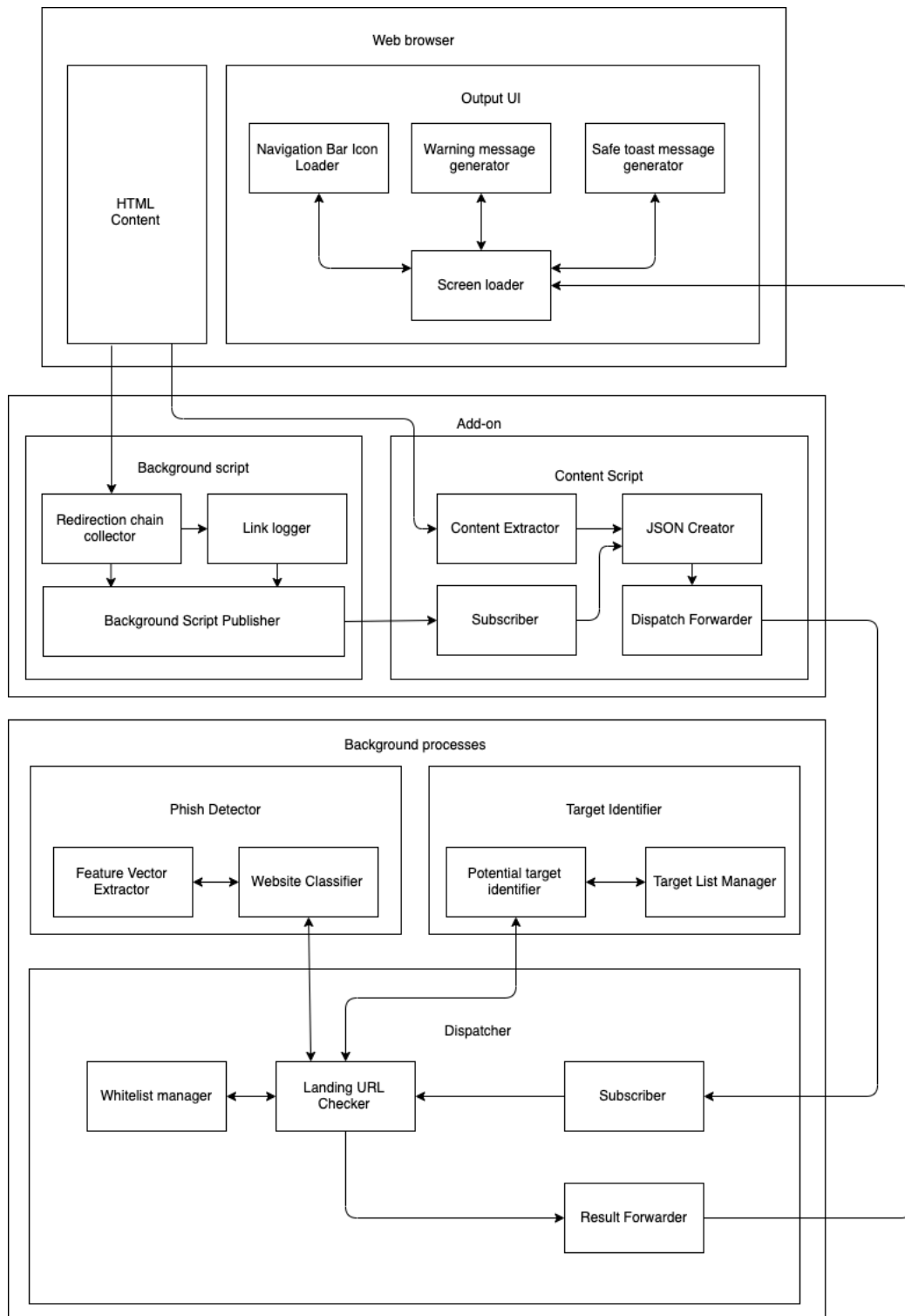


Figure 4.1 High Level Block Diagram

The above Figure 4.1 explains the detailed working of the project with the pipelines between the multiple components. To get an easier understanding, a simpler diagram has also been drawn as Figure 4.2

