

# PHISHING DETECTION SYSTEM

## GROUP MEMBERS:

Kevin John

Neeraj Sachi

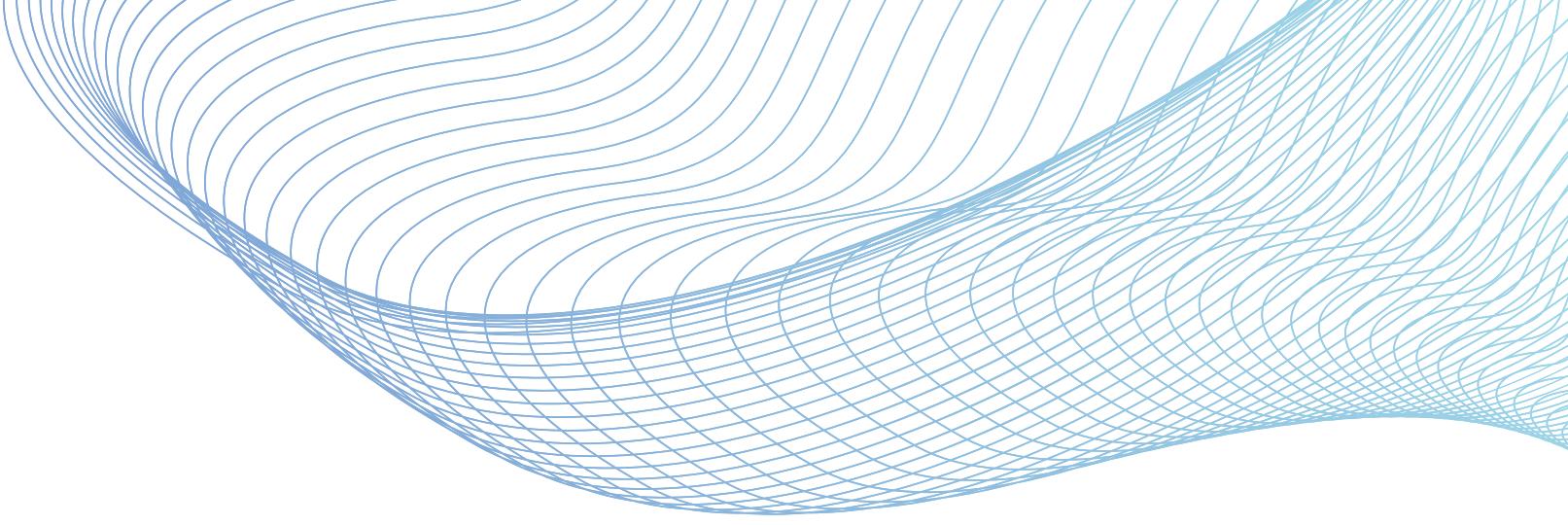
Mohammed Midhlaj V

Rumaisa Shajahan

GUIDED BY:  
Ms. Asha Raj



# CONTENTS

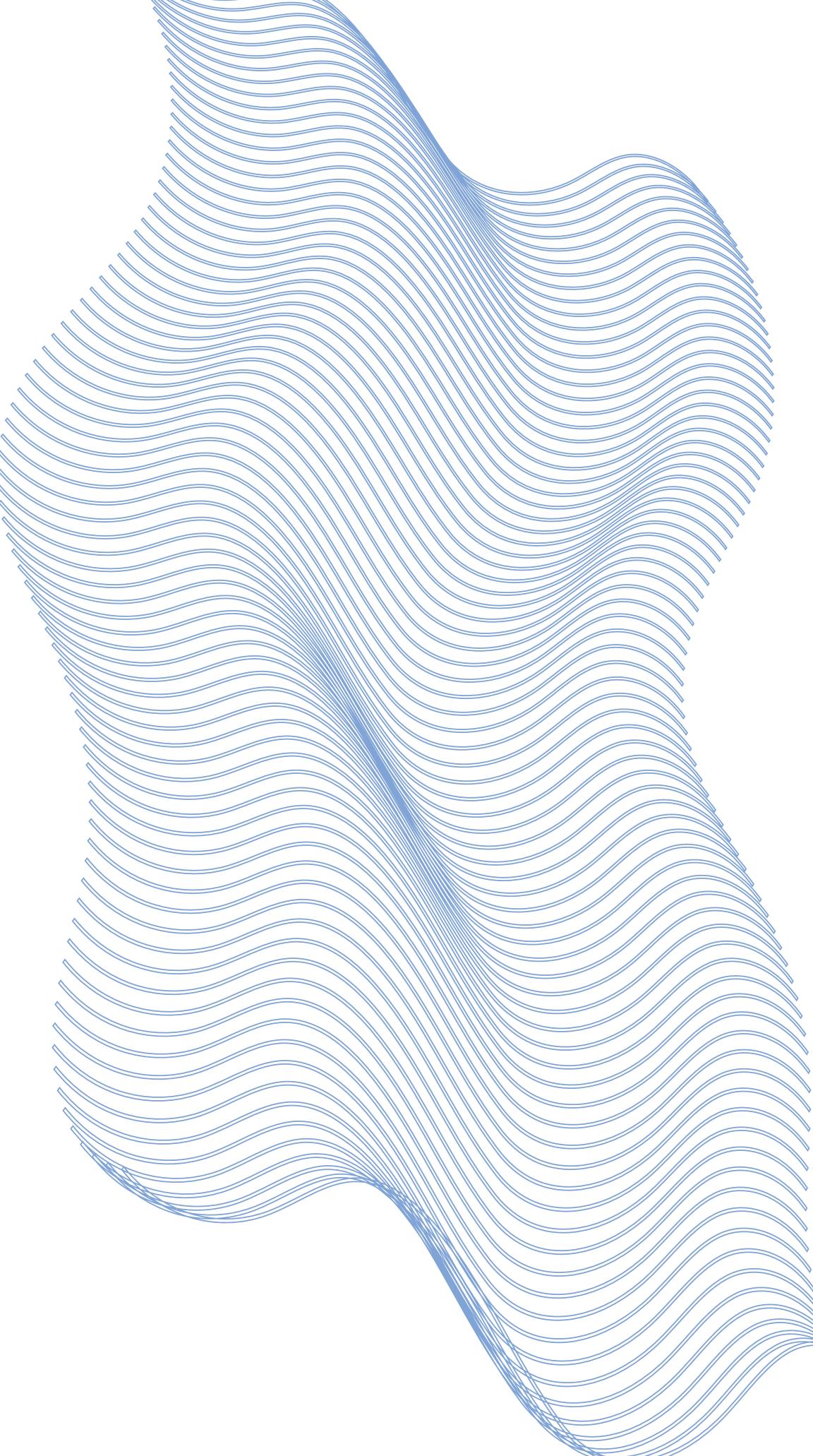
- 
1. Introduction
  2. Problem Statement
  3. Abstract
  4. Current System
  5. Proposed System
  6. Objectives
  7. Datasets
  8. Architecture
  9. Modules
  10. Methodology
  11. Conclusion
  12. Reference

