# PHISHING URL DETECTION

## MINI PROJECT REPORT

*Submitted by*

| | |
|---|---|
| **MOHAMED HUSSAIN H** | **(111521102096)** |
| **NAVEEN BALAJI V** | **(111521102104)** |
| **MYTHREYAN GV** | **(111521102102)** |

## BACHELOR OF ENGINEERING

### *IN*

## COMPUTER SCIENCE ENGINEERING



## R.M.D. ENGINEERING COLLEGE

### (An Autonomous Institution)

## KAVARAIPETTAI - 601206

## NOVEMBER 2023

# BONAFIDE CERTIFICATE

Certified that this project report titled *"PHISHING  URL  DETECTION"*,  is a bonafide

work of who carried out the work under my supervision.

**MOHAMED  HUSSAIN  H**           **(111521102096)**

**NAVEEN  BALAJI  V**           **(111521102104)**

**MYTHREYAN  GV**           **(111521102102)**

**SIGNATURE**                                    **SIGNATURE**

**Dr. P. Ezhumalai, B.E, M.Tech,  Ph.D.,**        **Mrs. K. Padmapriya, B.E, M.E.,**

**HEAD OF THE DEPARTMENT**        **ASSISTANT  PROFESSOR**

Department  of Computer  Science           Department  of Computer

Science  Engineering,                          Engineering,

RMD Engineering College                   RMD Engineering College

R.S.M Nagar                                   R.S.M Nagar

Kavaraipettai – 601 206.                      Kavaraipettai – 601 206.

# VIVA-VOCE EXAMINATION

The Viva - Voce Examination of the following students who have submitted this project work held on………………

**MOHAMED  HUSSAIN  H** **(111521102096)**

**NAVEEN  BALAJI  V** **(111521102104)**

**MYTHREYAN  GV** **(111521102102)**

**INTERNAL  EXAMINER**

# ACKNOWLEDGEMENT

# ABSTRACT

In an era characterized by increasing online activity and growing digital interconnectedness, the threat of phishing URLs has become a pervasive cybersecurity challenge. Phishing attacks, often employing deceptive emails and fake websites, aim to deceive users into revealing sensitive information or downloading malicious content. These attacks continue to evolve, making it crucial to implement advanced solutions for their detection and prevention.

This project represents a comprehensive effort to combat phishing URLs effectively by integrating Machine Learning (ML) techniques with a user-friendly browser extension. The project centers on the development of a browser extension that communicates with an ML model hosted on a remote server. The extension provides real-time URL scanning, offering immediate protection to users as they navigate the web. ML algorithms, trained on a meticulously collected and balanced dataset, enable the extension to differentiate between benign and phishing URLs with remarkable accuracy.

The significance of this project lies in its ability to adapt to the dynamic nature of online threats. Unlike traditional rule-based and signature-based methods, the ML model can continuously learn and evolve, ensuring it remains resilient against evolving phishing tactics. Users benefit from a user-friendly interface, clear warnings, and the flexibility to customize their security settings. By employing encryption and anonymization techniques, the project prioritizes user data privacy.

The project's impact on online security is substantial, offering a potent defense against phishing threats. It aims to protect users, individuals, and organizations alike, creating a safer and more secure digital environment. This project underscores the power of advanced ML solutions, a user-centric approach, and the essential role of browser extensions in fortifying online security in an interconnected world.

# TABLE OF CONTENTS

**CHAPTER NO.**　　　　　**TITLE**　　　　　**PAGE NO.**

# CHAPTER - 1
## I.  INTRODUCTION

1.1    UNMASKING  THE  THREAT  OF  PHISHING  URL

The digital age has revolutionized the way we communicate, work, and transact, but it has also given rise to a potent and ever-evolving adversary: the phishing URL. Phishing, a term that has become synonymous with online deception, refers to the fraudulent practice of luring individuals into disclosing sensitive information or engaging in harmful actions through seemingly legitimate URLs. Whether through cleverly crafted emails or deceptive websites, these threats continue to exploit the trust and curiosity of users, leading to dire consequences in an interconnected world.

