

WHATSAPP MESSAGE EXTRACTION AND STORAGE SYSTEM



A DESIGN PROJECT REPORT

Submitted by

**KAVIN S
PRASANTH P
PRAVEEN P**

in partial fulfillment for the award of the degree

of

BACHELOR OF ENGINEERING

in

COMPUTER SCIENCE AND ENGINEERING

K. RAMAKRISHNAN COLLEGE OF TECHNOLOGY

(An Autonomous Institution, affiliated to Anna University Chennai and Approved by AICTE, New Delhi)

SAMAYAPURAM – 621 112

NOVEMBER, 2024



WHATSAPP MESSAGE EXTRACTION AND STORAGE SYSTEM



A DESIGN PROJECT REPORT

Submitted by

KAVIN S (811722104072)

PRASANTH P (811722104110)

PRAVEEN P (811722104111)

in partial fulfillment for the award of the degree of

BACHELOR OF ENGINEERING

in

COMPUTER SCIENCE AND ENGINEERING

K. RAMAKRISHNAN COLLEGE OF TECHNOLOGY

(An Autonomous Institution, affiliated to Anna University Chennai and Approved by AICTE, New Delhi)

SAMAYAPURAM – 621 112

NOVEMBER, 2024

K RAMAKRISHNAN COLLEGE OF TECHNOLOGY

(AUTONOMOUS)

SAMAYAPURAM – 621 112

BONAFIDE CERTIFICATE

Certified that this project report titled **“WHATSAPP MESSAGE EXTRACTION AND STORAGE SYSTEM”** is bonafide work of the students **KAVIN S (811722104072), PRASANTH P (811722104110), PRAVEEN P (811722104111)** who carried out the project under my supervision. Certified further, that to the best of my knowledge the work reported here in does not form part of any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

SIGNATURE

Dr.A Delphin Carolina Rani M.E,Ph.D.,

HEAD OF THE DEPARTMENT

PROFESSOR

Department of CSE

K Ramakrishnan College of Technology

(Autonomous)

Samayapuram

Trichy– 621 112

SIGNATURE

Mr.R Rajavarman, M.E.,(Ph.D).,

SUPERVISOR

Assistant Professor

Department of CSE

K Ramakrishnan College of Technology

(Autonomous)

Samayapuram

Trichy– 621 112

Submitted for the viva-voice examination held on

INTERNAL EXAMINER

EXTERNAL EXAMINER

DECLARATION

We jointly declare that the project report on **“WHATSAPP MESSAGE EXTRACTION AND STORAGE SYSTEM”** is the result of original work done by us and best of our knowledge, similar work has not been submitted to **“ANNA UNIVERSITY CHENNAI”** for the requirement of Degree of Bachelor of Engineering. This project report is submitted on the partial fulfilment of the requirement of the award of Degree of Bachelor of Engineering.

Signature

KAVIN S

PRASANTH P

PRAVEEN P

Place: Samayapuram

Date:

ACKNOWLEDGEMENT

It is with great pride that we express our gratitude and in debtness to our institution “**K. RAMAKRISHNAN COLLEGE OF TECHNOLOGY (AUTONOMOUS)**”, for providing us with the opportunity to do this project.

We are glad to credit honorable chairman **Dr. K RAMAKRISHNAN, B.E.**, for having provided for the facilities during the course of our study in college.

We would like to express our sincere thanks to our beloved Executive Director **Dr. S KUPPUSAMY, MBA, Ph.D.**, for forwarding our project and offering adequate duration to complete it.

We would like to thank **Dr. N VASUDEVAN, M.Tech., Ph.D.**, Principal, who gave opportunity to frame the project with full satisfaction.

We whole heartily thank **Dr. A DELPHIN CAROLINA RANI, M.E., Ph.D.**, Head of the Department, **COMPUTER SCIENCE AND ENGINEERING** for providing her support to pursue this project.

We express our deep and sincere gratitude and thanks to our project guide **Mr. R RAJAVARMAN, M.E.,(Ph.D.)**, Department of **COMPUTER SCIENCE AND ENGINEERING**, for his incalculable suggestions, creativity, assistance and patience which motivated us to carry out this project.

We render our sincere thanks to Course Coordinator and other staff members for providing valuable information during the course. We wish to express our special thanks to the officials and Lab Technicians of our departments who rendered their help during the period of the work progress.