# INTRODUCTION

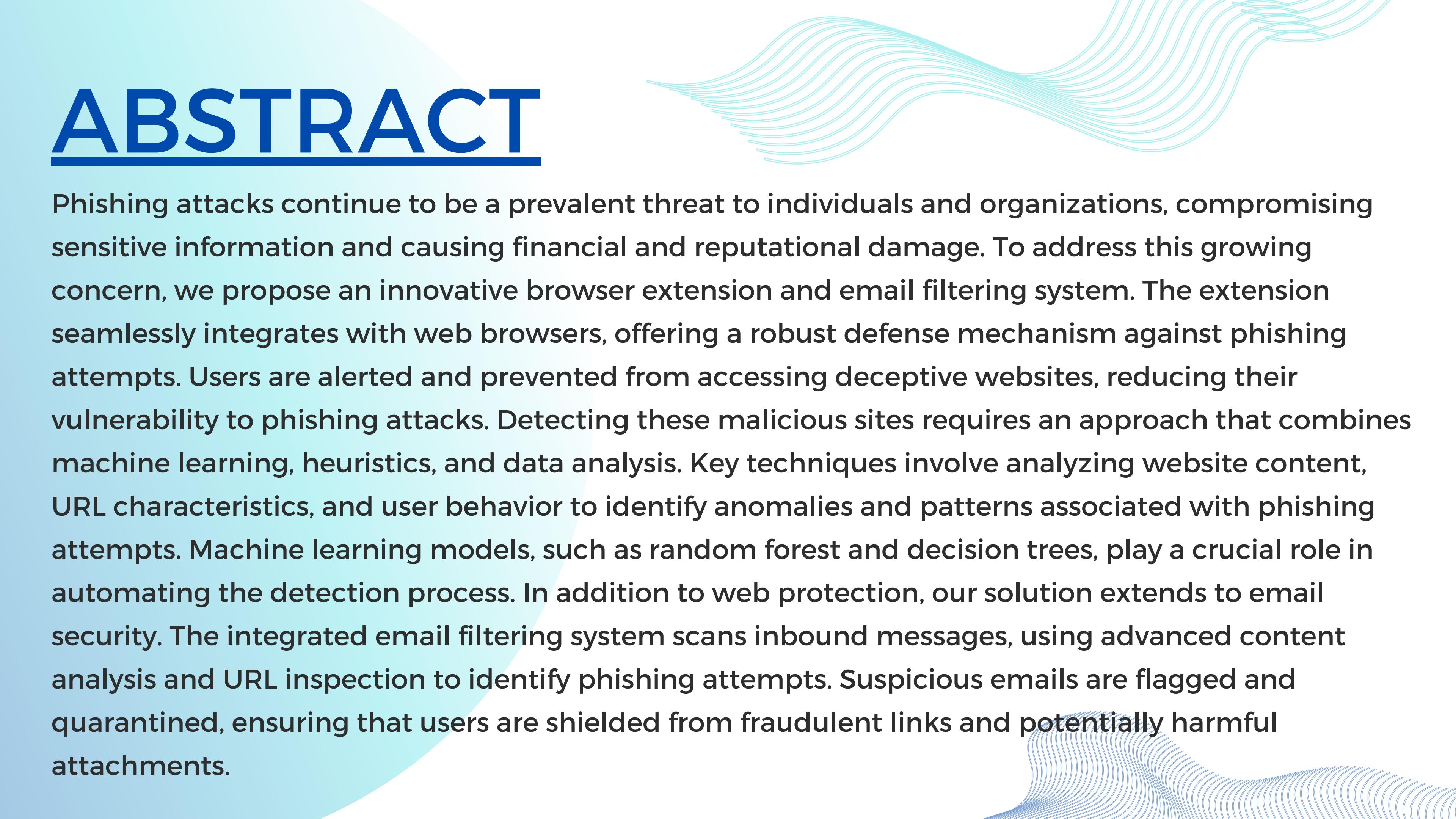
Phishing is a cybercrime tactic used to deceive individuals and organizations into revealing sensitive information. This presentation explores the use of machine learning for phishing detection.

# PROBLEM STATEMENT

To develop an extension for detecting phishing websites and emails that can accurately identify and flag potentially malicious websites and emails in real-time.