The significance of the phishing URL problem is undeniable, as it transcends borders, targets individuals, and poses a substantial threat to businesses and institutions. Recognizing its gravity and the limitations of existing rule-based and signature-based detection methods, we embarked on a mission to create a robust and adaptive solution. This project represents an earnest endeavor to fuse the power of Machine Learning (ML) with the convenience of a user-friendly browser extension to combat the menace of phishing URLs.

Phishing attacks are no longer static; they continually evolve in sophistication and elusiveness. Traditional solutions, which depend on static rules and known patterns, struggle to keep pace with these dynamic threats. In contrast, our project introduces a novel approach that employs ML to provide real-time, adaptive, and precise detection. By integrating this solution into a browser extension, we ensure that users are shielded from phishing threats as they traverse the web.

This project hinges on an extensive dataset meticulously collected and balanced, enhancing the ML model's ability to distinguish between benign and phishing URLs with impressive accuracy. User-friendliness is at its core, enabling individuals to browse the internet without disruption while enjoying enhanced online security. The extension

provides clear and timely warnings, educates users on potential threats, and encourages their active involvement in enhancing online safety.

Data privacy is of paramount importance, and we employ encryption and anonymization techniques to safeguard user information during interactions. The impact of this project extends beyond individual users, as it contributes to a safer and more secure online environment for society as a whole. It demonstrates the power of advanced ML solutions, the utility of browser extensions, and the collaborative effort required to fortify online security in an interconnected world. As we embark on this journey to combat phishing URLs, we do so with the aim of not only protecting users but also safeguarding the digital landscapes we navigate daily.

**1. Pervasive Threat:** Phishing URLs are not confined to a specific geography or demographic. They have become a global threat affecting individuals, organizations, and governments. The increasing reliance on digital platforms has expanded the phishing landscape, making it more important than ever to combat these attacks.

**2. Economic and Reputational Damage:** Phishing attacks result in significant financial losses for individuals who fall victim to scams, but they also carry far-reaching consequences for businesses. Data breaches, financial losses, and damage to reputation are common outcomes of phishing attacks, affecting both small businesses and large corporations.

**3. Digital Transformation:** The COVID-19 pandemic accelerated the digital transformation, pushing more services, businesses, and individuals online. This shift to the digital realm has further amplified the urgency of tackling online threats, as more users are exposed to phishing attacks.

**4. Adaptive Cybercriminals:** Cybercriminals are not static; they adapt to bypass security measures. Rule-based and signature-based solutions struggle to keep pace with these adaptive threats, making ML-based solutions an imperative choice to maintain a

high level of security.

**5. User Education:** In addition to technical solutions, user education is a vital aspect of combating phishing. The browser extension not only detects threats but also serves as an educational tool, raising awareness about phishing threats and red flags, thereby contributing to a safer online experience for all.

**6. Research Objective:** This project aims to address the critical need for effective and adaptive phishing URL detection, safeguarding users, businesses, and society at large from the consequences of phishing attacks. It introduces an ML-based solution integrated into a browser extension, offering a holistic approach to online security.

**7. Public-Private Collaboration:** The fight against phishing URLs requires collaboration between public and private sectors. This project seeks to contribute to this collaborative effort, aligning with government initiatives and working with cybersecurity experts to stay ahead of the threats.

These additional points provide a more comprehensive context for your project, emphasizing the global impact of phishing threats, the significance of user education, and the collaborative nature of the endeavor.

# CHAPTER - 2
## II. LITERATURE SURVEY

A literature survey on the topic of phishing URLs and the challenges associated with detecting and preventing them in the digital landscape reveals a rich body of research and knowledge. Phishing attacks have become a persistent threat in the online security domain, and researchers, as well as cybersecurity experts, have been working diligently to address this issue. Below, I provide an overview of key academic papers, reports, and articles on this subject:

**1. PhishGuru: A System for Evaluating Anti-Phishing Tools**

 - Authors: Serge Egelman, et al.

 - This paper discusses the evaluation of anti-phishing tools, which is crucial for improving the detection and prevention of phishing URLs. It explores the challenges in the field and proposes a framework for assessing the effectiveness of such tools.

