# MALICIOUS URL DETECTION USING MACHINE LEARNING

## Team - 38

CHIDURALA ASHWINI (2103A52028)

TAKKARSU SAI PRIYA (2103A52110)

KORUKANTI MADHUSRI (2103A52090)

KATHA SRAVANI (2103A52085)

**SR UNIVERSITY**

**SCHOOL OF COMPUTER SCIENCE & ARTIFICIAL INTELLIGENCE**
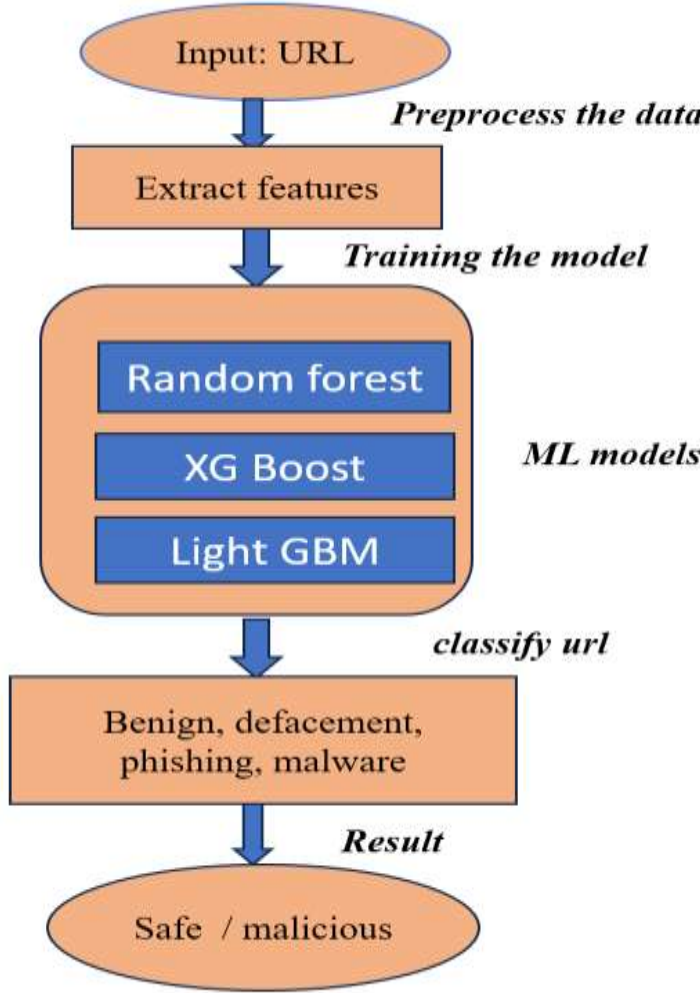
**sru**

## Abstract

Malicious URL detection is a critical aspect of cybersecurity aimed at protecting users from phishing attacks, malware, and other online threats. This project leverages machine learning techniques to build a robust system capable of identifying and classifying URLs as malicious or benign. By analyzing various features of URLs, such as length, structure, domain information, and special character usage, machine learning models learn patterns commonly associated with harmful websites. Algorithms like LightGBM

, XGBoost , Gradient Boosting, and Random Forest are employed to enhance detection accuracy and efficiency. Through data preprocessing, feature extraction, and model training, our approach ensures a high-performance system for real-time malicious URL detection. This work provides a scalable and automated solution to strengthen online security and mitigate the risks associated with malicious URLs.

## Introduction

As the internet becomes essential for personal, business, and government activities, malicious URLs pose a growing cybersecurity threat, often used in phishing and malware attacks. Detecting and blocking these harmful URLs in real-time is crucial to protect users and prevent damage.

Traditional signature-based methods struggle with new threats, but machine learning offers a promising solution. This project aims to build a system that classifies URLs as malicious or benign by analyzing their features, such as structure, domain, and special characters. Using advanced machine learning models like LightGBM, XGBoost, and Random Forest, the project focuses on enhancing detection accuracy and efficiency. The goal is to develop a robust, automated system for real-time malicious URL detection to improve online security.