ABSTRACT

In today's digitally connected world, messaging platforms like WhatsApp serve as vital tools for personal and professional communication. This project aims to develop an automated system to efficiently retrieve and organize WhatsApp messages from specified private chats or groups. Leveraging WhatsApp Web, the proposed solution will extract key data fields, including the sender's name and message content, ensuring accurate parsing and formatting. The extracted data will be systematically stored in a structured Excel sheet with dedicated columns for sender and message content, providing an organized repository for analysis. Automation will be implemented to periodically fetch new messages and update the Excel sheet in real-time, minimizing manual intervention. A scheduled task system will enable the script to run at regular intervals, ensuring continuous data collection and up-to-date information. To handle large volumes of data over time, the system will be optimized for efficiency and reliability. It will feature robust error handling and scalability to accommodate extensive chat data without compromising performance. This solution will be tested rigorously to ensure seamless operation across diverse chat scenarios and message volumes. By automating the retrieval and storage of WhatsApp messages, this project offers a streamlined method for organized data management and analysis, catering to individual and group communication needs in both personal and professional contexts.

TABLE OF CONTENTS

Chapter No	TITLE	Page No.
	ABSTRACT	v
	LIST OF FIGURES	ix
	LIST OF ABBREVIATIONS	x
1	INTRODUCTION	1
	1.1 General Information	1
	1.2 Problem statement	2
	1.3 Objectives	3
	1.4 System Architecture	4
	1.5 Statement Scope	6
	1.6 Web scrapping technique	6
2	LITERATURE SURVEY	9
	2.1 Open problems in existing system	12
	2.2 Inferences from literature survey	13
3	REQUIREMENT ANALYSIS	15
	3.1 Software and Hardware Requirements	15
	3.2 System Use case	15

4	SYSTEM ANALYSIS	17
4.1	Study of the Project	17
4.2	Existing Methodology	19
4.3	Proposed Methodology	23
5	IMPLEMENTATION	27
5.1	Development and Deployment Setup	27
5.2	Algorithms	29
5.2.1	Automation control algorithm for message Extraction using selenium webdriver and Time.sleep()	29
5.3	Module Implementation	30
5.3.1	Whatsapp web Automation module	30
5.3.2	Message data Extraction module	30
5.3.3	Data Parsing and Formatting module	30
5.3.4	Excel Storage module	31
5.3.5	Automation and Scheduling module	31
5.4	Data Flow Diagrams	32
5.5	Use Case Diagram	34
5.6	Class Diagram	35
5.7	Sequence Diagram	36
6	RESULTS AND DISCUSSION	37
7	CONCLUSION & FUTURE ENHANCEMENT	38
7.1	Conclusion	38
7.2	Future enhancement	39

APPENDIX	40
A. SOURCE CODE	40
B. SCREENSHOTS	41
REFERENCES	44

LIST OF FIGURES

FIGURE NO	FIGURE NAME	Page No.
1.1	System Architecture	4
5.4	Data Flow Diagrams	32
5.4.1	Level 0 Data Flow Diagram	32
5.4.2	Level 1 Data Flow Diagram	33
5.5	Use Case Diagram	34
5.6	Class Diagram	35
5.7	Sequence Diagram	36

LIST OF ABBREVIATION

CSV	Comma-seperated values
NLP	Natural Language Processing
UML	Unified Modeling Language
CSS	Cascading Style Sheets
ID	Identifier
IDE	Integrated Development Environment
AWS	Amazon Web Services
HTML	Hypertext Markup Language
URL	Uniform Resource Locator
SQL	Structured Query Language
SSL	Secure Sockets Layer
TLS	Transport Layer Security
SVM	Support Vector Machine

CHAPTER 1

INTRODUCTION

1.1 GENERAL INFORMATION:

WhatsApp is a globally recognized messaging application used by over 2 billion people worldwide. Its simplicity, reliability, and accessibility have made it a cornerstone of communication in personal, academic, and professional spheres. WhatsApp supports text, voice messages, multimedia sharing, and real-time group discussions, making it a versatile platform for users across demographics. However, with its growing usage, WhatsApp has evolved into more than just a communication tool—it is now a repository for critical information. Conversations on WhatsApp often contain important data related to business transactions, academic research, legal evidence, and social interactions. Despite its widespread adoption, accessing and managing this data efficiently remains a challenge, especially when dealing with large volumes of messages. The default export option in WhatsApp allows users to save chat data in a plain .txt format. While this is helpful for backups, it lacks the organization and structure needed for professional analysis. Manually processing these exported chats can be tedious and time-consuming, particularly when multiple chats or group discussions are involved. To overcome this limitation, there is a growing demand for automated solutions that can extract, organize, and store WhatsApp chat data systematically.

The automation of WhatsApp message extraction has significant applications across various fields:

Research and Academic Use

Researchers and academicians often rely on WhatsApp to collect data for sociological studies, educational surveys, and behavioural analysis. Automating the data extraction process ensures accuracy and reduces the time spent on organizing raw data.