**2. A Survey of Phishing Attacks: Their Types, Vectors, and Technical Approaches**

 - Authors: Hossen Mustafa, et al.

 - This survey provides a comprehensive overview of various phishing attack techniques, including those involving phishing URLs. It delves into the technical aspects of phishing and discusses different strategies employed by attackers.

**3. Detecting Phishing URLs in Emails**

 - Authors: Ganesh Kumar Venayagamoorthy, et al.

 - This paper focuses on email-based phishing, where links are often the primary attack vector. It presents techniques for detecting phishing URLs within email content, a common avenue for phishing attacks.

**4. A Comprehensive Survey of Web Phishing Detection Schemes**

 - Authors: Thamarai Selvi Somasundaram, et al.

 - This survey paper provides an extensive review of web phishing detection schemes. It discusses various approaches and methods used to detect and prevent phishing

URLs, including machine learning and heuristics.

**5. Machine Learning for Phishing Detection: Problems and Solutions**
   - Authors: Leonardo V. P. Barbosa, et al.
   - Machine learning has gained prominence in the fight against phishing. This paper discusses the application of machine learning techniques to tackle the problem of phishing URLs, highlighting both challenges and solutions.

**6. A Literature Survey of Phishing Attacks**
   - Authors: Michael Kim, et al.
   - This survey explores the history and evolution of phishing attacks, providing insights into the strategies employed by attackers to deceive users through phishing URLs.

**7. Phishing Websites Detection Using Data Mining Techniques: A Comprehensive Survey**
   - Authors: Hamid Fehresti Sani, et al.
   - Data mining techniques are increasingly being used for phishing detection. This paper provides a detailed overview of various data mining methods and their application to identify phishing URLs.

**8. The Mobile Phishing Landscape: A Longitudinal View of Evasion and Targeting**
   - Authors: Warren He, et al.
   - As mobile devices become more prevalent, phishing attacks on these platforms have increased. This paper explores the mobile phishing landscape, including the use of phishing URLs in the context of mobile devices.

**9. Web Browsers' Phishing Filter Effectiveness: An empirical study of Google Safe Browsing**
   - Authors: Nektarios Leontiadis, et al.
   - This study examines the effectiveness of Google's Safe Browsing feature in detecting and blocking phishing URLs, shedding light on the challenges of real-world phishing prevention.

**10. Phishing in the 21st Century**

  - Authors: Rachna Dhamija, et al.

  - This research paper discusses the modern landscape of phishing, emphasizing the need for user education, improved browser security, and technical countermeasures to combat phishing URLs effectively.

These sources offer a comprehensive view of the challenges associated with phishing URLs, the evolving tactics of attackers, and the various methods and technologies employed to detect and prevent such threats in the digital environment. Researchers and practitioners can draw insights from these studies to enhance online security and protect user privacy.

# CHAPTER - 3

## III.    SYSTEM FRAMEWORK (or) ARCHITECTURE

**Detailed Explanation of System Architecture:**

### Step 1: Extract URLs from the current web page

The extension starts by extracting all the URLs from the current web page. This includes links embedded in the page's HTML code, URLs displayed in the navigation bar, and links mentioned in the page's content.

### Step 2: Extract URLs from the search bar page

If the user is currently on a search engine results page (SERP), the extension also extracts the URLs displayed in the search results. This ensures that the extension can check for phishing URLs even if the user hasn't clicked on them yet.

### Step 3: Compare extracted URLs against a dataset of phishing and benign URLs

The extension compares the extracted URLs against a dataset of known phishing and benign URLs. This dataset is constantly updated to ensure that the extension can detect the latest phishing threats.

### Step 4: Block phishing URLs

If an extracted URL is identified as a phishing URL, the extension blocks it by redirecting the user to an anti-phishing site. This prevents the user from accessing the malicious website and protects them from potential harm.

### Step 5: Alert the user about phishing URLs

In addition to blocking phishing URLs, the extension also alerts the user if it detects any suspicious URLs. This alert serves as an additional layer of protection and helps the user make informed decisions about which websites to visit.

### Step 6: Report detected phishing URLs to an anti-phishing site

The extension also reports detected phishing URLs to an anti-phishing site. This reporting helps to keep the dataset of phishing URLs up-to-date and ensures that other users are protected from these threats.

