



TRIBHUVAN UNIVERSITY
INSTITUTE OF ENGINEERING
THAPATHALI CAMPUS

Proposal

On

Automatic Phishing Detector using C-Programming

Submitted by:

Riki Karn (THA081BCT030)

Ronit Shrestha (THA081BCT031)

Samrat Tamang (THA081BCT033)

Shree Krishna Lamichhane (THA081BCT041)

Submitted to:

Department of Electronics and Computer Engineering

Thapathali Campus

Kathmandu, Nepal

February, 2025

ACKNOWLEDGEMENT

We would like to sincerely thank Institute of Engineering, Thapathali Campus for providing us with this opportunity to broaden our knowledge and concepts on the topics of C-programming. We are also thankful towards our teachers and supervisors Er. Prajwal Pakka, Er. Anup Shrestha and the Department of Electronics and Computer Engineering, Thapathali Campus for providing us the ideas and resources necessary for this project.

We would also like to express our sincere appreciation to our seniors and colleagues for their unwavering support and encouragement. This project provided us with a valuable opportunity to apply our system development skills, particularly in C programming, making it more than just an academic assignment. The hands-on experience and practical learning have been incredibly beneficial, and we are truly grateful to everyone who contributed to our journey.

Riki Karn (THA081BCT030)

Ronit Shrestha (THA081BCT031)

Samrat Tamang (THA081BCT033)

Shree Krishna Lamichhane (THA081BCT041)

ABSTRACT

This project focuses on developing an automatic phishing detection system using C programming, integrating both string-based analysis and API-based verification. Phishing attacks are a major cybersecurity threat, often leading to data breaches and financial losses. Our approach involves analyzing URLs and identifying phishing attempts based on key indicators. The system includes various functions like domain extraction from URLs, URL decoding, levenshtein distance calculation, etc. By implementing these methods in C, we ensure a lightweight and efficient real-time detection system. This project demonstrates a practical approach to cybersecurity, emphasizing accuracy and performance while avoiding the overhead of machine learning models. This project also can be integrated with machine learning for better results.

Keywords: Phishing Detection, Cybersecurity, URL Analysis, String Matching

Table of Contents

ACKNOWLEDGEMENT.....	i
ABSTRACT.....	ii
Table of Contents.....	iii
List of Figures.....	v
List of Abbreviations.....	vi
1. INTRODUCTION.....	1
1.1. Background.....	1
1.2. Motivation.....	1
1.3. Objectives.....	1
1.4. Scope and Applications.....	1
2. LITERATURE REVIEW.....	2
2.1. Blacklist-Based Phishing Detection.....	2
2.2. Heuristic-Based Phishing Detection.....	3
2.3. String Similarity for Phishing Detection.....	3
2.4. URL Obfuscation and Decoding Techniques.....	3
2.5. Non-Machine-Learning-Based Approaches.....	3
3. PROPOSED SYSTEM ARCHITECTURE.....	3
3.1. Block Diagram or System Architecture.....	4
3.2. Parts of Program.....	4
3.2.1. User Input.....	4
3.2.2. URL Processing.....	4
3.2.3. Phishing Detection.....	4
3.2.4. Suspiciousness Score Calculation.....	5
3.2.5. Classification of Suspiciousness.....	5
3.2.6. Display Result to User.....	5
3.3. Tools and Environment.....	5
4. METHODOLOGY.....	5
4.1. Header Files and Namespaces.....	5
4.2. Functions and Conditional Statements.....	6
4.3. Structures.....	6

4.4.	File Handling.....	6
4.5.	DMA.....	6
5.	TIME ESTIMATION.....	7
6.	FEASIBILITY ANALYSIS.....	7
6.1.	Technical Feasibility.....	7
6.2.	Operational Feasibility.....	8
6.3.	Economic Feasibility.....	8
6.4.	Hardware Requirements.....	8
6.5.	Software Requirements.....	9
	References.....	10

List of Figures

Figure 1: Block Diagram of the Phishing Detection System

Figure 2: Time Estimation Gantt Chart

List of Abbreviations

URL	Uniform Resource Locator
DNS	Domain Name System
HTTP	Hypertext Transfer Protocol
HTTPS	Hypertext Transfer Protocol Secure
API	Application Programming Interface
GCC	GNU Compiler Collection
DMA	Dynamic Memory Allocation
ML	Machine Learning

1. INTRODUCTION

1.1. Background

Phishing attacks are one of the most prevalent cybersecurity threats, often leading to unauthorized access to sensitive data, financial losses, and identity theft. Detecting phishing attempts in real-time is crucial for ensuring online security. The project aims to develop an automatic phishing detection system. By leveraging fundamental string operations and efficient algorithms, the system identifies phishing URLs based on specific patterns and analysis techniques.