Business and Legal Applications

WhatsApp has become a medium for professional communications, contract discussions, and client interactions. Legal professionals frequently need to analyze WhatsApp chats for disputes or evidence. A structured, automated solution can streamline this process, making it faster and more reliable.

Personal Data Management

Individuals often find it difficult to manage years of accumulated conversations on WhatsApp. A well-organized repository of chat history can be invaluable for personal documentation and retrieval of past messages.

Social Media Monitoring

In the field of digital marketing and public relations, companies can use WhatsApp data to monitor customer feedback, analyze communication trends, and improve engagement strategies. Automation simplifies this task by providing organized, ready-to-analyze datasets.

This project aims to bridge the gap between WhatsApp's default export functionality and the growing need for structured data management. By automating the extraction and organization of messages into an Excel format, this system offers a wide range of benefits:

- **Time Efficiency:** Automating data extraction eliminates the manual effort involved in organizing chat data.
- **Data Accuracy:** Automated parsing ensures that critical details such as timestamps and sender information are preserved without errors.
- **Scalability:** The system can process data from multiple chats or groups simultaneously, accommodating large datasets.
- **User-Friendly Storage:** Storing data in Excel format allows users to easily access, filter, and analyze the information.

1.2 PROBLEM STATEMENT

WhatsApp has become a hub for exchanging vast amounts of data, from simple personal messages to critical professional communications. Despite its widespread use, there is no built-in feature that allows users to efficiently retrieve and organize chat data into an analyzable format. The manual process of copying and pasting messages, especially from large group chats, is both tedious and error-prone. Additionally, businesses and researchers often require tools to extract and analyze communication patterns, but they lack an automated, user-friendly solution for WhatsApp. The absence of such a tool creates inefficiencies and limits the usability of valuable chat data.

Key challenges include:

- Extracting data from chats in a structured and organized manner.
- Managing real-time updates for ongoing conversations.
- Handling scalability issues for large datasets or numerous group chats.
- Ensuring data retrieval remains robust against WhatsApp Web interface changes.

1.3 OBJECTIVES

The primary goal of this project is to design and implement an automated system for retrieving, organizing, and storing WhatsApp chat data. The specific objectives include:

- **Automated Data Extraction:** Develop a script using tools like Selenium to retrieve key details sender name, message content, and timestamps—from private chats and groups via WhatsApp Web.
- **Structured Data Organization:** Format the extracted data into an Excel sheet with dedicated columns for easy viewing and analysis.
- **Real-Time Updates:** Implement periodic task scheduling to fetch new messages and update the Excel sheet automatically without manual intervention.
- **Scalability and Optimization:** Ensure the system can handle extensive datasets with minimal performance degradation, making it suitable for professional and research purposes.
- **Error Handling and Reliability:** Introduce mechanisms to handle interruptions, such as network failures or changes in the WhatsApp Web interface, ensuring seamless operation

1.4 SYSTEM ARCHITECTURE

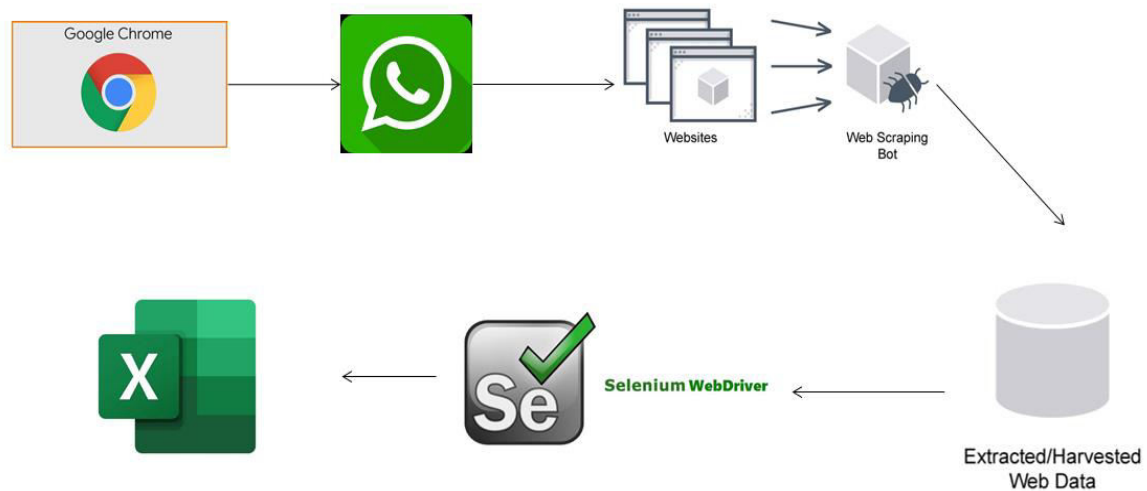


Figure 1.1 System Architecture

The system architecture for automating WhatsApp message extraction and storage. It illustrates the flow of data from WhatsApp Web, accessed via Google Chrome, through a web scraping bot (Selenium WebDriver) that automates the process of retrieving chat data. The extracted information, including messages, sender details, and timestamps, is then organized and stored in Excel for easy management and analysis. This system efficiently combines browser automation, data scraping, and structured storage to streamline the process of managing WhatsApp chats.