# ABSTRACT



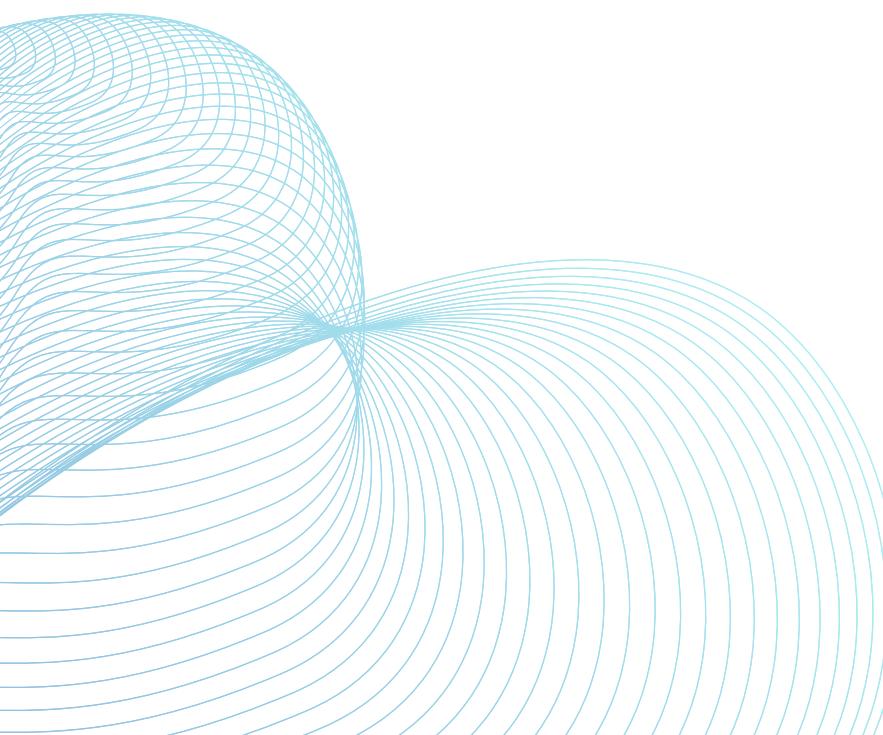
Phishing attacks continue to be a prevalent threat to individuals and organizations, compromising sensitive information and causing financial and reputational damage. To address this growing concern, we propose an innovative browser extension and email filtering system. The extension seamlessly integrates with web browsers, offering a robust defense mechanism against phishing attempts. Users are alerted and prevented from accessing deceptive websites, reducing their vulnerability to phishing attacks. Detecting these malicious sites requires an approach that combines machine learning, heuristics, and data analysis. Key techniques involve analyzing website content, URL characteristics, and user behavior to identify anomalies and patterns associated with phishing attempts. Machine learning models, such as random forest and decision trees, play a crucial role in automating the detection process. In addition to web protection, our solution extends to email security. The integrated email filtering system scans inbound messages, using advanced content analysis and URL inspection to identify phishing attempts. Suspicious emails are flagged and quarantined, ensuring that users are shielded from fraudulent links and potentially harmful attachments.

# CURRENT SYSTEM

- **Blacklist-Based Detection:**
  - Maintains a list of known phishing URLs or domains.
  - Ineffective against zero-day phishing attacks and frequently changing domains.
- **High False Positive Rates:**
  - Current systems often generate false positives, causing inconvenience to legitimate users.
- **Dependency on User Reports:**
  - Often relies on user-generated reports of suspected phishing sites, leading to delayed responses.

# PROPOSED SYSTEM

- Utilizes machine learning algorithms like XGBoost and Random Forest .
- Integrates real-time analysis of user behavior and website content.
- Incorporates deep URL analysis.



# OBJECTIVES

The primary objective of the extension is to protect individuals and organizations from the significant financial, reputational, and security risks associated with phishing attacks by identifying and blocking fraudulent websites and educating users about the threat.

- Protection Against Fraud
- Data Security
- Prevent unauthorised access
- Reduced Cybersecurity risks
- Early Threat Detection

# DATASET

**Source:**

**<https://www.kaggle.com/datasets/sid321axn/malicious-urls-dataset/data>**

# DATASET

## Raw Dataset

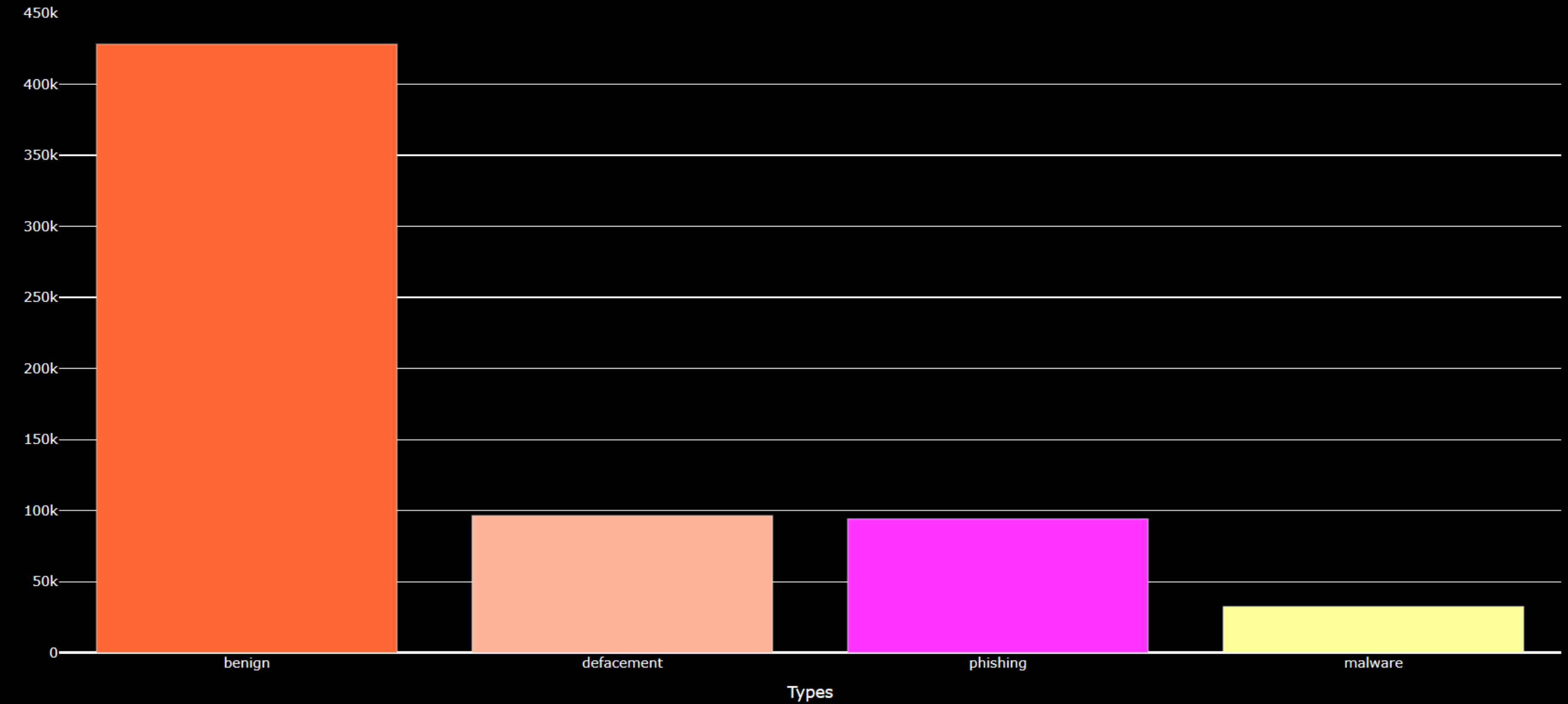
- br-icloud.com.br,phishing
- mp3raid.com/music/krizz\_kaliko.html,benign
- bopsecrets.org/rexroth/cr/1.htm,benign
- [http://www.garage-pirenne.be/index.php?option=com\\_content&view=article&id=70&vsig70\\_0=15](http://www.garage-pirenne.be/index.php?option=com_content&view=article&id=70&vsig70_0=15),defacement
- myspace.com/video/vid/30602581,benign
- <http://www.lebensmittel-ueberwachung.de/index.php/aktuelles.1>,defacement
- <http://www.szabadmunkaero.hu/cimoldal.html?start=12>,defacement
- <http://buzzfil.net/m/show-art/ils-etaient-loin-de-s-imaginer-que-le-hibou-allait-faire-ceci-quand-ils-filmaient-2.html>,benign
- espn.go.com/nba/player/\_/id/3457/brandон-rush,benign
- <http://www.vnic.co/khach-hang.html>,defacement