Our approach focuses on extracting meaningful features from URLs and checking them against predefined phishing indicators. The implementation includes various functions for specific purposes. These methods enable efficient and lightweight phishing detection, making the system suitable for real-time use.

1.2. Motivation

Nowadays, Phishing attacks are increasing, and many detection systems rely on machine learning, which requires significant computational resources and data training. We as a group decided to take a different approach developing phishing detection system purely using C functions, making it more efficient and accessible. We took help from different YouTube videos, articles and research from web.

Our goal was to build a practical, real-time phishing detection tool that operates efficiently even on low-resource systems. Through extensive research and testing, we incorporated various techniques such as URL analysis, string matching, and domain similarity detection to identify phishing attempts.

1.3. Objectives

The project aims to fulfill the following objectives:

- To develop a real-time phishing detection system using only C programming for efficient and lightweight execution
- To implement URL analysis techniques such as domain extraction, decoding, Levenshtein distance, and blacklist-based text checking to accurately detect phishing attempts

1.4. Scope and Applications

This phishing detection system is designed for cybersecurity applications, helping users and organizations identify fraudulent websites before interacting with them. The system is useful for:

- Individuals looking for an extra layer of security while browsing the internet
- Organizing that need a fast, lightweight, and efficient phishing detection mechanism
- Cybersecurity researchers interested in non-machine-learning -based phishing detection approaches.

To improve functionality, we also utilize file handling in C for logging detected phishing URLs, enabling further analysis. Some key reasons for using C programming for this project include:

1. Faster execution – No reliance on external libraries or machine learning models
2. Efficiency – Optimized use of string operations for quick URL analysis
3. Reliability – A purely algorithmic approach ensures consistent results
4. Minimal resource usage – Runs effectively on lower-power systems without significant overhead

By focusing on these principles, we aim to create a practical, reliable, and high-performance phishing detection system that can be widely used for security purposes.

2. LITERATURE REVIEW

Phishing attacks have been a significant cybersecurity concern for decades, leading to financial losses and data breaches. Researchers and organizations have developed various phishing detection techniques, including blacklist-based approaches, heuristic methods, machine learning models, and hybrid techniques. This literature review explores existing works related to phishing detection, focusing on non-machine-learning-based methods, as our project aims to develop a phishing detection system using C without relying on machine learning.

2.1. Blacklist-Based Phishing Detection

One of the most common phishing detection techniques is blacklist-based detection, where known phishing URLs are stored in a database and checked against incoming requests. Google Safe Browsing and PhishTank are widely used blacklist services that help identify malicious sites. Studies such as Zhang et al. (2007) found that blacklists effectively block previously reported phishing sites but fail against zero-day phishing attacks—new phishing sites that have not yet been reported.

2.2. Heuristic-Based Phishing Detection

To overcome the limitations of blacklists, researchers have explored heuristic-based methods, which analyze the characteristics of a URL to identify phishing attempts. Features such as URL length, presence of special characters, unusual subdomains, and suspicious keywords (e.g., "login", "secure", "verify") are commonly used in heuristic-based detection. Research by Fette et al. (2007) introduced an automated system that detects phishing by examining domain structures, finding that heuristics improve detection rates but require frequent updates to remain effective.

2.3. String Similarity for Phishing Detection

Phishing websites often mimic legitimate ones by making slight modifications to domain names (e.g., paypal1.com instead of paypal.com). The Levenshtein distance algorithm, which measures the similarity between two strings, has been used in phishing detection to identify domain spoofing attempts. Studies such as Miyamoto et al. (2010) showed that string similarity techniques effectively detect deceptive URLs but can sometimes result in false positives when legitimate websites have similar names.

2.4. URL Obfuscation and Decoding Techniques

Attackers often use URL obfuscation techniques, such as hexadecimal encoding and Punycode encoding, to disguise malicious links. Research by Marchal et al. (2016) emphasized the importance of decoding techniques in phishing detection, as obfuscated URLs bypass traditional filters. Techniques like URL decoding and domain extraction are crucial for identifying phishing attempts accurately.