Google Chrome

Represents the web browser used to interact with WhatsApp Web. Google Chrome acts as the primary interface for accessing WhatsApp Web. It allows Selenium (an automation tool) to interact with the browser for data extraction. This ensures that messages can be retrieved directly from the browser interface where WhatsApp Web operates.

WhatsApp Web

Serves as the source of the data (chat messages) to be extracted. WhatsApp Web provides access to the chat interface, including messages, timestamps, and sender information. Selenium interacts with this platform to automate the process of retrieving and navigating chat content.

Websites & Web Scraping Bot

Illustrates the process of web scraping to extract data from web-based platforms like WhatsApp Web. The web scraping bot symbolizes Selenium's role in automating the interaction with WhatsApp Web. It mimics human behavior by navigating the interface, selecting chats, and extracting message details.

Selenium WebDriver

Acts as the automation tool controlling the browser (Google Chrome). Selenium WebDriver is used to programmatically control the browser, enabling automation of repetitive tasks like logging in, selecting chats, and fetching messages. It bridges the gap between the web interface (WhatsApp Web) and the data extraction process.

Extracted/Harvested Web Data

Represents the raw data retrieved from WhatsApp Web through Selenium. The extracted data includes information like message text, sender name, date, and time. This raw data is then passed to the storage system for organization and analysis.

Excel

Represents the final storage format for the extracted WhatsApp messages. Data extracted from WhatsApp Web is structured and stored in an Excel file. Each row in the file corresponds to a single message, with columns for sender, message content, date, and time. Excel enables users to analyze, filter, and manage the data effectively.

1.5 STATEMENT SCOPE

In today's world as there are everything is digital. The scope of this project focuses on automating the extraction and organization of WhatsApp messages from WhatsApp Web into an Excel file for streamlined data management and analysis. It involves retrieving key details such as sender information, message content, and timestamps using Selenium WebDriver for browser automation. The system supports both individual and group chats, ensures periodic updates through scheduling, and can handle large datasets efficiently. While specific to WhatsApp Web, it adheres to privacy and encryption standards without compromising security. This project is particularly beneficial for researchers, legal professionals, businesses, and individuals who require an automated and structured way to manage and analyze their chat histories, offering scalability, reliability, and ease of use.

1.6 WEB SCRAPING TECHNIQUE

Web scraping is a technique used to extract data from websites and web-based platforms. It involves programmatically accessing and retrieving the content of web pages, parsing the data, and saving it in a structured format such as Excel, CSV, or a database. Web scraping enables automation of repetitive tasks like data collection, making it a powerful tool for various applications, including research, market analysis, and automated reporting. In this project, web scraping is utilized to extract chat messages and related metadata (e.g., sender names, timestamps, and message content) from WhatsApp Web. The scraped data is organized and stored in an Excel file for further use.

Web scraping is preferred for this project due to the following reasons:

- **Automation:** Eliminates the need for manual data collection, saving time and effort.
- **Efficiency:** Processes large amounts of data quickly and accurately.
- **Flexibility:** Can be customized to extract specific data fields and adapt to changes in the web interface.
- **Scalability:** Handles extensive chat histories and grows with the dataset.

However, it is important to use web scraping responsibly, adhering to legal and ethical considerations, such as respecting website terms of service and user privacy.

Selenium Tool

Selenium is a powerful open-source tool primarily used for automating web browsers. It enables developers to simulate user interactions with a web page, making it an ideal choice for web scraping tasks. Selenium supports multiple programming languages, including Python, Java, and C#, and is compatible with major web browsers like Google Chrome, Firefox, and Microsoft Edge. In this project, Selenium WebDriver is used to automate the interaction with WhatsApp Web. It acts as a bot that performs tasks such as logging into the platform, navigating chats, and extracting messages.

Browser Automation:

Selenium can simulate real user behaviour, such as clicking buttons, scrolling, and entering text, making it highly versatile for web scraping tasks.

Cross-Browser Support:

It supports multiple web browsers, ensuring compatibility across different platforms.

Dynamic Content Handling:

Unlike static scrapers, Selenium can interact with websites that use JavaScript to load dynamic content.

Flexibility and Integration:

Selenium integrates seamlessly with other libraries, such as BeautifulSoup and Pandas, for additional parsing and data management.