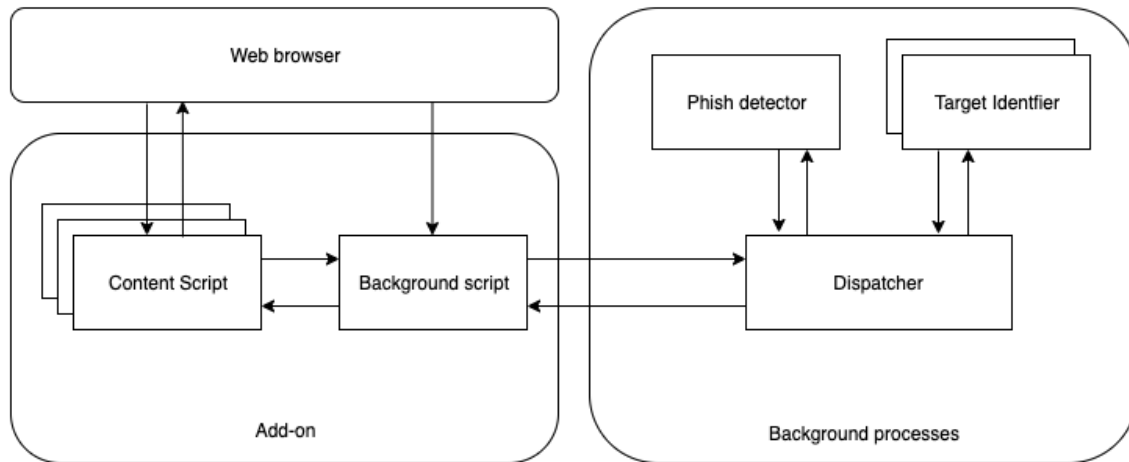


Figure 4.2 System Architecture

5 DETAILED MODULE DESIGN

5.1 MODULES

Our system has three modules that complete its functionality. They are as follows.

1. Add on
 - a. Background script
 - b. Content script
2. Background process
 - a. Dispatcher
 - b. Phish Detector
 - c. Target Identifier
3. Web Browser
 - a. HTML content
 - b. Output UI

5.2 ADD ON

An add on is required in the form of an extension because of the reasons that the frameworks provided by Javascript for tasks related to machine learning are still in its infant stages. And so the addon performs the two important tasks of getting the page contents from the browsers and also making changes to the elements in the web browser to convey the threat level for phishing.

5.2.1 BACKGROUND SCRIPT

The task to be done by this component is to get the list of redirection URL chains from the user's web browser, and publish it to the content script while also logging the links for reference purposes.

Begin

For each page load redirect

Add listener to that event

Get the list of redirects from listener

If page is fully loaded

Send the list of redirects to content script

Done

End

5.2.2 CONTENT SCRIPT

The task to be done by this component is to get the landing URL of the page and the contents of the page from the browser and also the URL redirection list from the background script and then combine those and send it to the background process which will have to do further computation.

Begin

For each page load redirect

If page is fully loaded

Get the URL from the tab

Get the HTML content from innerHTML tag

Get redirection list from background script

Send them to the background process

Done

End

5.3 BACKGROUND PROCESS

The background process gets the contents from the content script and identifies if the site is used for phishing or not. It has the following subcomponents which are discussed in detail.

1. Dispatcher
2. Phish Detector
3. Target Identifier

5.3.1 DISPATCHER

The dispatcher is used for the performance enhancements it provides by using the whitelisted addresses that can be used without even having to run the model. It has direct control over the phish detector and target identifier.