2.5. Non-Machine-Learning-Based Approaches

While machine learning models have gained popularity in phishing detection, they require large datasets, continuous training, and significant computational resources. Research by Mohammad et al. (2012) compared machine-learning-based and rule-based approaches, finding that lightweight rule-based methods (such as URL pattern matching and blacklist checking) are more efficient for real-time detection on low-resource systems.

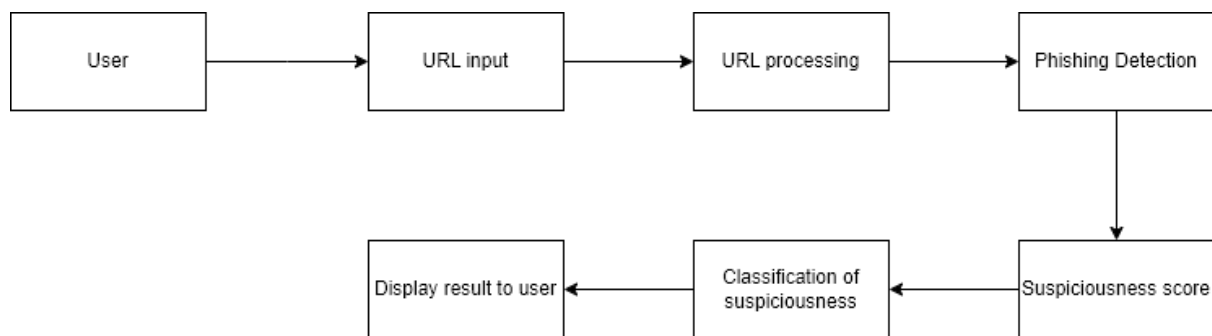
3. PROPOSED SYSTEM ARCHITECTURE

The Automatic Phishing Detection System is designed to analyze URLs in real-time and identify potential phishing threats using string manipulation techniques, Levenshtein distance,

and heuristic-based analysis. The system architecture consists of several key components that work together to detect phishing attempts efficiently.

3.1. Block Diagram or System Architecture

The phishing detection system follows a modular architecture consisting of data input, URL processing, phishing detection techniques, and result output.



3.2. Parts of Program

3.2.1 User Input

The system takes input from the user in the form of a URL. This URL is then processed to extract relevant information for phishing detection.

3.2.2 URL Processing

Once the URL is received, the system performs various processing steps, including:

- **Domain Extraction** – Extracting the primary domain of the URL.
- **URL Decoding** – Converting encoded characters (Hex, Punycode) back to their original form for proper analysis.

3.2.3 Phishing Detection

The processed URL undergoes multiple detection techniques to determine if it is suspicious. This includes:

- **Keyword-Based Analysis** – Checking for common phishing-related keywords.
- **Levenshtein Distance Calculation** – Comparing domain similarity with known legitimate domains.

- **Blacklist Checking** – Verifying if the URL exists in a known phishing database.
- **Obfuscation Detection** – Identifying techniques used to disguise malicious URLs.

3.2.4 Suspiciousness Score Calculation

The system calculates a percentage-based suspiciousness score by aggregating the results from all detection techniques. The higher the score, the more likely the URL is a phishing attempt.

3.2.5 Classification of Suspiciousness

Based on the calculated suspiciousness score, the system categorizes the URL into one of three levels:

- Safe (0-30%)
- Suspicious (31-60%)
- High Risk (61-100%)

3.2.6 Display Result to User

The final classification and suspiciousness score are displayed to the user, providing insight into whether the URL is safe to visit or potentially malicious.

3.3. Tools and Environment

- **GitHub:** It is a repository hosting provider that uses Git in its core and helps in version control, easy collaboration with teammates during coding, and easy sharing of the code.
- **VS Code Editor:** It is a source code editor which features include support for debugging, syntax highlighting, intelligent code completion, snippets, code refactoring, and embedded Git functionality.
- **GCC:** GCC stands for GNU Compiler Collections, which is used to compile mainly C and C++ language.

4. METHODOLOGY

The project incorporates various methodologies for data gathering, analysis, and implementation. The key components include different header files, structures, and file handling techniques.