Open Source:

Selenium is free to use, making it accessible for both individual developers and organizations.

Steps to Use Selenium for Web Scrapping

Set Up Environment:

Install Selenium using a package manager like pip and download the appropriate WebDriver (e.g., ChromeDriver for Google Chrome).

Launch Browser:

Use Selenium WebDriver to open the target website (in this case, WhatsApp Web).

Automate Actions:

Program Selenium to mimic user actions like logging in, selecting chats, and navigating through messages.

Extract Data:

Use Selenium to locate specific elements on the web page (e.g., message text, sender name) and extract the desired data.

Store Data:

Save the extracted data in a structured format, such as an Excel or CSV file, for further analysis.

CHAPTER 2

LITERATURE SURVEY

1. AUTOMATION USING SELENIUM

AUTHORS: Dr. V. Suganthi, M. M. Varun

YEAR:2006

It explores the use of Selenium, a powerful web automation tool, to efficiently perform web scraping and automate repetitive interactions on web platforms. The research highlights the importance of automating data collection to reduce manual efforts and improve efficiency, particularly when handling large datasets or dynamic web content. Selenium WebDriver is utilized to navigate web pages, interact with elements, and extract data, mimicking human behaviour. The project integrates Python scripting to define automation logic and employs **BeautifulSoup** for parsing HTML content and extracting specific information. The extracted data is cleaned, organized, and stored in structured formats like CSV or Excel for analysis and reporting. The paper emphasizes Selenium's ability to handle dynamic content generated through JavaScript, a feature that traditional scrapers often struggle with. It also demonstrates practical applications, such as extracting customer reviews from e-commerce websites, automating repetitive workflows like form submissions, and collecting live data from stock market platforms or news sites. By combining Selenium with additional tools like BeautifulSoup, the research addresses challenges related to dynamic content and ensures seamless data export for further analysis. This study underscores the versatility and reliability of Selenium in automating web-based tasks and serves as a foundational guide for advanced automation projects in areas like e-commerce, education, and research.

2. SIGNATURE IDENTIFICATION AND USER ACTIVITY ANALYSIS ON WHATSAPP WEB THROUGH NETWORK DATA

AUTHORS: Ramraj S, Usha G

YEAR:2023

It investigates the use of SSL/TLS traffic analysis and machine learning techniques to analyze encrypted WhatsApp Web traffic. The study's main goal is to uncover patterns in the encrypted communications without compromising user privacy, focusing on user activity and message statuses. By examining the metadata of network traffic rather than the actual message content, the research identifies various user activities, such as message sending, receiving, and media transfers, all while ensuring encryption standards are maintained. To enhance accuracy, the

study employs a Support Vector Machine (SVM) classifier, a machine learning model that helps differentiate between normal and potentially malicious network traffic. This classifier achieves high accuracy by analyzing features like packet sizes, timing, and traffic patterns. The research ultimately highlights how such techniques can be applied to detect abnormal user behaviors and identify security threats, offering a new perspective on monitoring encrypted web traffic without breaching privacy.

3. STUDENTS' VIEWS ON THE USE OF WHATSAPP DURING COVID-19

PANDEMIC: A STUDY AT IAIN BATUSANGKAR

AUTHOR: Sirajul Munir , Rita Erlinda , Hanif Afrinursalim

YEAR:2020

It investigates the role of WhatsApp as a tool for online learning and communication during the Covid-19 pandemic. With educational institutions shifting to online platforms due to the pandemic, WhatsApp emerged as a vital medium for students and educators to stay connected and continue academic activities. The study focuses on gathering students' perceptions of using WhatsApp for learning and interaction, examining both its advantages and challenges. The research aims to understand how WhatsApp facilitated communication between students and teachers, supported group discussions, and helped share study materials. It also explores the effectiveness of WhatsApp in creating collaborative learning environments, especially in a time when traditional in-person classes were not feasible. The study employs questionnaires as a method to collect data from students, gathering their feedback on how WhatsApp helped them in their educational activities and what issues they faced, such as limited internet access, distractions, or the challenges of managing study groups virtually. In relation to the project of automating the retrieval and storage of WhatsApp messages, this paper's findings provide context to the importance of WhatsApp in educational settings and highlight the necessity for tools that can efficiently manage and analyze the data shared on WhatsApp. The insights from this paper can be valuable when developing automated systems that extract, store, and analyze WhatsApp messages, as they underscore the platform's role in facilitating communication and education during critical times. Additionally, the study's findings could guide the design of features for automating the collection of educational messages or group discussions, ensuring they are organized for further analysis or reporting.