# DATASET

## Raw Dataset

- <http://www.824555.com/app/member/SportOption.php?uid=guest&langx=gb>,malware
- <http://www.raci.it/component/user/reset.html>,defacement
- <https://docs.google.com/spreadsheet/viewform?formkey=dGg2Z1ICUHISdjIITVNRUW50TFIzSkE6MQ>,phishing
- [psychology.wikia.com/wiki/Phonemes](http://psychology.wikia.com/wiki/Phonemes),benign
- [info.centriq.com/content/fivereasons](http://info.centriq.com/content/fivereasons),benign
- lowery,benign [infinitysw.com/](http://infinitysw.com/),benign [strawberrycreekgardens.com/](http://strawberrycreekgardens.com/),benign
- [peopleaz.org/firstname/Fowlkes/6](http://peopleaz.org/firstname/Fowlkes/6),benign
- [orensamtraktrip.wordpress.com/](http://orensamtraktrip.wordpress.com/),benign
- [retajconsultancy.com](http://retajconsultancy.com),phishing
- [alexa.com/whatshot?q=james+arness](http://alexa.com/whatshot?q=james+arness),benign

## Count of Different Types of URLs



# DATASET(BEFORE PREPROCESSING)

```
1 url,type
2 br-icloud.com.br,phishing
3 mp3raid.com/music/krizz_kaliko.html,benign
4 bopsecrets.org/rexroth/cr/1.htm,benign
5 http://www.garage-pirenne.be/index.php?option=com_content&view=article&id=70&vsig70_0=15,defacement
6 http://adventure-nicaragua.net/index.php?option=com_mailto&tmpl=component&link=aHR0cDovL2FkdmVudHVyZS1uaWNhcmFnWUbmV0L2luZGV4LnBocD9
7 http://buzzfil.net/m/show-art/ils-etaient-loin-de-s-imaginer-que-le-hibou-allait-faire-ceci-quand-ils-filmaient-2.html,benign
8 espn.go.com/nba/player/_/id/3457/brandon-rush,benign
9 yourbittorrent.com/?q=anthony-hamilton-soulife,benign
10 http://www.pashminaonline.com/pure-pashminas,defacement
11 allmusic.com/album/crazy-from-the-heat-r16990,benign
12 corporationwiki.com/Ohio/Columbus/frank-s-benson-P3333917.aspx,benign
13 http://www.ikenmijnkunst.nl/index.php/exposities/exposities-2006,defacement
14 myspace.com/video/vid/30602581,benign
15 http://www.lebensmittel-ueberwachung.de/index.php/aktuelles.1,defacement
16 http://www.szabadmunkaero.hu/cimoldal.html?start=12,defacement
17 http://larcadelcarnevale.com/catalogo/palloncini,defacement
18 quickfacts.census.gov/qfd/maps/iowa_map.html,benign
19 nugget.ca/ArticleDisplay.aspx?archive=true&e=1160966,benign
20 uk.linkedin.com/pub/steve-rubenstein/8/718/755,benign
21 http://www.vnic.co/khach-hang.html,defacement
22 baseball-reference.com/players/h/harrige01.shtml,benign
23 signin.eby.de.zukruygxctzmmqi.civpro.co.za,phishing
24 192.com/atoz/people/oakley/patrick/,benign
25 nytimes.com/1998/03/29/style/cuttings-oh-that-brazen-raucous-glorious-hibiscus.html,benign
```

# PREPROCESSING METHODOLOGY

## **Loading the Dataset:**

- The dataset is loaded using the `pd.read_csv` function from the Pandas library

## **Handling Missing Values:**

- The presence of missing values is checked using `urls_data.isnull().sum()`
- Missing values are replaced with 0 using `urls_data.fillna(0, inplace=True)`

## **URL Cleaning:**

- The 'www.' prefix is removed from the 'url' column using `urls_data['url'] = urls_data['url'].replace('www.', "", regex=True).`
- 

## **Mapping Types of Websites:**

- The 'type' column is mapped to numerical values using `urls_data["url_type"] = urls_data["type"].replace(...)`

