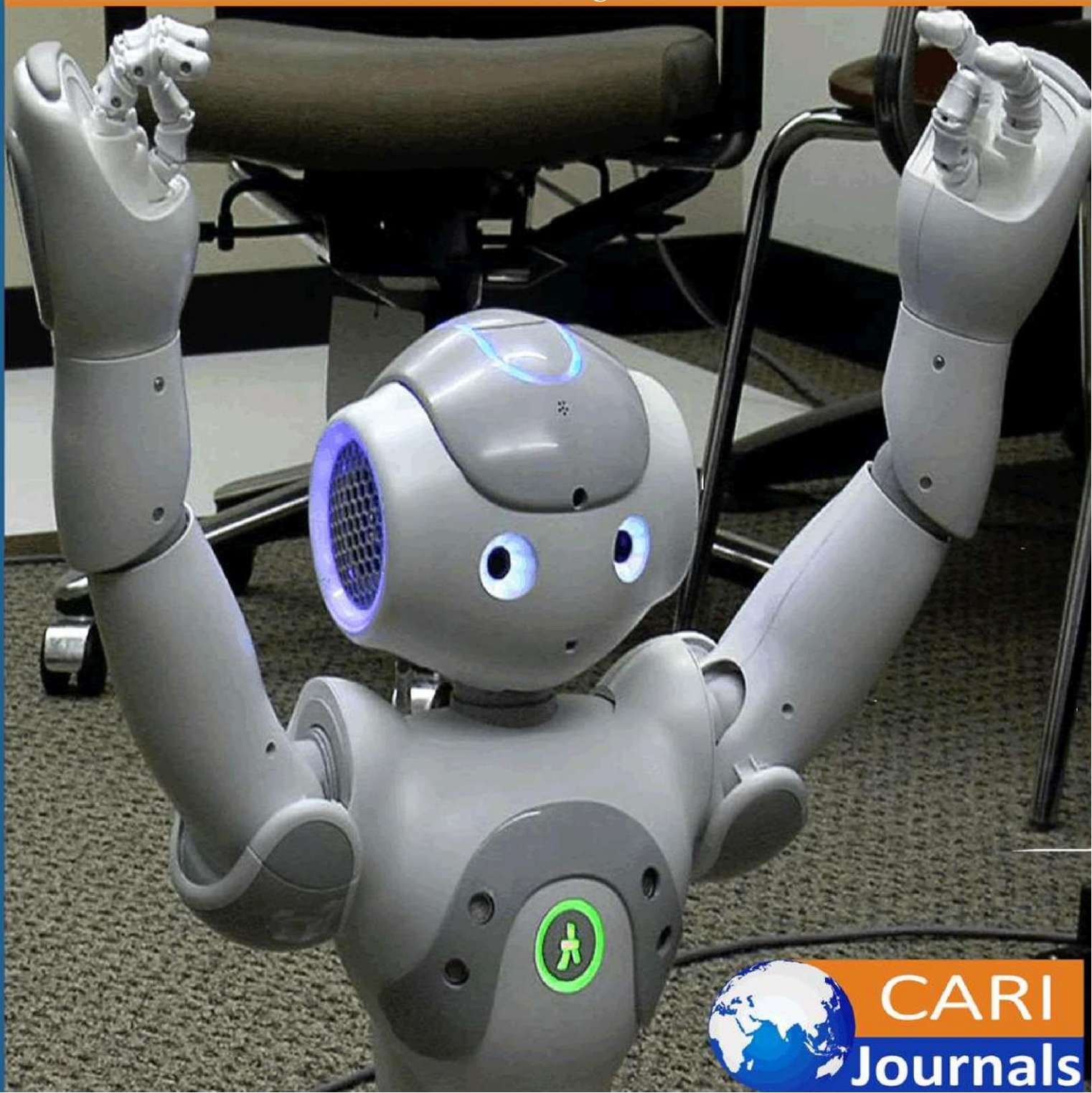


# International Journal of Computing and Engineering (IJCE)

Secure Browse: AI-Powered Phishing Defense for Browsers



CARI  
Journals

## Secure Browse: AI-Powered Phishing Defense for Browsers

 Santosh Kumar Kande

*Accepted: 25<sup>th</sup> Mar, 2024 Received in Revised Form: 25<sup>th</sup> Apr, 2024 Published: 25<sup>th</sup> May, 2024*

### Abstract

**Purpose:** With the rising threat of phishing attacks exploiting user naivety, this report introduces a novel approach to bolster web security. Traditional rule-based systems and existing solutions fall short in addressing sophisticated phishing attempts. The proposed solution entails a Chromium-based browser extension that leverages machine learning classification techniques.