Begin

If page address is in whitelist

Send the GREEN signal

Else

Send content to phish detector

Get results from phish detector

If phish is FALSE

Send the GREEN signal

Else

Send the RED signal

Send content to target identifier

If target is found

Publish target

Else

No target matched

End

5.3.2 PHISH DETECTOR

The phish detector gets the content from the dispatcher and gives the result which is either the site is a phish or not. It is done by using a machine learning model.

Begin

For each page URL

Get the feature values for the URL

Load the saved model

Publish the result

Done

End

The features used are as follows,

1. Have IP address
2. URL length
3. Shortening service
4. Having @ symbol
5. Double slash redirecting

6. Prefix suffix
7. Having sub domain
8. Domain registration length
9. Favicon
10. HTTPS token
11. Request URL
12. URL of anchor
13. Links in tags
14. Server form handler
15. Submitting to email
16. Abnormal URL
17. IFrame redirection
18. Age of domain
19. Web traffic
20. Google index
21. Statistical Reports

The model is a Random Forest Classifier which has the pseudocode as follows.

Begin

For each record in dataset

Get the feature values

Create an arff file to save results

Done

Train the dataset with at least 7 splits

Save the model as pkl file

End

5.3.3 TARGET IDENTIFIER

Once the dispatcher gets the signal that a site is phishing, it can be useful to find which site is being used as a template so that the unsuspecting user is fooled. This is done by using the similarity of hashes between the phishing and target website. The SHA algorithm is as follows.

Begin

Get the hash value for page content

Compare with values in hash list

If match

Display target

Else

No target found

End

5.4 WEB BROWSER

Though the above described components play major roles in this project, the one that the user will be able to view is this component. And so care has to be taken to make it look as professional as possible.

5.4.1 HTML CONTENT

The add on must be published as extensions for the browsers and so the UI of the components must be taken care of. And the tasks of the background scripts must run as helper tasks and not interfere with the main script, otherwise the UI of the extension will appear to be jittery.

5.4.2 OUTPUT UI

The two possible results for this project are that either the site is a phish or not. And if a site is not a phish no changes have to be made to notify the user. But if the site is found out to be a phish the changes made to the UI must meaningfully convey to the user that the site is a phish.

Begin

If site is phish

Change icon to red

Display warning message
If site has target
 Display target link
 Else
 Display no target
 Else
 Change icon to green
 Display safe to proceed message
 End

6 IMPLEMENTATION

6.1 CONTENT SCRIPT

This component takes the input which is the URL of the page and its contents and sends it to the background process along with the redirection URL from the background script. A design of this component is required as the Chrome Extension cannot get the tab HTML content in the background and so a script in php is used to get the content of the URL provided by the Javascript code segment.

tablink = tab.url;

The above javascript gives the URL of the tab which is used to get the content using the following php code

\$site=\$_POST['url'];
\$html = file_get_contents(\$site);

This component could not be activated once the tab is fully loaded and as a result a button to click has been provided.