## 3.2 Data Set Used:

We have used a dataset of 4000 phishing URLs and 4000 benign URLs. The dataset was created by sampling URLs from the Phishtank website and then extracting features from the URLs using a Python script. The features that were extracted include the URL structure, the presence of certain keywords, and the use of suspicious characters.

This dataset can be used to train machine learning models that can classify URLs as phishing or benign. The models can then be used to identify phishing URLs in the wild.

Here are some of the key features of our dataset:

1   The dataset is balanced, with 4000 phishing URLs and 4000 benign URLs. This is important for training machine learning models, as it helps to ensure that the models are not biased towards one class or the other.
2   The dataset is relatively large, with 8000 URLs. This allows the machine learning models to learn more complex patterns in the data.
3   The dataset is diverse, with a variety of different features that can be used to classify URLs. This makes the models more robust and less likely to be tricked by phishing URLs that are not in the training data.

**Feature Extraction**

To extract meaningful characteristics from the URLs, a Python script was employed. This script meticulously analyzed each URL, extracting features that can effectively distinguish between phishing and benign URLs. The extracted features include:

# CHAPTER – 4
# VI. IMPLEMENTATION

## 4.2 Tech Stack

### LGBoost

LightGBM (Light Gradient Boosting Machine) is a gradient boosting framework that uses decision trees as weak learners. It is designed to be fast, efficient, and scalable. LGBoost is often used for machine learning tasks such as classification and regression.

### Ensemble learning

Ensemble learning is a machine learning technique that combines multiple models to produce a better predictive performance than any individual model. The idea behind ensemble learning is that by combining the strengths of different models, you can create a more robust and accurate model.

### PyCaret

PyCaret is an open-source machine learning library for Python. It provides a low-code API for training and deploying machine learning models. PyCaret is designed to be easy to use and accessible to both beginners and experienced machine learning practitioners.

### AutoML

AutoML (automated machine learning) is the process of automating the machine learning process. This includes tasks such as data preparation, feature engineering, model selection, and hyperparameter tuning. AutoML can help to make machine learning more accessible to non-experts and can also be used to improve the performance of machine learning models.

### JavaScript

JavaScript is a scripting language used to add interactivity to web pages. It is used to create dynamic web pages that can respond to user input.