**Methodology:** A Python web server, utilizing decision trees, k- nearest neighbors, and random forests, assesses the legitimacy of a given URL. The extension communicates with the server, providing real-time notifications to users when visiting potential phishing sites.

**Findings:** Experimental results demonstrate the effectiveness of the ensemble model with an accuracy of 90.68%, marking a significant improvement over rule-based alternatives.

**Unique contribution to theory, policy and practice:** Future work includes refining models, incorporating user feedback, and expanding the application to diverse platforms and contexts.

**Keywords:** *Phishing Defense, Machine Learning Classification, Web Server Architecture, Ensemble Model*

## 1 Introduction

With the recent advancement in various cybersecurity technologies, the weakest link in cybersecurity happens to be the end users. Attackers utilize phishing, which exploits the naivety of users to trick them into handing out sensitive information. This poses a great risk not only to the users themselves but the organizations and institutions of which they are a part of. According to recent research from Proofpoint, 75% of organizations around the world experienced a phishing attack in 2020, and 74% of attacks targeting US businesses were *successful*.<sup>8</sup>

Most of the software interactions between users and organizations happen through websites on computers/laptops. Forcing users to install a custom application to be run is a rarity due to the convenience a browser provides. Web-based security for browsers is crucial today, and Phishing is one of the crimes that still exists.

Apart from increasing security awareness among users, the tools that complement that awareness to help users make safe decisions must be developed. This report proposes and demonstrates a Chromium-based browser extension to help mitigate the risk of phishing while browsing the web. The central idea of the browser extension is to notify the user whenever they open any potential phishing website.

The solution also includes a Python web server, which utilizes various Machine Learning classification techniques to determine the legitimacy of the web-page in question. The web server takes in a URL and returns a boolean value indicating if the given URL is part of a potential phishing attempt.

The browser extension monitors each URL that the user visits and tries to determine if the URL is malicious with the help of the web server. The web server exposes a REST API consumed by the extension for communication. The same API can also be reused to implement a similar phishing detection in a different context like a network-level application or a mobile application like Android.

The server implements both a mixture of rule-based approach and Machine Learning classification techniques. The rule-based approach is useful for weeding out obvious URLs and is inexpensive. The Machine Learning classification techniques is more expensive to do but help predict whether a URL can be a phishing site.

## 2 Related Work

*PhishDetector*<sup>7</sup> This Chromium extension has set out to solve the same problem as this project. Based on its description on the website and the analysis of its behavior, it can be concluded that this extension uses a rule-based system to determine if a webpage is a phishing attempt. It also seems to be particularly accurate when it comes to identifying illegitimate banking pages. Even though rule-based systems are great for detecting simple phishing attempts, they are not ideal for more sophisticated ones. Rule-based systems are also inherently complex to maintain - adding and modifying rules over time makes the system more complex and unmanageable with time. They

also demand more human intervention to define the rules and maintain them. Using Machine Learning instead of a rule-based system eliminates many of these limitations. As Machine Learning is data-centric, it doesn't require managing complex rules and makes it straightforward to tweak the algorithms.

*Cloudphish*<sup>3</sup> Cloudphish is a phishing detector for web-based email software. It monitors all emails users receive and checks each email for a phishing attack. Having a paid subscription model, it offers a decent service. However, its major limitation is that it works only with the email inbox. Even though many phishing attacks are carried over through email, phishing is as prominent on social media and messaging apps in this age. And that calls for a solution that monitors all web activity to identify phishing attacks regardless of their delivery method.

Various other browser extensions virtually have the same limitations as the aforementioned solutions<sup>2911</sup> As their limitations are encompassed in discussing other solutions above, their detailed discussion has been omitted for brevity.