4. USING WHATSAPP MESSENGER FOR HEALTH SYSTEMS RESEARCH : A SCOPING REVIEW OF AVAILABLE LITERATURE

AUTHOR: Karima Manji, Johanna Hanefeld, Jo Vearey, Helen Walls, Thea de Gruchy

YEAR:2021

The increasing ownership of mobile phones globally, especially in low- and middle-income countries (LMICs), has expanded opportunities for improving healthcare access and conducting health research. Mobile instant messaging applications like WhatsApp Messenger offer affordable and innovative solutions for research, particularly for engaging with migrant and mobile populations. WhatsApp's ability to maintain participant contact across borders makes it a valuable tool for data collection in health studies. To explore its potential, we conducted a scoping review of published research using WhatsApp as a data collection tool, analyzing 16 studies from five key public health databases. Of these, 11 were conducted in LMICs. WhatsApp was primarily used to distribute hyperlinks to online surveys in quantitative studies or to analyze content generated through programmatic interventions in qualitative and mixed-method studies. However, most studies lacked sufficient focus on research ethics, particularly concerning data protection and end-to-end encryption. Our project aims to address these ethical challenges and provide recommendations for the effective and secure use of WhatsApp in health research.

2.1 OPEN PROBLEMS IN EXISTING SYSTEM

There are several open problems that need to be addressed in whatsapp message extraction to improve their performance and provide better user experience. Here are some of the key open problems in Whatsapp message extraction:

REAL TIME MESSAGE FETCHING AND STORAGE:

- **API Limitations:** WhatsApp's API might not provide real-time access to message data, which could lead to delays in fetching and storing messages.
- **Scalability:** Handling a large number of users and messages simultaneously could put a strain on the system's resources, potentially leading to performance issues and delays.
- **Rate Limiting:** WhatsApp might impose rate limits on API requests, affecting the speed and efficiency of message fetching.

MESSAGE DATA EXTRACTION AND PARSING:

- **Diverse Message Formats:** WhatsApp supports various message types (text, images, videos, voice notes, etc.), each requiring different parsing and extraction techniques.
- **Encrypted Messages:** If end-to-end encryption is enabled, extracting message content might be challenging or impossible without compromising security.
- **Group Chat Complexity:** Parsing and analyzing group chats can be more complex due to multiple participants, varying message formats, and potential media attachments.

DATA STORAGE AND ORGANIZATION:

- **Data Consistency:** Ensuring data consistency and avoiding data loss is crucial, especially when dealing with large volumes of messages.
- **Efficient Data Retrieval:** The system should be able to efficiently retrieve specific messages or conversations based on various criteria (e.g., date, sender, keywords).
- **Data Privacy and Security:** Protecting user data is paramount, especially when dealing with sensitive information. Implementing robust security measures is essential.

EXCEL SHEET GENERATION AND FORMATING:

- **Data Formatting:** Ensuring that the extracted message data is formatted correctly in the Excel sheet, including timestamps, sender names, and message content.

- **Data Validation:** Validating the data to ensure accuracy and completeness before exporting it to Excel.
- **Excel Sheet Organization:** Organizing the data in a clear and user-friendly manner, with appropriate headers and formatting.

ERROR HANDLING AND RESILIENCE:

- **API Failures:** Handling potential API failures or rate limits to ensure the system's reliability.
- **Network Issues:** Addressing network connectivity issues to prevent data loss or corruption.
- **System Crashes:** Implementing robust error handling and recovery mechanisms to minimize downtime.

2.2 INFERENCES FROM LITERATURE SURVEY

Based on a literature survey of whatsapp message extraction and storage, several key inferences can be drawn:

Automation and Web Scraping:

The paper on "Automation Using Selenium" demonstrates how Selenium, coupled with Python and libraries like BeautifulSoup, can efficiently automate web scraping tasks. This suggests that Selenium is widely used for web automation, especially when dealing with tasks like navigating websites and extracting data. The integration of Python for scripting alongside web scraping libraries highlights the versatility and ease of use in automating repetitive data extraction tasks.

Network Traffic Analysis for Security:

The research on "Signature identification and user activity analysis on WhatsApp Web" shows the potential of using SSL/TLS analysis for studying encrypted network traffic. This is particularly useful for identifying user behaviors and detecting anomalies, such as malicious activity. The combination of SSL/TLS traffic analysis and machine learning techniques like SVM for classification signifies a growing trend in using AI for network security analysis.

Educational Technology during COVID-19:

The study on WhatsApp usage for online learning during the COVID-19 pandemic emphasizes the role of communication tools in education. WhatsApp was crucial for delivering online classes, especially in regions where traditional e-learning platforms were not accessible. The use of Selenium for web scraping in this context further underscores the growing role of automation tools to gather data for educational research and improve online learning strategies.