## Methodology



1. Input: URL dataset (collected from Google, Kaggle)
2. Feature Extraction: URL structure, length, domain, special characters
3. Machine Learning Models:
   - XGBoost (Extreme Gradient Boosting)
   - Random Forest
   - LightGBM (Light Gradient Boosting Machine)
4. Training & Testing: Preprocessing, training on dataset, evaluating models
5. Classification: Output as Benign, Defacement, Phishing, Malware
6. Result: Safe/Malicious classification
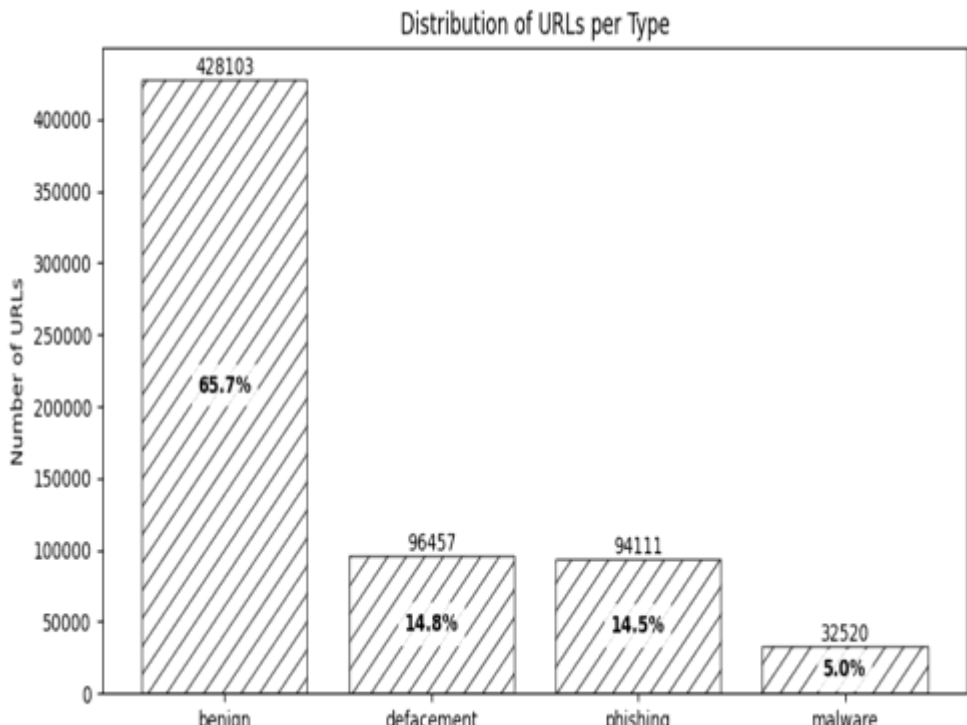
## Results
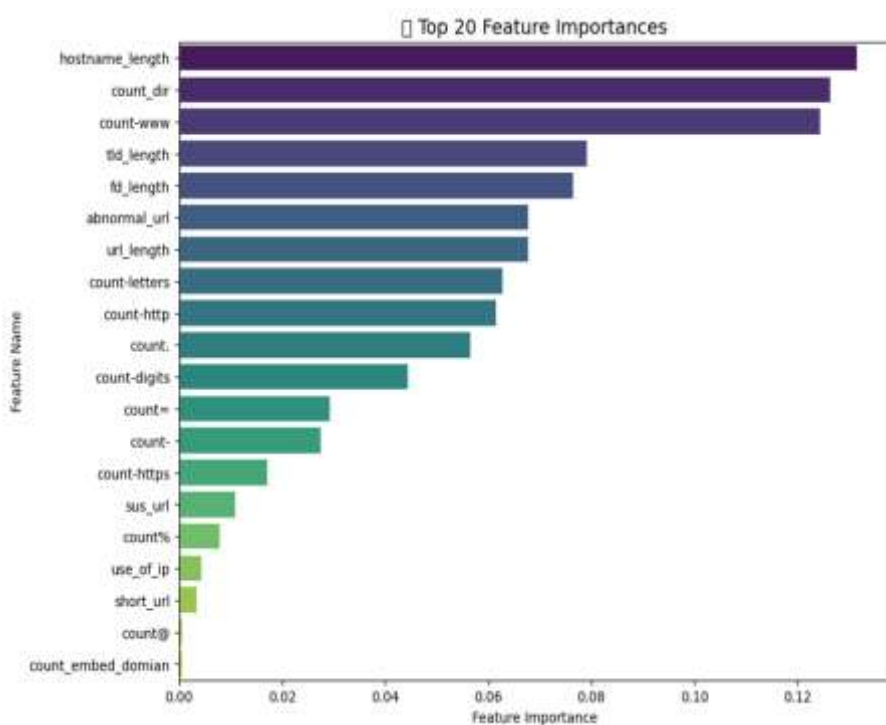