To summarize, there are various browser plugins consisting of rule-based systems, simple whitelists-blacklists, and some even using Machine Learning and Artificial Intelligence. However there needs to be a solution that utilizes all available phishing detection methods to protect the average internet user from criminals.

### 3 Approach

The architecture of this project consists of two primary components: The browser extension and the web server.

*Browser Extension* The extension is developed for Chromium-based browsers using JavaScript with HTML and CSS. Therefore, it is compatible with any Chromium implementation, including Google Chrome, Microsoft Edge, Brave, etc. The extension monitors each web page that the user visits and fetches the URL of that web page. It then communicates with the web server, which tells if the URL is part of a phishing attempt. The user is notified of the analysis results based on the server's response. The extension will stay silent in the background while the user is visiting websites deemed safe. It only bothers them when there is a potential of phishing on the site, they are currently visiting to help them make a safer decision.

*Python Web Server* The web server has a REST endpoint that takes in a URL and uses Machine Learning techniques to classify it. It extracts relevant "features" from the URL and feeds them to *Classification Models* to determine the legitimacy of the URL. This process is expanded upon further down in this section.

Classification is a process of categorizing a given set of data into classes. There would be two classification labels in this case: "spam" and "not spam". The input data would be values of various URL features that are deemed effective for the high-quality classification of any URL into one of the classes.

The process of creating a classification model consists of two primary stages. The first is the training stage, where a classifier is fed a large amount of input along with their respective class labels. That creates a classification model with a set of inputs without their respective labels to let it classify each input into a class. That constitutes the second stage, where the correctness of the newly created classification model is compared against the actual labels. The input set given to the model for the second stage is often referred to as *testing dataset*, and the data used for the first stage is called *training dataset*. As a common practice, the dataset on hand is split into training and testing datasets for creating and testing the classification models, respectively.

The classifiers used in this project are described below:

1. **Decision Trees:** Decision Trees belong to the family of supervised machine-learning algorithms. Decision trees classify the input by running them down the tree from the root node to some leaf node, whereas the leaf node provides the classification of the input.
2. **K-Nearest Neighbors (KNN):** The KNN algorithm assumes that similar things are near to each other. Based on this assumption, it classifies all nearby data points into one. Then, it classifies the given input by locating N-nearest neighbors and finding a mode of their labels, which is predicted to be the label of the input set.
3. **Random Forests:** Random Forests consist of many Decision Trees that operate as an ensemble. Each tree in a forest classifies the given input set to a label, and the label with the highest number of votes is considered the Random Forest prediction.

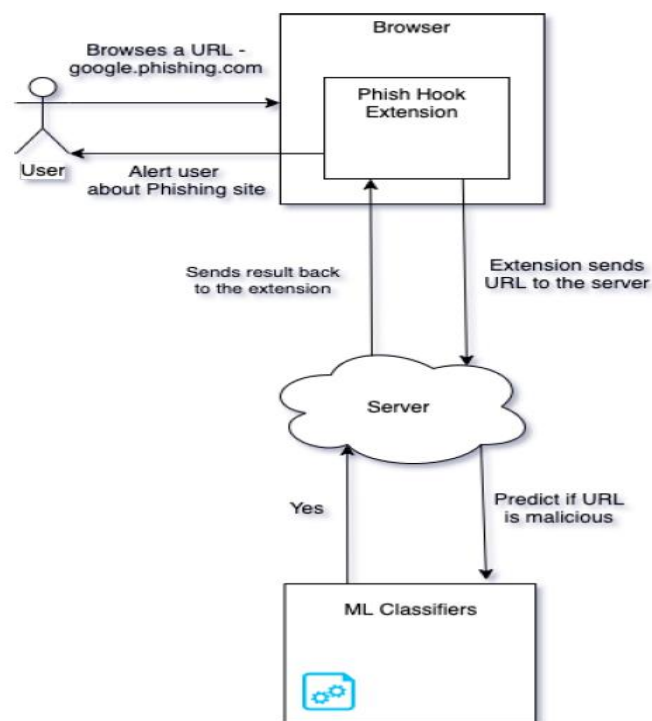


Figure 1. Architecture overview of the application

The individual trees in a random forest are relatively uncorrelated, and they perform better as an ensemble than they would on their own. To put it in layman's terms, the trees safeguard each other from their errors.