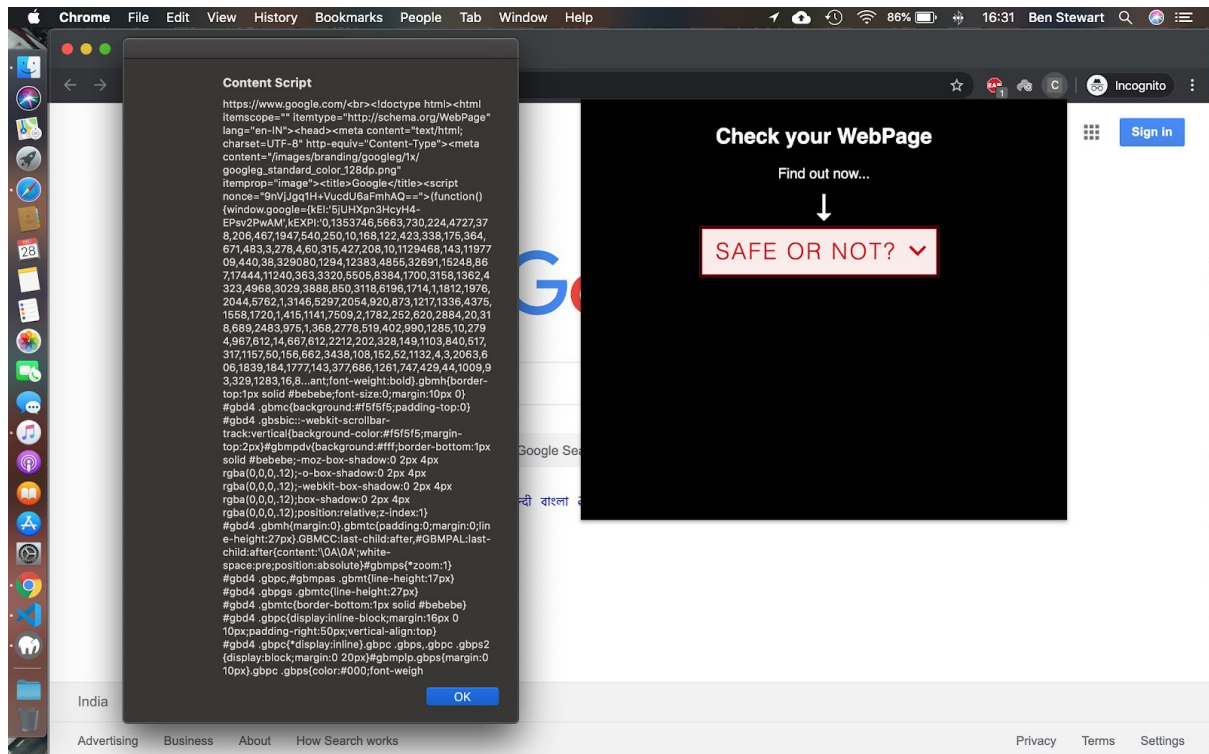


Figure 6.1 Content Script

6.2 BACKGROUND SCRIPT

This script is based on the open sourced code on GitHub that demos how to get the background URL redirects of the current tab. It handles multiple types of redirects and also their security levels based on the URL redirects. This component finally returns the list of path components that the page had traversed through. The path item is as follows,