*Fig: Distribution of type of url's*



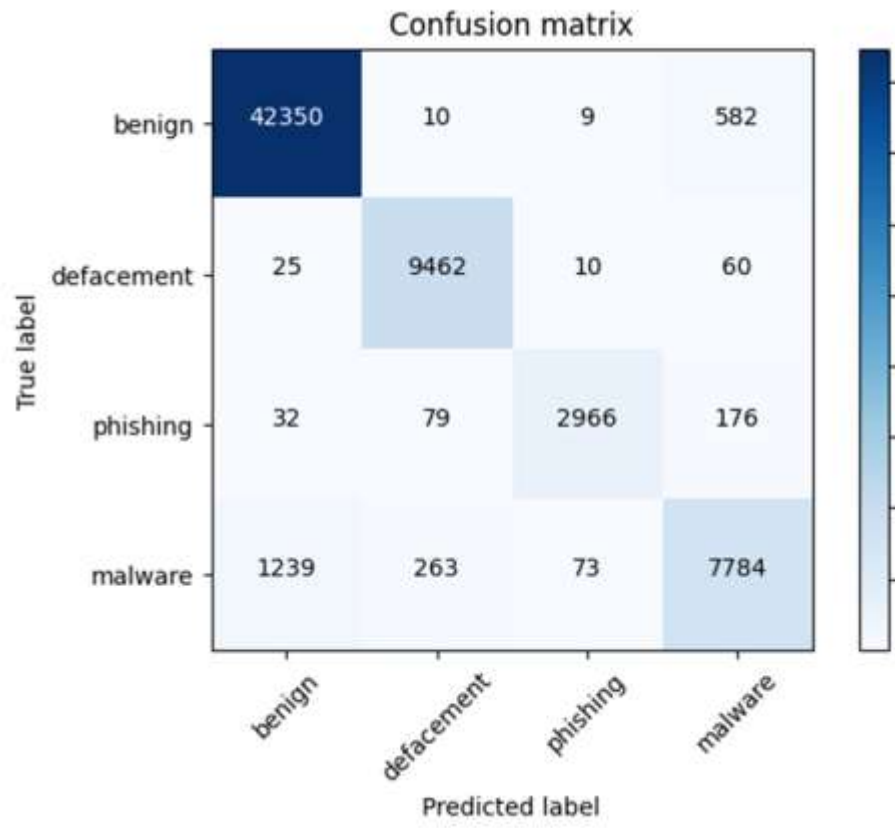*Fig: top 20 important features for url classification*
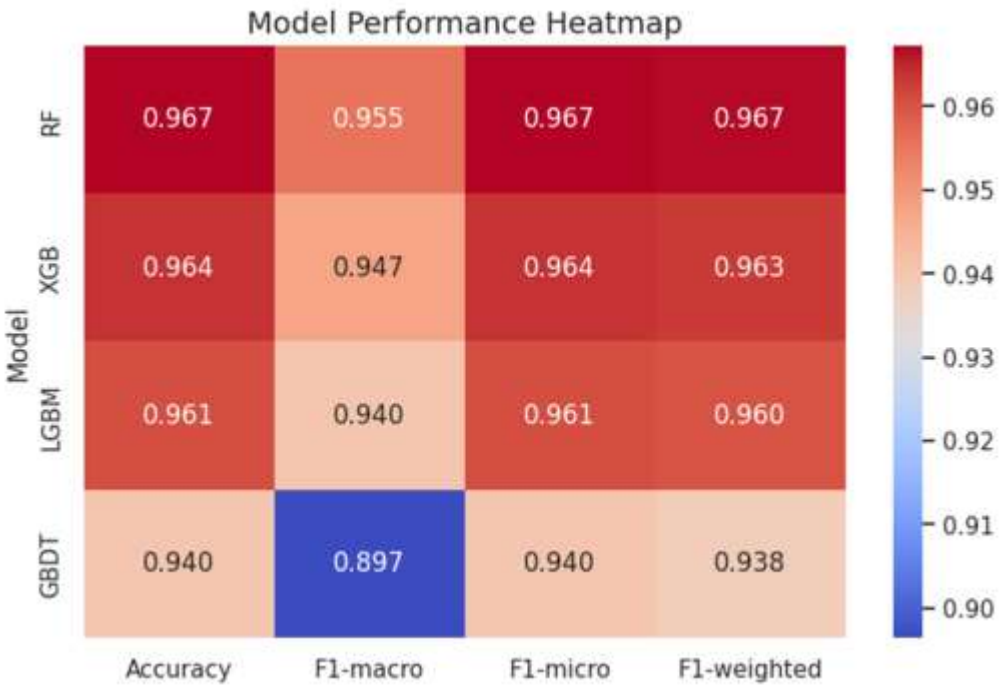