# DATASET(AFTER PREPROCESSING)

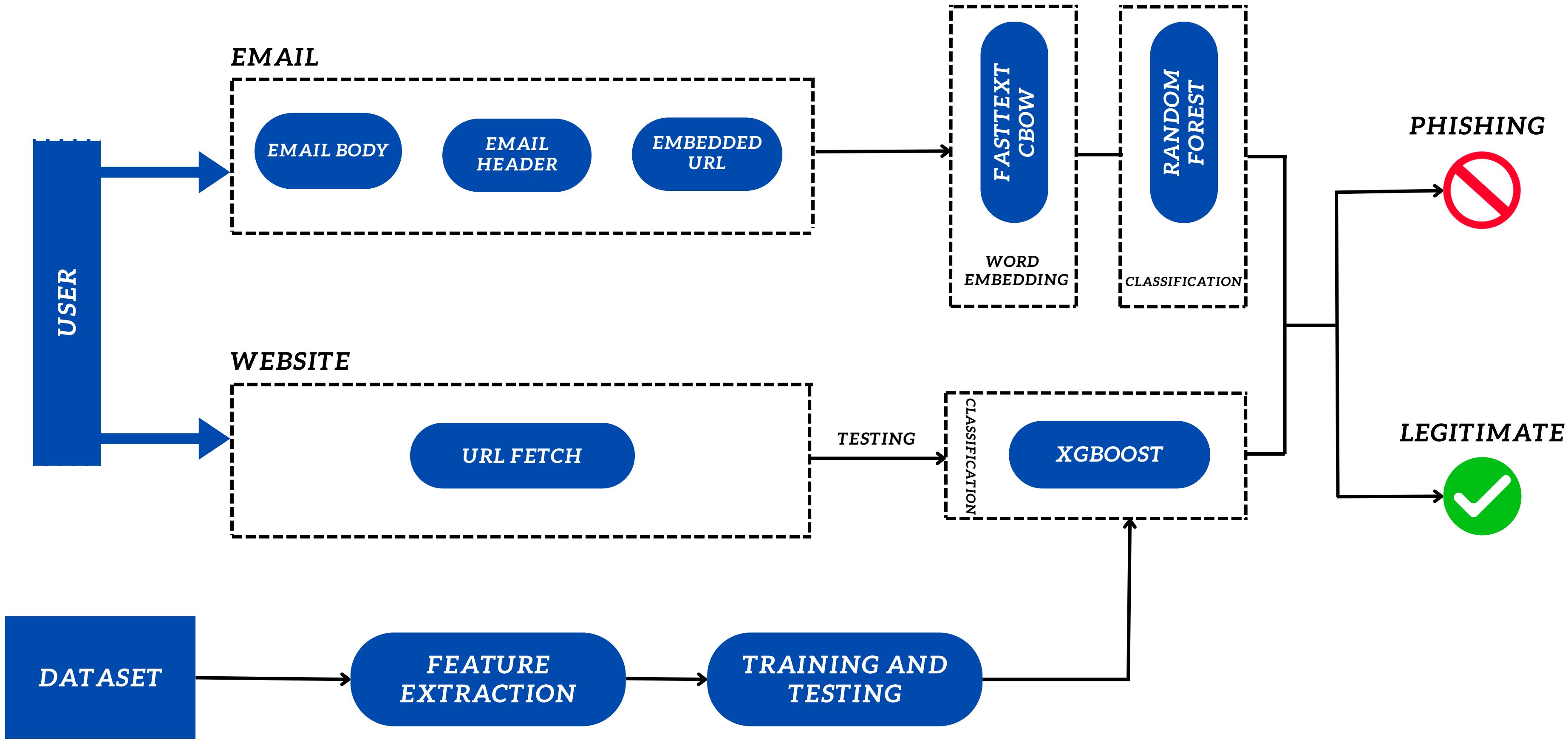
```
1 url,type,url_type,url_len
2 br-icloud.com.br,phishing,2,16
3 mp3raid.com/music/krizz_kaliko.html,benign,0,35
4 bopsecrets.org/rexroth/cr/1.htm,benign,0,31
5 http://garage-pirenne.be/index.php?option=com_content&view=article&id=70&vsig70_0=15,defacement,1,77
6 http://adventure-nicaragua.net/index.php?option=com_mailto&tmpl=component&link=aHR0cDovL2FkdmVudHvyZS1uaWNhcmFndWEubmV0L2luZGV4LnBocD9
7 http://buzzfil.net/m/show-art/ils-etaient-loin-de-s-imaginer-que-le-hibou-allait-faire-ceci-quand-ils-filmaient-2.html,benign,0,111
8 espn.go.com/nba/player/_/id/3457/brandon-rush,benign,0,45
9 yourbittorrent.com/?q=anthony-hamilton-soulife,benign,0,46
10 http://pashminaonline.com/pure-pashminas,defacement,1,33
11 allmusic.com/album/crazy-from-the-heat-r16990,benign,0,45
12 corporationwiki.com/Ohio/Columbus/frank-s-benson-P3333917.aspx,benign,0,62
13 http://ikenmijnkunst.nl/index.php/exposities/exposities-2006,defacement,1,53
14 myspace.com/video/vid/30602581,benign,0,30
15 http://lebensmittel-ueberwachung.de/index.php/aktuelles.1,defacement,1,50
16 http://szabadmunkaero.hu/cimoldal.html?start=12,defacement,1,40
17 http://larcadelcarnevale.com/catalogo/palloncini,defacement,1,41
18 quickfacts.census.gov/qfd/maps/iowa_map.html,benign,0,44
19 nugget.ca/ArticleDisplay.aspx?archive=true&e=1160966,benign,0,52
20 uk.linkedin.com/pub/steve-rubenstein/8/718/755,benign,0,46
21 http://vnic.co/khach-hang.html,defacement,1,23
22 baseball-reference.com/players/h/harrige01.shtml,benign,0,48
23 signin.eby.de.zukruygxctzmmqi.civpro.co.za,phishing,2,42
24 192.com/atoz/people/oakley/patrick/,benign,0,35
25 nytimes.com/1998/03/29/style/cuttings-oh-that-brazen-raucous-glorious-hibiscus.html,benign,0,83
```