The following section describes the datasets used for training and testing the models with the aforementioned classifiers.

1. Phishing Websites Dataset:<sup>1</sup> This dataset has 11,055 data points with 6,157 legitimate URLs and 4,898 phishing URLs. And it contains 30 features sub- divided into three categories:

- (a) Features based on the domain and subdomains
- (b) Features derived from the other parts of the URL
- (c) Features derived from the webpage HTML and JavaScript

2. Datasets for phishing websites detection:<sup>10</sup> This dataset consists of 111 features, of which 97 are based on the URL. For the phishing websites, the list was extracted from the PhishTank registry which are verified by multiple users. And for sets of legitimate websites, the publicly available and community labeled lists are utilized.<sup>5</sup>

When the server receives the URL before any of the Machine Learning models can test it, it is checked against a whitelist. The whitelist is extracted from The Majestic Million dataset,<sup>6</sup> which maintains a list of top 1 million domains on the internet. The whitelist helps reduce the processing and network overhead of checking the most popular sites that an average user visits daily. The experiments section further elaborates on the need for doing this.

The machine learning models are already trained with both datasets, are precomputed, and stored on the server. If the given domain is not part of the whitelist, the URL is passed on to these models. Whenever the server receives a URL to be tested, it extracts all the features from the given URL. Many of the features are extracted through string parsing, and the rest of them require the use of external APIs and libraries. For example, PageRank of a given domain is fetched using Open PageRank API.<sup>4</sup>

After extracting the relevant features for each dataset, they are passed on to the classification models and their ensembles. The result is sent back based on the outputs of the various models. One advantage of using a server-based approach is that essentially, results for different URLs from different users can be cached. This can save recurring computation for the same URL across users and decrease response times on repeated requests.

#### **4 Experiments and Results**

The experiments in this report are mainly concerned with the effectiveness and efficiency of phishing detection. Due to the nature of the problem, the focus was more on effectiveness since there were a lot of parameters to tune. That being said, efficiency concerns are legitimate and were explored as well.

One of the important characteristics of a machine learning model is the data set used to train it. The nature of the data set greatly impacts the effectiveness of the result. Initially, Experiments were done on various data sets available across the web using simple classification methods, and the accuracy results were compared to filter out the ineffective collections.

The decision to use two data sets as a part of the final application was based on the observation that "Phishing Websites Dataset"<sup>1</sup> had a lot of non-url based features that worked very well based on the content of the website. In contrast, "Datasets for Phishing websites detection"<sup>10</sup> data set contains 97 features that depended only on the URL structure, which effectively detect phishing websites. Including both data sets, with the result from their combined models, effectively acts as a way for the result to have legitimate checks and balances as well. If both models predict that a website might be a phishing website, there is a high probability that it is true. Whenever there is a partial agreement, it is prudent to let the user know and be the final arbiter of this conflict. This exploration allowed the feature of "Caution" vs "Alert" notifications to come up organically.

The models were built by splitting a portion of the input as training and test data. Using a Decision Tree classifier and cross-validated into 10 splits. It was verified manually by using it as an extension in regular browsing to flag obvious errors. For the model from the dataset "*Datasets for phishing websites detection*"<sup>10</sup> it was noticed that the model's accuracy was high with the test data but performed rather poorly with random samples from the real world. This led to the hypothesis that the model was overfitting the data set. This led to the introduction of the ensembles for this model in particular. The combination of additional classifiers was tried out (Decision Tree, k-Nearest Neighbors and Random Forest) and verified with the same methods, until a minimum accuracy of 80% was achieved.

Tuning individual classifiers is also detrimental to the overall accuracy measurement of a model. The classification ensemble contained a K-Nearest Neighbor classifier for the "Datasets for phishing websites detection"<sup>10</sup> model. The number of neighbors for the classification was increased to the effect, that it over-fit the data. This was fruitful since the ensemble of different classifiers helped reduce the overall variance of the model.