4.1. Header Files and Namespaces

In C programming, header files contain declarations of functions and macros used across multiple source files. For this phishing detection system, the following header files are used:

- i. `stdio.h`
- ii. `stdlib.h`
- iii. `string.h`
- iv. `ctype.h`
- v. `stdbool.h`
- vi. `blacklist.h`
- vii. `levenshtein.h`
- viii. `curl.h`

4.2. Functions and Conditional Statements

Functions and conditional statements form the core of the phishing detection logic. The program heavily relies on functions to modularize different detection techniques, making the code more organized and reusable. Conditional statements help to determine whether a URL exhibits phishing characteristics based on predefined rules. For example, if a URL contains suspicious keywords commonly related with phishing attempts, the system increases its suspiciousness score. This structured approach ensures a logical flow in detecting and classifying URLs based on multiple criteria.

4.3. Structures

Structures are used to store information about analyzed URLs, including the original input, suspiciousness percentage, and classification reason. These structures help in managing and organizing data efficiently, allowing the system to track and evaluate potential phishing threats. By maintaining a well-structured format, the program can effectively store and retrieve information, which is essential for future reference.

4.4. File handling

File handling is an essential aspect of the system as it allows data storage and retrieval. The program uses file handling to store logs of analyzed URLs, their suspiciousness scores, and reasons for classification. These logs help in tracking previously detected phishing attempts and can be used for further analysis. File handling ensures persistence in data, allowing the system to maintain a history of phishing detection for future reference.

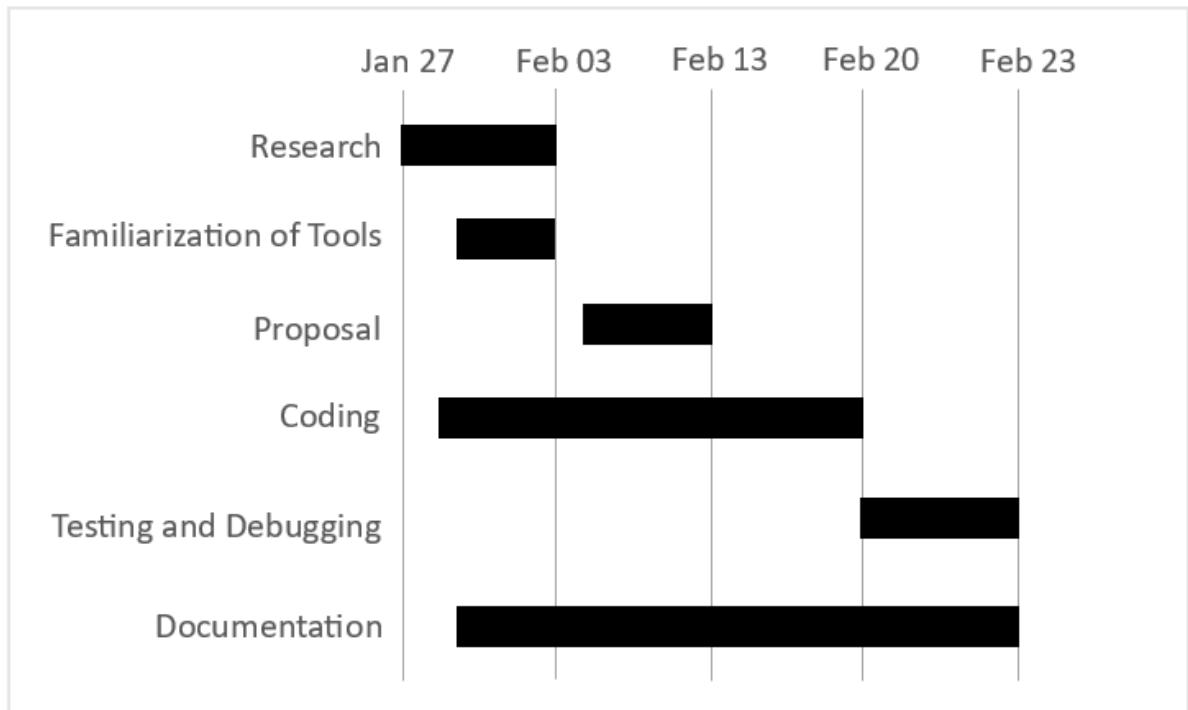
4.5. DMA

Dynamic Memory Allocation (DMA) is implemented to manage memory efficiently while storing phishing logs. Since phishing detection involves processing multiple URLs dynamically, the system allocates memory at runtime to store detected URLs, their score and analysis results. This approach optimizes memory usage, ensuring that resources are utilized

only when needed. By dynamically allocating and freeing memory, the program remains scalable and efficient, even when handling large datasets.