*Fig: confusion matrix*



*Fig: Heatmap for model performance*
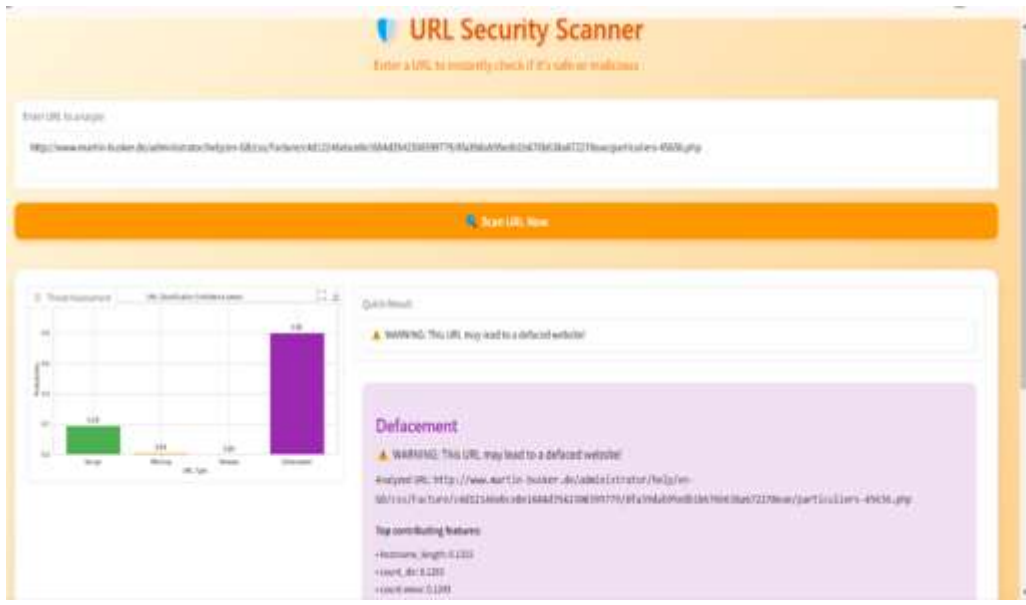


*Fig: User interface :entering url in gui*



*Fig: url detection and classification of url type*

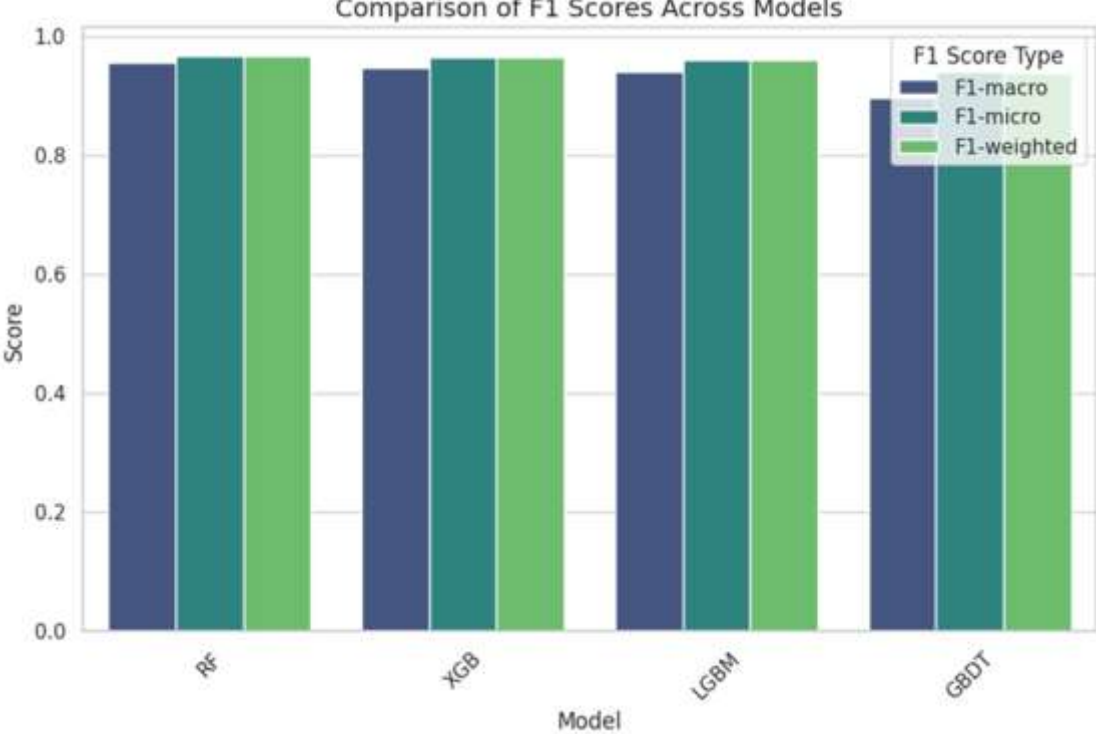## Comparision and conclusion



*Fig: comparision of F1 score across different models*

**Conclusion**

Highest accuracy achieved: 96.73% (Random Forest). The evaluation was based on accuracy, F1-macro, F1-micro, and F1-weighted scores .Random Forest achieved the highest accuracy (96.73%) and best overall performance, making it the most effective model for detecting malicious URLs. Overall, ensemble learning models (such as RF and XGB) proved to be highly effective in detecting malicious URLs, demonstrating that combining multiple decision trees leads to better generalization and robustness.

## References

1) Mankar, N. P., Sakunde, P. E., Zurange, S., Date, A., Borate, V., & Mali, Y. K. (2024, April). Comparative Evaluation of Machine Learning Models for Malicious URL Detection. In 2024 MIT Art, Design and Technology School of Computing International Conference (MITADTSoCiCon) . IEEE.
2)Vanhoenshoven, F., Nápoles, G., Falcon, R., Vanhoof, K., & Köppen, M. (2016, December). Detecting malicious URLs using machine learning techniques. In 2016 IEEE Symposium Series on Computational Intelligence (SSCI) . IEEE.

## Acknowledgements