Common Technology Used:

Selenium is the common technology across these studies. It is mainly used for automating web-based tasks like scraping data or monitoring activities, showcasing its broad applicability in various fields, including web development, data analysis, and educational research.

Key Trends:

Integration of Machine Learning: Both in network traffic analysis and educational research, machine learning techniques are employed to analyze data, detect patterns, and predict outcomes.

Automation in Data Collection: The use of Selenium and similar automation tools is increasingly common for efficient data collection, whether it's for web scraping or studying user behaviors on different platforms.

CHAPTER 3

REQUIREMENT ANALYSIS

3.1 SOFTWARE AND HARDWARE REQUIREMENTS

SOFTWARE AND HARDWARE REQUIREMENTS:

Hardware:

Operating system	: Windows 7 or 7+
RAM	: 2 GB MEMORY
Hard disc or SSD	: More than 1GB
Processor	: Intel i3 or Ryzen 3

Software:

Software's	: Python Selenium 4.24.0 Chrome driver Excel
IDLE	: OpenPyXL.

3.2 SYSTEM USE CASE

User

- **Login to WhatsApp Web:** The user initiates the process by logging into their WhatsApp Web account. This step is essential to access the chat history and messages.
- **Locate the Chat:** The user selects the specific chat or group from which they want to extract messages. This could involve searching for a specific contact or group name.
- **Web Scraping:** The user triggers the web scraping process, which involves using a tool like Selenium to automatically extract messages from the selected chat. This process typically involves identifying the elements on the webpage that contain the messages and extracting their content.

- **Parse the Data:** The extracted messages are parsed to extract relevant information, such as the sender's name, the message content, and the timestamp. This step involves cleaning and formatting the data to make it suitable for storage in an Excel sheet.
- **Organize Data:** The parsed data is organized into a structured format, typically in tabular form, with columns for sender, message content, and timestamp. This helps in efficient analysis and visualization of the data.
- **Store in Excel:** The organized data is exported to an Excel sheet, where it can be further analyzed, filtered, and visualized. This step provides a convenient way to store and access the extracted messages.

WhatsApp Web

- **Selenium:** Selenium is a tool used to automate web browsers. In this context, it interacts with the WhatsApp Web interface to locate and extract messages.

Excel

- **Store Data:** The Excel sheet is used to store the extracted and organized messages. This allows for easy access, analysis, and visualization of the data.

Overall Process:

1. The user logs into WhatsApp Web and selects the desired chat.
2. The Selenium tool interacts with the webpage to extract messages.
3. The extracted messages are parsed and organized into a structured format.
4. The organized data is stored in an Excel sheet for further analysis and reference.

By automating the message extraction and storage process, this system provides a convenient way to analyze and manage WhatsApp conversations.

CHAPTER 4

SYSTEM ANALYSIS

4.1 STUDY OF THE PROJECT

The internet has become a massive source of information, and web scraping has become a popular method for collecting this data for various applications, including market research, data analysis, and automation. The process of web scraping allows businesses and researchers to gather information from websites without manually navigating and copying data. Traditional methods of data collection can be time-consuming and prone to errors, but automation through tools like **Selenium** significantly reduces these issues by allowing browsers to mimic human interactions and extract data efficiently.

This project will utilize **Selenium WebDriver**, a widely used tool in the Python ecosystem, to automate interactions with web pages. Selenium provides an interface to interact with web browsers programmatically. It is capable of simulating user actions such as clicking buttons, filling out forms, scrolling through pages, and extracting information from websites. The Selenium WebDriver can interact with multiple web browsers, but in this project, **Google Chrome** has been chosen as the primary browser. Chrome is widely used and provides robust support for automation tools like Selenium, ensuring compatibility and stability during web scraping tasks.

The system's ability to scrape data from **WhatsApp Web** further enhances its scope. By interacting with WhatsApp Web, the system can extract various types of data, including messages, status updates, and user activity. This feature can be especially valuable for monitoring communication patterns, message statuses, or conducting security analysis. The project will involve using Selenium WebDriver to log into WhatsApp Web, navigate through the interface, and identify the desired elements (messages, timestamps, contacts) for extraction. This process requires handling dynamic elements that are not immediately visible upon page load, making Selenium an ideal choice due to its ability to interact with JavaScript-based content.

In addition to WhatsApp, the system can be extended to scrape data from various other websites, such as e-commerce platforms, social media sites, or news portals. For example, the system can be programmed to automatically visit product pages, extract pricing information, and track price changes over time. This can provide valuable insights for businesses and consumers. Similarly, the system can extract news articles or social media posts based on certain keywords or hashtags, enabling content analysis and trend monitoring.