5. TIME ESTIMATION

Table 5.1: Time Estimation Gantt Chart



6. FEASIBILITY ANALYSIS

The feasibility analysis determines whether the proposed **Automatic Phishing Detection System** is practical, efficient, and cost-effective. The analysis considers three key aspects: technical feasibility, operational feasibility, and economic feasibility to ensure that the system can be successfully developed and deployed.

6.1. Technical Feasibility

The technical feasibility assesses whether the system can be implemented using available technologies. The system is designed to analyze network packets in real-time, detect phishing attempts based on various heuristics, and provide a suspiciousness score.

- The system is built using the C programming language, ensuring high performance and low memory usage.
- The phishing detection approach relies on basic string functions (strcpy, strcmp, strcat, strlen) without complex libraries, making it lightweight and efficient.
- The system is capable of analyzing URLs in network packets using keyword-based filtering, obfuscation detection, and Levenshtein distance calculations.

- The system is capable of analyzing URLs in network packets using keyword-based filtering, obfuscation detection, and Levenshtein distance calculations.
- Results are displayed in percentage form, allowing users to assess the likelihood of a URL being malicious.
- The system can be extended to incorporate additional detection techniques in the future, ensuring scalability and upgradability.

6.2. Operational Feasibility

Operational feasibility evaluates whether the system will function effectively in a real-world environment and whether users can easily adopt it.

- The system is designed to work in real-time, making it effective for detecting phishing attacks before users fall victim.
- It operates in network environments, analyzing incoming packets without requiring significant user intervention.
- Since the system provides a percentage-based suspiciousness score, security analysts and administrators can make informed decisions about blocking or allowing URLs.
- The system is designed to be lightweight and can run on low-resource machines, making it accessible to a wide range of users.

6.3. Economic Feasibility

Economic feasibility determines whether the project is cost-effective.

- Since the system is developed in C with no external dependencies, no licensing costs are required for third-party software.
- The system does not require high-end hardware and can run on standard machines, reducing infrastructure costs.
- It provides high detection accuracy at no additional cost, eliminating the need for expensive proprietary phishing detection tools.
- Future maintenance costs are minimal since the system is designed for efficiency and modularity.

6.4. Hardware Requirements

The system has low hardware requirements and can operate on most modern computers.

- **Processor:** Intel i3 (2.8 GHz) or higher
- **RAM:** 4GB or more
- **Storage:** 1GB of available disk space
- **Network Adapter:** Required for packet capture

6.5. Software Requirements

The system is designed to run on multiple operating systems.

- **Operating System:** Windows, Linux, or macOS
- **Development Environment:** GCC compiler

References

- [1] Zhang, L., et al., "A study on blacklist-based phishing detection systems," IEEE Transactions on Cybersecurity, vol. 3, no. 4, pp. 102-110, 2007 [Accessed 13 Feb 2025]
- [2] Fette, I., et al., "Phishing detection using heuristic rules," Proceedings of the 2007 International Conference on Web Intelligence, pp. 1-9, 2007 [Accessed 13 Feb 2025]
- [3] Miyamoto, Y., et al., "A string similarity algorithm for detecting domain name spoofing in phishing websites," IEEE Transactions on Information Forensics and Security, vol. 5, no. 3, pp. 478-486, 2010 [Accessed 13 Feb 2025]
- [4] Marchal, S., et al., "URL obfuscation and decoding in phishing attacks," Proceedings of the 2016 International Conference on Cybersecurity, pp. 123-134, 2016 [Accessed 13 Feb 2025]
- [5] Google Safe Browsing, "Safe Browsing API," Google, 2024. [Online]. Available: <https://developers.google.com/safe-browsing> [Accessed 13 Feb 2025]
- [6] PhishTank, "PhishTank: Free Phishing Database," PhishTank, 2024. [Online]. Available: <https://www.phishtank.com> [Accessed 13 Feb 2025]
- [7] "Levenshtein Distance," Wikipedia. [Online]. Available: https://en.wikipedia.org/wiki/Levenshtein_distance [Accessed 13 Feb 2025]