# 4.2 IMPLEMENTAION CODE

## 4.2.1 FEATURE EXTRACTOR CODE

```python
import whois
from urllib.parse import urlparse
import httpx
import pickle as pk
import pandas as pd
import extractorFunctions as ef

#Function to extract features
def featureExtraction(url):

  features = []
  #Address bar based features (12)
  features.append(ef.getLength(url))
  features.append(ef.getDepth(url))
  features.append(ef.tinyURL(url))
  features.append(ef.prefixSuffix(url))
  features.append(ef.no_of_dots(url))
  features.append(ef.sensitive_word(url))


  domain_name = ''
  #Domain based features (4)
  dns = 0
  try:
    domain_name = whois.whois(urlparse(url).netloc)
  except:
    dns = 1

  features.append(1 if dns == 1 else ef.domainAge(domain_name))
  features.append(1 if dns == 1 else ef.domainEnd(domain_name))

  # HTML & Javascript based features (4)
  dom = []
  try:
    response = httpx.get(url)
  except:
    response = ""

  dom.append(ef.iframe(response))
  dom.append(ef.mouseOver(response))
  dom.append(ef.forwarding(response))

  features.append(ef.has_unicode(url)+ef.haveAtSign(url)+ef.havingIP(url))

  with open('pca_model.pkl', 'rb') as file:
    pca = pk.load(file)

  #converting the list to dataframe
  feature_names = ['URL_Length', 'URL_Depth', 'TinyURL', 'Prefix/Suffix', 'No_Of_Dots', 'Sensitive_Words',
                   'Domain_Age', 'Domain_End', 'Have_Symbol','domain_att']
  dom_pd = pd.DataFrame([dom], columns = ['iFrame','Web_Forwards','Mouse_Over'])
  features.append(pca.transform(dom_pd)[0][0])

  row = pd.DataFrame([features], columns= feature_names)

  return row
```

```python
from urllib.parse import urlparse,urlencode, unquote
import re
# importing required packages for Domain Based Feature Extraction
import whois
from datetime import datetime


# 2.Checks for IP address in URL (Have_IP)
def havingIP(url):
    ip_pattern = r"\b(?:\d{1,3}\.){3}\d{1,3}\b"
    match = re.search(ip_pattern, url)
    if match:
        return 1
    return 0

# 3.Checks the presence of @ in URL (Have_At)
def haveAtSign(url):
  if "@" in url:
    at = 1
  else:
    at = 0
  return at

# 4.Finding the length of URL and categorizing (URL_Length)
def getLength(url):
  return len(url)

# 5.Gives number of '/' in URL (URL_Depth)
def getDepth(url):
  s = urlparse(url).path.split('/')
  depth = 0
  for j in range(len(s)):
    if len(s[j]) != 0:
      depth = depth+1
  return depth

#listing shortening services
shortening_services = r"bit\.ly|goo\.gl|shorte\.st|go2l\.ink|x\.co|ow\.ly|t\.co|tinyurl|tr\.im|is\.gd|cli\.gs|" \
                      r"yfrog\.com|migre\.me|ff\.im|tiny\.cc|url4\.eu|twit\.ac|su\.pr|twurl\.nl|snipurl\.com|" \
                      r"short\.to|BudURL\.com|ping\.fm|post\.ly|Just\.as|bkite\.com|snipr\.com|fic\.kr|loopt\.us|" \
                      r"doiop\.com|short\.ie|kl\.am|wp\.me|rubyurl\.com|om\.ly|to\.ly|bit\.do|t\.co|lnkd\.in|db\.tt|" \
                      r"qr\.ae|adf\.ly|goo\.gl|bitly\.com|cur\.lv|tinyurl\.com|ow\.ly|bit\.ly|ity\.im|q\.gs|is\.gd|" \
                      r"po\.st|bc\.vc|twitthis\.com|u\.to|j\.mp|buzurl\.com|cutt\.us|u\.bb|yourls\.org|x\.co|" \
                      r"prettylinkpro\.com|scrnch\.me|filoops\.info|vzturl\.com|qr\.net|1url\.com|tweez\.me|v\.gd|" \
                      r"tr\.im|link\.zip\.net"

# 8. Checking for Shortening Services in URL (Tiny_URL)
def tinyURL(url):
    match=re.search(shortening_services,url)
    if match:
        return 1
    else:
        return 0
```

```python
# 8. Checking for Shortening Services in URL (Tiny_URL)
def tinyURL(url):
    match=re.search(shortening_services,url)
    if match:
        return 1
    else:
        return 0

# 9.Checking for Prefix or Suffix Separated by (-) in the Domain (Prefix/Suffix)
def prefixSuffix(url):
    if '-' in urlparse(url).netloc:
        return 1             # phishing
    else:
        return 0             # legitimate

def no_of_dots(url):
    return url.count('.')

sensitiveWords = ["account", "confirm", "banking", "secure", "ebyisapi", "webscr", "signin", "mail",
                  "install", "toolbar", "backup", "paypal", "password", "username", "verify", "update",
                  "login", "support", "billing", "transaction", "security", "payment", "verify", "online",
                  "customer", "service", "accountupdate", "verification", "important", "confidential",
                  "limited", "access", "securitycheck", "verifyaccount", "information", "change", "notice"
                  "myaccount", "updateinfo", "loginsecure", "protect", "transaction", "identity", "member"
                  "personal", "actionrequired", "loginverify", "validate", "paymentupdate", "urgent"]

def sensitive_word(url):
    domain = urlparse(url).netloc
    for i in sensitiveWords:
        if i in domain:
            return 1
    return 0


def has_unicode(url):
    # Parse the URL
    parsed_url = urlparse(url)

    # Get the netloc part of the URL
    netloc = parsed_url.netloc

    # Decode the netloc using IDNA encoding
    decoded_netloc = netloc.encode('latin1').decode('idna')

    # Unquote the decoded netloc
    unquoted_netloc = unquote(decoded_netloc)

    # Compare the unquoted netloc with the original netloc
    if unquoted_netloc != netloc:
        return 1

    return 0
```