# DATASET(AFTER FEATURE EXTRACTION)

url	type	url_type	url_len	pri_domain	letters_count	digits_count	special_chars_shortened	abnormal_url	secure_http	have_ip
<a href="http://br-icloud.com.br">br-icloud.com.br</a>	phishing		2	<a href="http://br-icloud.com.br">br-icloud.com.br</a>	13	0	3	0	0	0
<a href="http://mp3raid.com/mu">mp3raid.com/mu</a>	benign		0	<a href="http://mp3raid.com">mp3raid.com</a>	29	1	5	0	0	0
<a href="http://bopsecrets.org/r">bopsecrets.org/r</a>	benign		0	<a href="http://bopsecrets.org">bopsecrets.org</a>	25	1	5	0	0	0
<a href="http://garage-pire">http://garage-pire</a>	defacement		1	<a href="http://garage-pire">garage-pirennne.k</a>	60	7	17	0	1	0
<a href="http://adventure-">http://adventure-</a>	defacement		1	<a href="http://adventure-nicara">adventure-nicara</a>	199	22	14	0	1	0
<a href="http://buzzfil.net/">http://buzzfil.net/</a>	benign		0	<a href="http://buzzfil.net">buzzfil.net</a>	93	1	24	0	1	0
<a href="http://espn.go.com/nba">espn.go.com/nba</a>	benign		0	<a href="http://espn.go.com">espn.go.com</a>	31	4	10	0	0	0
<a href="http://yourbittorrent.co">yourbittorrent.co</a>	benign		0	<a href="http://yourbittorrent.co">yourbittorrent.co</a>	40	0	6	0	0	0
<a href="http://pashminaonline">http://pashminaonline</a>	defacement		1	<a href="http://pashminaonline.">pashminaonline.</a>	34	0	6	0	1	0
<a href="http://allmusic.com/alb">allmusic.com/alb</a>	benign		0	<a href="http://allmusic.com">allmusic.com</a>	33	5	7	0	0	0
<a href="http://corporationwiki.c">corporationwiki.c</a>	benign		0	<a href="http://corporationwiki.c">corporationwiki.c</a>	47	7	8	0	0	0
<a href="http://ikenmijnkunst.nl">http://ikenmijnkunst.nl</a>	defacement		1	<a href="http://ikenmijnkunst.nl">ikenmijnkunst.nl</a>	47	4	9	0	1	0
<a href="http://myspace.com/vic">myspace.com/vic</a>	benign		0	<a href="http://myspace.com">myspace.com</a>	18	8	4	0	0	0
<a href="http://lebensmittel-ueb">http://lebensmittel-ueb</a>	defacement		1	<a href="http://lebensmittel-ueb">lebensmittel-ueb</a>	47	1	9	0	1	0
<a href="http://szabadmunkaerde">http://szabadmunkaerde</a>	defacement		1	<a href="http://szabadmunkaerde">szabadmunkaerde</a>	37	2	8	0	1	0

# DATASET(AFTER FEATURE EXTRACTION)

<a href="http://szabadmunkaerd">http://szabadmunkaerd</a>	defacement	1	40	<a href="http://szabadmunkaerd">szabadmunkaerd</a>	37	2	8	0	1	0	0	0
<a href="http://larcadelcarneval">http://larcadelcarneval</a>	defacement	1	41	<a href="http://larcadelcarneval">larcadelcarneval</a>	42	0	6	0	1	0	0	0
<a href="http://quickfacts.census">quickfacts.census</a>	benign	0	44	<a href="http://quickfacts.census">quickfacts.census</a>	37	0	7	0	0	0	0	0
<a href="http://nugget.ca/Article">nugget.ca/Article</a>	benign	0	52	<a href="http://nugget.ca">nugget.ca</a>	38	7	7	0	0	0	0	0
<a href="http://uk.linkedin.com/">uk.linkedin.com/</a>	benign	0	46	<a href="http://uk.linkedin.com">uk.linkedin.com</a>	31	7	8	0	0	0	0	0
<a href="http://vnic.co/kha">http://vnic.co/kha</a>	defacement	1	23	<a href="http://vnic.co">vnic.co</a>	23	0	7	0	1	0	0	0
<a href="http://baseball-references.com">baseball-references.com</a>	benign	0	48	<a href="http://baseball-references.com">baseball-references.com</a>	40	2	6	0	0	0	0	0
<a href="http://signin.eby.de.zul">signin.eby.de.zul</a>	phishing	2	42	<a href="http://signin.eby.de.zul">signin.eby.de.zul</a>	36	0	6	0	0	0	0	0
<a href="http://192.com/atoz/people">192.com/atoz/people</a>	benign	0	35	<a href="http://192.com">192.com</a>	26	3	6	0	0	0	0	0
<a href="http://nytimes.com/1993">nytimes.com/1993</a>	benign	0	83	<a href="http://nytimes.com">nytimes.com</a>	62	8	13	0	0	0	0	0
<a href="http://escholarship.org">escholarship.org</a>	benign	0	33	<a href="http://escholarship.org">escholarship.org</a>	24	5	4	0	0	0	0	0
<a href="http://songfacts.com/d">songfacts.com/d</a>	benign	0	33	<a href="http://songfacts.com">songfacts.com</a>	23	5	5	0	0	0	0	0
<a href="http://casamanana.org">casamanana.org</a>	benign	0	30	<a href="http://casamanana.org">casamanana.org</a>	26	0	4	0	0	0	0	0
<a href="http://hollywoodlife.com">hollywoodlife.com</a>	benign	0	78	<a href="http://hollywoodlife.com">hollywoodlife.com</a>	58	12	15	0	1	0	0	0
<a href="http://marketingbyinter">http://marketingbyinter</a>	phishing	2	60	<a href="http://marketingbyinter">marketingbyinter</a>	43	17	7	0	1	0	0	0
<a href="http://en.wikipedia.org">en.wikipedia.org</a>	benign	0	34	<a href="http://en.wikipedia.org">en.wikipedia.org</a>	29	0	5	0	0	0	0	0
<a href="http://soaps.sheknows">soaps.sheknows</a>	benign	0	81	<a href="http://soaps.sheknows">soaps.sheknows</a>	62	5	14	0	0	0	0	0



# MODULES

Module 1:

## Feature Extraction:

-In Website Detection  
Extract

- url

-In Email Detection  
Extract

- The header of the email
- The body of the email
- Embedded urls

# MODULES

**Module 2:**

## **Word Embedding:**

- Word embedding is done using FastText-CBOW

**Module 3:**

## **Prediction Analysis:**

- Website and Email content- Random Forest Algorithm.
- Embedded URLs and website URL-URL Analysis

https://zenodo.org/records/8041387

Lenovo Support | Lenovo | McAfee | New folder | Major project | Watch Anime Online... | Watch

Other favorites

zenodo

Search records... 