The above experiments were predominantly effectiveness-based approaches. Although using a combination of two models helped improve the overall goal of detecting phishing websites, it came at the cost of the delay in prediction. Steps were taken to parallelize the flow of aggregating the predictions from the two models. That being said, there is no denying that using a classifier to predict outcomes is inherently expensive. In a resource-constrained environment, it must be called upon only when truly necessary. Once it's called upon, it is prudent to reuse its results to save the cost of classifying again. To this effect, a "Whitelist" of frequently used websites was added to skip prediction entirely for overtly obvious websites. Since the intended application of this classification is in browser extensions, the fact that users will browse popular websites frequently is not an entirely unfair assumption. Caching the results on the server also goes well with this approach. This rule-based counterpart for the Machine Learning classifier helps improve the overall efficiency of the server.

**Table 1.** Performance Results of the Models

| Model                               | Accuracy | Precision | Recall | F1 Score |
|-------------------------------------|----------|-----------|--------|----------|
| <i>Ensemble model</i> <sup>10</sup> | 90.68    | 92.69     | 89.20  | 90.92    |
| Decision tree Model <sup>1</sup>    | 80.37    | 84.03     | 80.82  | 82.40    |

## 5 Conclusion and Future Work

Phishing is one of the most widespread security attacks on the internet of this age. As the attacks become increasingly sophisticated, the solutions to prevent them need to keep with them. Traditional methods of detecting attacks must be combined with the growing fields of machine learning and artificial intelligence. It is apparent from various experiments and observations that the ideal solution to such issues is to combine the best approaches and make them work with each other.

For the enhancement of this implementation, there seems to be a potential to improve the machine learning models and the features that are being used for classification. The extension and the server can be updated to also take in user feedback regarding a misclassification of any site. And the feedback would be considered by the classification models to improve their performance. Although, it might introduce a fresh set of challenges such as abuse of the feedback system to pass off an illegitimate site as safe (adversarial attacks). Potential solution to this problem would be to have a threshold on number of feedbacks until which the system does not consider them for the particular website.

The solution proposed and implemented in this project can be a genesis to a series of tools and products that strive to solve the issue of phishing attacks on the internet. The existing rule-based solutions that use HTML properties of the webpage can be incorporated into the machine learning models for potential improvement in performance. The existing REST APIs and the Python interfaces can also be reused to monitor phishing in other contexts. There can be an Android application that monitors the URLs visited by the user and notifies them if any of them have a potential to be a threat. And the same use case can be applied to other platforms, operating systems and at different abstraction levels like network or system.

## References

- [1] et al, M.: Phishing websites data set by uci,  
<https://archive.ics.uci.edu/ml/datasets/Phishing+Websites>
- [2] Acra, B.: Blue arca anti-phishing extension,  
<https://chrome.google.com/webstore/detail/blue-arca-anti-phishing-e/>



- [3] Cloudphish: Cloudphish anti-phishing extension, <https://chrome.google.com/webstore/detail/cloudphish-anti-phishing/fcahokjdmffdhglnlhgbceafccfdjkd?hl=en>
- [4] Initiative, O.P.: Open pagerank api, <https://www.domcop.com/openpagerank/>
- [5] Lab, T.C.: Test lists, <https://github.com/citizenlab/test-lists>
- [6] Majestic: Majestic top 1 million websites, <https://majestic.com/reports/majestic-million>
- [7] Moghimi, M.: Phishdetector – True phishing detection, <https://chrome.google.com/webstore/detail/phishdetector-true-phishi/kgecldbalfgmgelepbbloodfoogmjdgmj>
- [8] Proofpoint: Threat report: 2021 state of the phish report, <https://www.proofpoint.com/us/resources/threat-reports/state-of-phish>
- [9] Retruster: Retruster phishing protection, <https://chrome.google.com/webstore/detail/retruster-phishing-protec/akcpbmbdplmbhlpeglpbghnkcbbhiapil?hl=en>
- [10] Vrbančič, G.: Datasets for phishing websites detection, [Phishing Websites Dataset - Mendeley Data](#)
- [11] Zuelsdorf, A.: Phishing boat, <https://chrome.google.com/webstore/detail/phishing-boat/ljaiihgfejaggbjfieldfnjdckomlfop?hl=en>



©2023 by the Authors. This Article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>)