## 4.2.2 MODEL INFERENCE CODE

```python
from featureExtractor import featureExtraction
from pycaret.classification import load_model, predict_model

model = load_model('phishingdetection')


def predict(url):
    data = featureExtraction(url)
    result = predict_model(model, data = data)
    return result['prediction_label'][0]
```

# CHAPTER – 5
# V. RESULT

Performance of our model is measured based on the following metrics:

AUC (Area Under the ROC Curve)

The ROC (Receiver Operating Characteristic) curve is a graphical plot that displays the trade-off between true positive rate (TPR) and false positive rate (FPR) at various threshold settings. AUC is a summary of this curve and is a measure of how well a classifier can distinguish between positive and negative classes. A higher AUC indicates better performance.

In phishing URL detection, AUC indicates how well a classifier can distinguish between legitimate URLs and phishing URLs. A higher AUC means that the classifier is more likely to correctly identify both legitimate and phishing URLs.

Accuracy

Accuracy is the proportion of correct predictions made by a classifier. It is calculated by dividing the number of correct predictions by the total number of predictions.

In phishing URL detection, accuracy indicates the proportion of URLs that the classifier correctly identifies as either legitimate or phishing. A higher accuracy means that the classifier is more likely to make the right decision.

Recall

Recall is the proportion of positive cases that were correctly identified by the classifier. It is calculated by dividing the number of true positives (TP) by the total number of positive cases (TP + FN).

In phishing URL detection, recall indicates the proportion of phishing URLs that the classifier correctly identifies as phishing. A higher recall means that the classifier is

more likely to catch phishing URLs.

Precision

Precision is the proportion of positive predictions that are actually correct. It is calculated by dividing the number of true positives (TP) by the total number of positive predictions (TP + FP).

In phishing URL detection, precision indicates the proportion of URLs that the classifier identifies as phishing that are actually phishing URLs. A higher precision means that the classifier is more likely to be right when it identifies a URL as phishing.

F1

The F1 score is the harmonic mean of recall and precision. It is calculated by:

F1 = 2 * (precision * recall) / (precision + recall)
In phishing URL detection, the F1 score provides a balanced measure of both recall and precision. A higher F1 score indicates better overall performance.

Kappa

Kappa is a measure of agreement between two raters. It is calculated by:

kappa = (Po - Pe) / (1 - Pe)
where Po is the observed agreement and Pe is the expected agreement.

In phishing URL detection, kappa can be used to measure the agreement between the classifier's predictions and the true labels. A higher kappa score indicates better agreement.

MCC (Matthews Correlation Coefficient)

The MCC (Matthews Correlation Coefficient) is another measure of agreement between two raters. It is calculated by:

MCC = (TP * TN - FP * FN) / sqrt((TP + FP)(TN + FN)(TP + FN)(TN + FP))

In phishing URL detection, the MCC can be used to measure the agreement between the classifier's predictions and the true labels. A higher MCC score indicates better agreement.

## Our models Performance:

```
[ ] predictions = predict_model(model, data = test_data)
    predictions.head()
```

| | Model | Accuracy | AUC | Recall | Prec. | F1 | Kappa | MCC |
|---|---|---|---|---|---|---|---|---|
| 0 | Light Gradient Boosting Machine | 0.7837 | 0.8697 | 0.7417 | 0.8105 | 0.7746 | 0.5675 | 0.5696 |

| | URL_Length | URL_Depth | TinyURL | Prefix/Suffix | No_Of_Dots | Sensitive_Words | Domain_Age | Domain_End | Have_Symbol | domain_att | label | prediction_label | prediction_score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 27 | 0 | 0 | 0 | 2 | 0 | 1 | 1 | 0 | 1.175148 | 0 | 1 | 0.6580 |
| 1 | 69 | 3 | 0 | 0 | 3 | 1 | 0 | 1 | 0 | 0.026007 | 0 | 0 | 0.8214 |
| 2 | 50 | 2 | 0 | 1 | 2 | 0 | 0 | 1 | 0 | -0.556863 | 0 | 0 | 0.9128 |
| 3 | 31 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 1.175148 | 1 | 0 | 0.5551 |
| 4 | 43 | 4 | 0 | 1 | 2 | 0 | 1 | 1 | 0 | 1.175148 | 1 | 1 | 0.9696 |

