

## Research Proposal

**Team members:** Bar Goldman 208785154 , Kirill Perevalov 336485636

**Paper:** Highly accurate phishing URL detection based on machine learning.

Sajjad Jalil · Muhammad Usman · Alvis Fong © The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2022

### **Abstract:**

Phishing is a type of cyber-attack, when phisher mimics a legitimate website page to harvest victim's sensitive information and misuse it.

the paper presents that Various detection methods exist, including AI-based, third-party, heuristic, and content-based techniques, but they have limitations like

1. features extracted in the past are extensive, with a limitation that it takes a considerable amount of time to extract such features.
2. Different methods, including statistical and custom feature proposals, can yield incorrect results without sufficient domain knowledge.
3. Most research uses small, pre-classified datasets, which do not accurately reflect real-world efficiency and precision.
4. Previous approaches are limited in their evaluation metrics.

The paper tries to solve those problems by introduces an efficient machine learning framework for predicting phishing URLs without the need to visit the webpage or rely on external services. The proposed technique utilizes various components of the URL, including the full URL, protocol scheme, hostname, and path area. It incorporates features such as entropy, suspicious words, and brand name matching using the TF-IDF technique to classify phishing URLs.

In this project we plan to enhance the methodology of detecting phishing URLs using advanced machine learning techniques.

This will involve a deeper analysis of the existing datasets, identifying potential areas for improvement in feature extraction and model training.

We want to try incorporate more sophisticated algorithms.

And we plan to extend the dataset by including more recent and diverse phishing URL examples, ensuring the model remains effective against evolving phishing tactics.

### **Related work,**

Specific studies and surveys related to phishing detection are referenced in the paper, such as the works of Dou et al. (2017), Alsharnouby et al. (2015), Chiew et al. (2018), and Jalil and Usman (2020). These studies focus on different aspects of phishing attacks, including the process of phishing, the difficulty for individuals in identifying phishing URLs, and comprehensive reviews of phishing detection techniques.

The paper also compares its proposed framework against past approaches. For instance, it mentions the work of Shahrivari et al. (2020), Feng et al. (2018), Zhu et al. (2019), and others, highlighting their achieved accuracies in phishing URL detection using various machine learning

algorithms. This comparison allows for a better understanding of how the proposed framework stands in relation to existing methodologies

**Dataset:**

The datasets they used in the research are collected from various sources, including Kaggle and CatchPhish.

we will use [Kaggle.](#)