url: pathItem.url,

status: pathItem.status_line,

redirect_type: pathItem.redirect_type,

redirect_url: pathItem.redirect_url,

meta_timer: pathItem.meta_timer

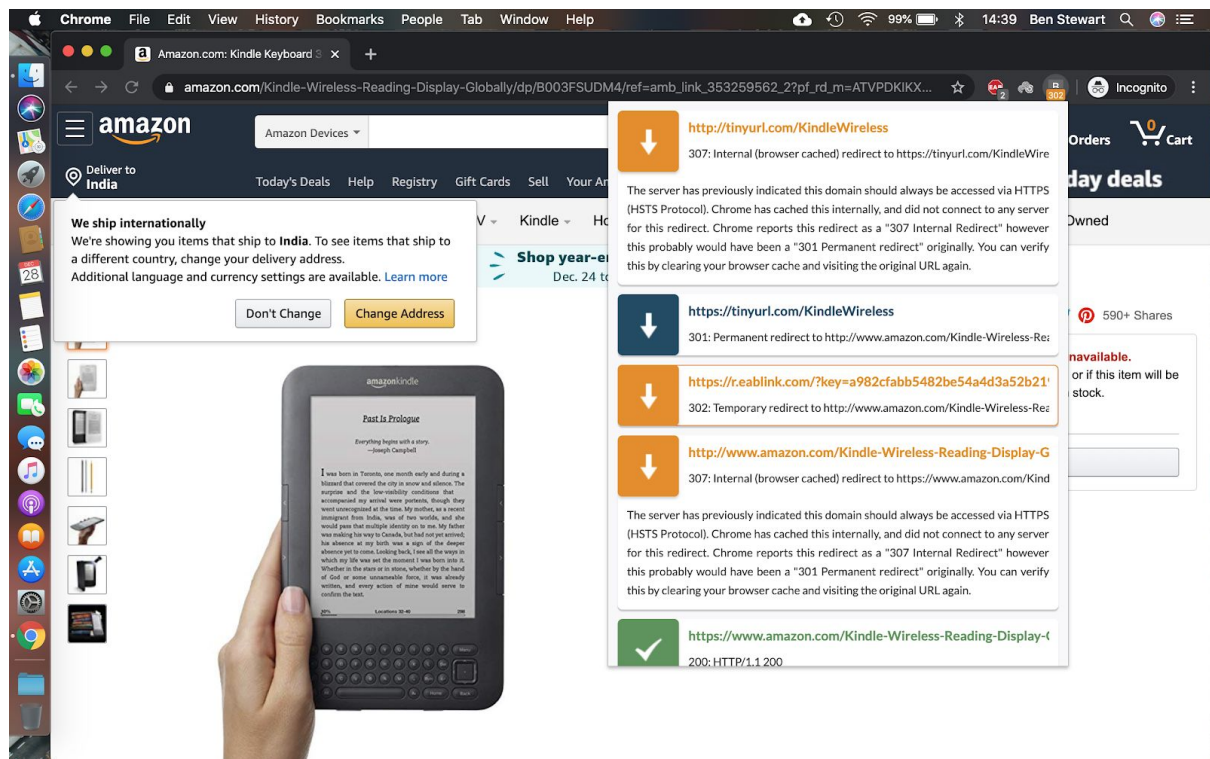


Figure 6.2 Background Script

7 EVALUATION METRICS

This application must be evaluated for the following performance metrics, which become important as this is meant for users who will never need to know about computer programming.

Phish detection accuracy: This metric is really important because the main task for this application is to notify the people if the site is a phish or not. It is defined as the ratio of the total number of correct classifications to the total number of classifications.

Target detection ratio: This metric is to measure the ease of use for the user by providing them with the original site which is being mimicked by the phish. This will be the ratio of phish sites whose target has been found to that of the total number of phish sites.

Memory usage profiling: Since this is an application to be used by many people who might have different configurations of machines, we must take into consideration that the application must use as little memory as possible.

Addon rendering time: Since the addon has a background component which has to collect the data from multiple tabs that could be running simultaneously, the addon must be tested to check if it is stable and does not take time to render and thereby blink.

Temporal resilience accuracy: This is a metric which says that the application must be resilient to adaptations by the malicious actor, over time. Thus this can be measured by checking if the accuracy does not reduce with the passage of time.

REFERENCES

- [1] S. Marchal, G. Armano, T. Gröndahl, K. Saari, N. Singh, N. Asokan, "Off-the-hook: An efficient and usable client-side phishing prevention application", *IEEE Trans. Comput.*, vol. 66, no. 10, pp. 1717-1733, Oct. 2017.
- [2] A. K. Jain, B. B. Gupta, "A novel approach to protect against phishing attacks at client side using auto-updated white-list", *EURASIP J. Inf. Secur.*, vol. 2016, no. 1, Dec. 2016.
- [3] G. Xiang, J. Hong, C. P. Rosé, L. Cranor, "CANTINA: A feature-rich machine learning framework for detecting phishing Web sites", *ACM Trans. Inf. Syst. Secur.*, vol. 14, no. 2, 2011.
- [4] *Implementation for the Usage of Google Safe Browsing APIs (v4)*, 2019, [online] Available: <https://github.com/google/safebrowsing>.
- [5] C. Whittaker, B. Ryner, and M. Nazif, "Large-scale automatic classification of phishing pages," in *Proc. Netw. Distrib. Syst. Security Symp.*, 2010, pp. 1–14.