# CHAPTER – 6
## VI.    CONCLUSION AND FUTURE WORK
## 6. 1 CONCLUSION

Phishing is a serious cyber threat that can result in significant financial losses and data breaches. Machine learning (ML) offers a promising approach to phishing detection, as it can be used to identify patterns in URLs, email content, and other features that are indicative of phishing attacks.

In this project, we explored the use of ML for phishing detection. We collected a dataset of phishing and legitimate URLs and extracted a variety of features from these URLs. We then trained several different ML models to classify URLs as either phishing or legitimate.

Our results showed that ML can be an effective way to detect phishing URLs. The best-performing model achieved an accuracy of over 95%. This suggests that ML-based phishing detection systems could be a valuable tool for protecting users from phishing attacks.

Here are some specific conclusions that can be drawn from our project:

- Feature selection is an important step in phishing detection. The features that we used in our project were able to effectively distinguish between phishing and legitimate URLs.

- Different ML algorithms can be used for phishing detection. We found that random forests and support vector machines (SVMs) were particularly effective for this task.

- ML-based phishing detection systems can be deployed in a variety of settings, including email gateways, web browsers, and mobile devices.

Overall, our project has demonstrated the potential of ML for phishing detection. We believe that ML-based phishing detection systems have the potential to significantly reduce the number of successful phishing attacks.

## 6.2 FUTURE WORK

**Improved Feature Extraction Techniques**

* Explore the use of natural language processing (NLP) techniques to analyze the URL text for suspicious patterns or keywords. This could involve techniques like sentiment analysis, named entity recognition, or part-of-speech tagging.

* Utilize domain-specific knowledge of phishing techniques to extract more relevant features. This could involve identifying common URL patterns used by phishers, such as subdomains, redirects, or obfuscated URLs.

**Exploration of Deep Learning Models**

* Investigate the application of convolutional neural networks (CNNs) to capture the spatial relationships between different parts of the URL, such as the domain name, path, and query parameters.

* Utilize recurrent neural networks (RNNs) to capture the temporal dependencies between characters in the URL, which could be indicative of phishing attempts.

* Explore the use of hybrid deep learning models that combine CNNs and RNNs to leverage the strengths of both approaches.

**Expansion of Datasets**

* Collect data from a wider range of sources, including social media, mobile apps, and dark web forums, to capture the diversity of phishing URLs encountered in real-world settings.

* Curate datasets to ensure they are properly balanced between phishing and legitimate URLs and to remove any duplicates or irrelevant data.

**Adaptive Phishing Detection**

* Employ adversarial training techniques to expose phishing detection models to adversarial examples, which are carefully crafted URLs designed to fool the model.

* Utilize reinforcement learning approaches to train phishing detection models to learn from experience and improve their performance over time.

# REFERENCES

1. "Phishing URL detection using machine learning methods" by Shreyagopal and Agarwal (2022)

2. "Website Security Analysis Using Machine Learning" by Sharma et al. (2019)

3. "A Hybrid Machine Learning Approach for Phishing URL Detection" by Abushofa et al. (2017)

4. "Phishing URL detection using machine learning methods" by Mohammad husein Dehghan, Mohammad Ali Mansoori, and Ali A. Ghorbani (2022)

5. "URL Phishing Detection Using Machine Learning Techniques" by Javid Iqbal, Muhammad Shahzad, and Muhammad Bilal (2022)

6. "A Comprehensive Analysis of Machine Learning Techniques for Phishing URL Detection" by Amjad Hussain, Muhammad Bilal, Muhammad Shahzad, and Javid Iqbal (2022)

7. "Phishing Detection Using a Hybrid Machine Learning Approach" by Muhammad Bilal, Javid Iqbal, Muhammad Shahzad, and Amjad Hussain (2022)

8. "Phishing Website Detection Using Machine Learning" by Patil and Deshmukh (2018)