Communities

Published July 2, 2023 | Version 1

## Phishing Website Dataset

Putra, I Kadek Agus Ariesta 

This dataset contains a collection of legitimate and phishing websites, along with information on the brands impersonated in the phishing attacks. The dataset includes a total of 10,395 websites, 5,244 of which are phishing websites. These websites impersonate a total of 86 different target brands.

For phishing datasets, the files can be downloaded in a zip file with a "phishing" prefix, while for non-phishing datasets, they can be downloaded in a zip file with a "not-phishing" prefix.

In addition, the dataset includes features such as screenshots, text, CSS, and HTML structure for each website, as well as domain information (WHOIS data), IP information, and SSL information. Each website is labeled as either legitimate or phishing and includes additional metadata such as the date it was discovered, the target brand being impersonated, and any other relevant information.

The dataset has been curated for research purposes and can be used to analyze the effectiveness of phishing attacks, develop and evaluate anti-phishing solutions, and identify trends and patterns in phishing attacks. It is hoped that this dataset will contribute to the advancement of research in the field of cybersecurity and help improve our understanding of phishing attacks.

### Files

category	description	identifier	name	website
	1&1 IONOS is a			

**Check your WebPage**

Find out now... 

**SAFE OR NOT?**

1K DOWNLOADS

Jul 2, 2023

all versions by using the DOI 10.5281/zenodo.8041386. This DOI represents all versions, and will always resolve to the latest one. [Read more](#).

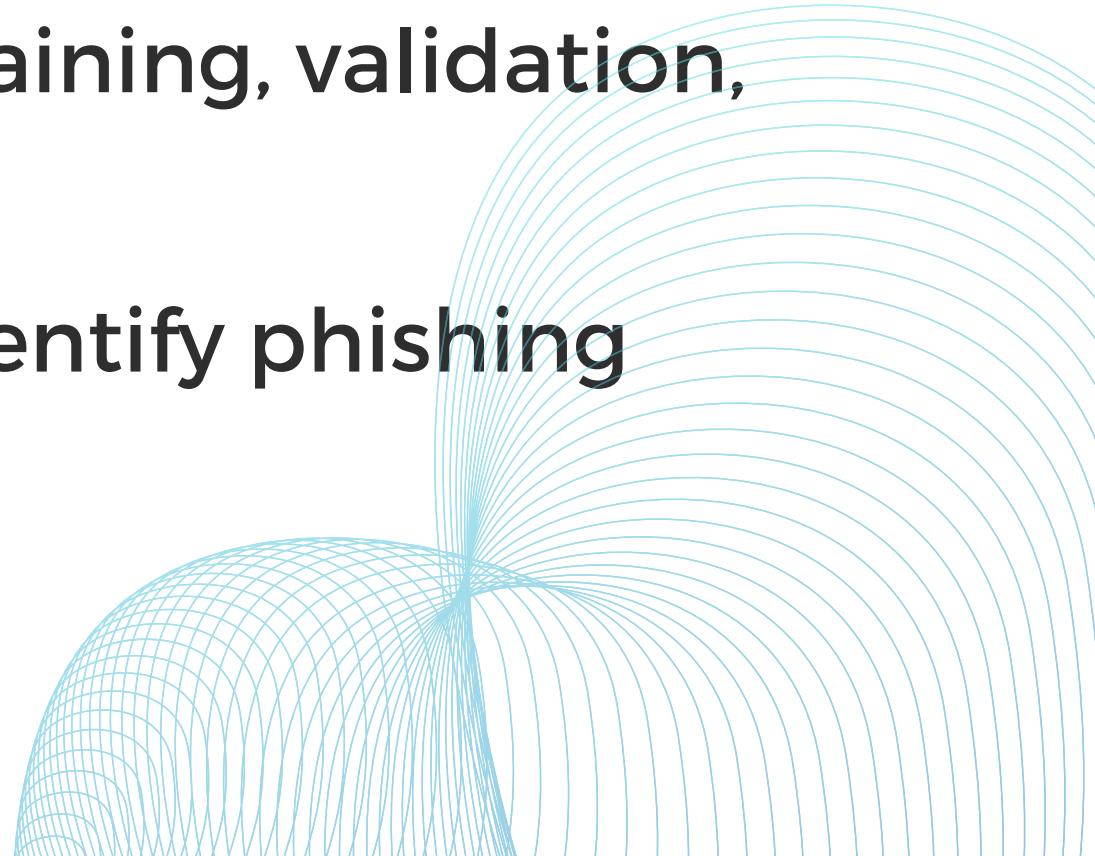
External resources

Indexed in

 OpenAIRE

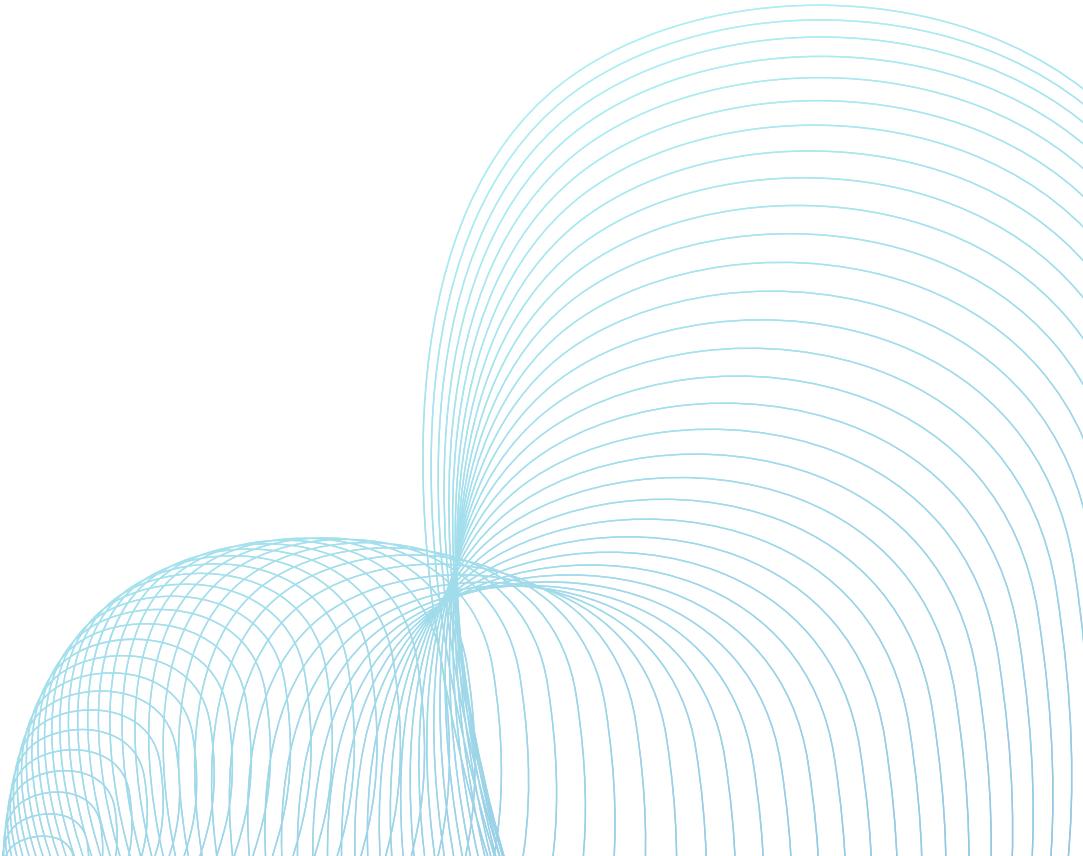
# METHODOLOGY

- **Data Collection and Preparation:** Gather a diverse dataset of both phishing and legitimate websites.
- **Feature Engineering:** Create meaningful features that can help distinguish between phishing and legitimate websites
- **Machine Learning Model Selection:** Choose appropriate machine learning algorithms or architectures and split the dataset into training, validation, and testing sets.
- **URL Analysis:** Implement URL analysis algorithms to identify phishing indicators in URLs.



# CONCLUSION

In conclusion, the proposed Phishing Detection Extension integrates advanced technology, user education, and real-time threat intelligence to provide effective protection against phishing attacks. It prioritizes user awareness, privacy, and continuous improvement to mitigate risks and enhance cybersecurity in the digital age.



# REFERENCE

- [1]Gururaj Harinahalli Lokesh & Goutham BoreGowda (2020): Phishing website detection based on effective machine learning approach, Journal of Cyber Security Technology. Available: <https://doi.org/10.1080/23742917.2020.1813396>
- [2 ]ABDUL KARIM, MOBEEN SHAHROZ, KHABIB MUSTOFA,SAMIR BRAHIM BELHAOUARI, AND S. RAMANA KUMAR JOGA,Phishing Detection System Through Hybrid Machine Learning Based on URL,Volume 11,2023, 10.1109/ACCESS.2023.3252366
- [3 ]MARIA SAMEEN, KYUNGHYUN HAN , ANDSEONG OUN HWANG,PhishHaven—An Efficient Real-Time AI Phishing URLs Detection System, VOLUME 8, 2020, 10.1109/ACCESS.2020.2991403
- [4 ]MANUEL SÁNCHEZ-PANIAGUA , EDUARDO FIDALGO FERNÁNDEZ,ENRIQUE ALEGRE ,WESAM AL-NABKI , AND VÍCTOR GONZÁLEZ-CASTRO,Phishing URL Detection: A Real-Case Scenario Through Login URLs,VOLUME 10, 2022,10.1109/ACCESS.2022.3168681

# REFERENCE

- [5]Ehsan Nowroozi , Senior Member, IEEE, Abhishek , Member, IEEE,Mohammadreza Mohammadi , Member, IEEE, and Mauro Conti ,An Adversarial Attack Analysis on Malicious Advertisement URL Detection Framework,, VOL. 20, NO. 2, JUNE 2023, doi: 10.1109/TNSM.2022.3225217
- [6]E. S. Gualberto, R. T. De Sousa, T. P. De Brito Vieira, J. P. C. L. Da Costa and C. G. Duque, "The Answer is in the Text: Multi-Stage Methods for Phishing Detection Based on Feature Engineering," in IEEE Access, vol. 8, pp. 223529-223547, 2020, doi: 10.1109/ACCESS.2020.3043396.
- [7]N. Q. Do, A. Selamat, O. Krejcar, E. Herrera-Viedma and H. Fujita, "Deep Learning for Phishing Detection: Taxonomy, Current Challenges and Future Directions," in IEEE Access, vol. 10, pp. 36429-36463, 2022, doi: 10.1109/ACCESS.2022.3151903.
- [8]K. Althobaiti, M. K. Wolters, N. Alsufyani and K. Vaneia, "Using Clustering Algorithms to Automatically Identify Phishing Campaigns," in IEEE Access, vol. 11, pp. 96502-96513, 2023, doi: 10.1109/ACCESS.2023.3310810.

# REFERENCE

- [9] Mittal, Apurv; Engels, Dr Daniel; Kommanapalli, Harsha; Sivaraman, Ravi; and Chowdhury, Taifur (2022) "Phishing Detection Using Natural Language Processing and Machine Learning," SMU Data Science Review: Vol. 6: No. 2, Article 14.  
<https://scholar.smu.edu/datasciencereview/vol6/iss2/14>
- [10] Umer Ahmed, Rashid Amin, Hamza Aldabbas, Senthikumar Mohan, Bader Alouffi, Ali Ahmadian,"Cloud-based email phishing attack using machine and deep learning algorithm"Complex Intell. Syst. 9, 3043-3070 (2023). <https://doi.org/10.1007/s40747-022-00760>
- [11] Somesh m, Alwyn R Pais,"Classification of Phishing Email Using Word Embedding and Machine Learning Techniques"2022 vol 11 issue 3 <https://doi.org/10.13052/jcsm2245-1439.1131>
- [12] Yong Fang, Cheng Zhang, Cheng Huang, Liang Liu, Yue Yang,"Phishing Email Detection Using Improved RCNN Model With Multilevel Vectors and Attention Mechanism"IEEE May 9,2019.

# **THANK YOU**