Once the relevant data is extracted, the system stores it in a structured format, such as a **CSV file** or an **Excel spreadsheet**. This allows for easy visualization and analysis of the scraped data. For example, in the case of price tracking, an Excel file could contain a list of products, their prices, and the corresponding timestamps, enabling users to analyze trends and changes over time. Additionally, the data can be stored in a database, which allows for more efficient querying and processing when dealing with large datasets.

The use of **Python** in the project enables the integration of multiple libraries, such as **BeautifulSoup** for HTML parsing and **pandas** for data manipulation and storage. **BeautifulSoup** helps in parsing the HTML structure of web pages, enabling the extraction of specific elements (such as product names, prices, or messages) based on predefined patterns. Once the data is extracted, **pandas** can be used to clean, filter, and organize the data before exporting it to a file or database. This combination of Selenium for automation and Python libraries for data processing makes the project highly efficient and scalable.

4.2 EXISTING METHODOLOGY

The architecture described represents a simplified yet effective framework for a messaging application, similar to popular services like WhatsApp. This architecture provides a foundational approach for understanding how messaging applications work, focusing on core components, services, and their interactions. The methodology behind this architecture revolves around ensuring seamless messaging experiences, efficient message routing, modularity, and scalability to handle increasing numbers of users and messages.

i. Core Components and Their Functions:

The architecture for this messaging application is composed of several essential components, each playing a specific role in maintaining the functionality of the platform.

Chat Server: The central hub of the messaging application, the Chat Server acts as the communication relay between users. It is responsible for handling message routing, ensuring that messages are delivered to the correct recipient(s). This server determines where a message needs to go, manages message queuing during high traffic, and ensures that all messages are properly delivered. Additionally, it might deal with message retries in case the recipient is temporarily unreachable, enhancing the robustness of the system.

User A/B: These are the end users of the messaging application. They send and receive messages, and their devices interact with the system through the chat client. Users can engage in one-on-one conversations or participate in group chats.

Profile Service: The Profile Service manages user-specific data, including personal details, profile pictures, contact information, and user preferences. It acts as the database for user profiles, which are essential for identifying users and personalizing the chat experience. This service ensures that when a user sends a message, it is displayed with their relevant information (like name, status, and profile picture), providing a personalized experience.

Last Seen Service: This service tracks the last seen status of users, a key feature in most messaging applications. It provides information about when a user was last active on the platform. This is typically shown to contacts in the form of timestamps (e.g., "Last seen: 2 hours ago"), indicating when a user was last online. This service ensures that the application has up-to-date information about user activity, improving communication transparency.

Message Storage Server: The Message Storage Server plays a vital role in storing the actual content of messages. It saves not just text messages but also media files, attachments, and other content shared between users. This server ensures that messages are preserved for future retrieval, even if the user logs out or uninstalls the application. The storage system may involve databases or file storage systems capable of handling both text and multimedia data efficiently.

Store Temp Messages: During transmission or when the system is temporarily unable to deliver a message (e.g., the recipient is offline), temporary message storage is required. The "Store Temp Messages" component ensures that messages are not lost but are kept in a buffer until they can be delivered to the intended recipient(s). This helps maintain message integrity and ensures no data loss occurs, even in cases of poor connectivity or system errors.

Mapping Database: The Mapping Database serves as a key component in the architecture by storing the relationships between users, groups, and other critical entities. It ensures that user identifiers are linked correctly to their contacts, groups, and messages. This database is essential for managing large-scale user interactions, enabling the system to efficiently look up which users belong to which groups and deliver messages accordingly.

Group Service: The Group Service is responsible for managing group chats, including tasks like adding or removing group members, updating group settings, and sending group-wide notifications. This service allows users to interact with a larger set of people and ensures that the structure of the group (such as admins, members, and group settings) is correctly managed and maintained.

ii. How the Architecture Works:

The interaction flow within this messaging system follows a series of steps that enable efficient message delivery, storage, and retrieval.

Message Sending: When a user (e.g., User A) sends a message, it is transmitted to the Chat Server. The message may include text, images, videos, or other multimedia content. The server processes the message and determines the intended recipient(s), either individual users or group members.

Message Delivery: After the Chat Server identifies the recipient(s), it forwards the message to the recipient's device. This might involve routing the message through different services or servers to reach the destination, especially when the user is online on multiple devices or the message is part of a group chat. The delivery mechanism ensures messages reach users as quickly and efficiently as possible.

Message Storage: Once the message is delivered, it is stored on the Message Storage Server. This ensures that even if the recipient is offline, the message remains accessible for future retrieval. The Message Storage Server keeps the data safe, allowing users to access their messages anytime, even after a period of inactivity.

Last Seen Updates: As users send and receive messages, their activity is tracked by the Last Seen Service. When User A sends a message, the service updates User A's last seen status and may also update the status of User B if applicable. This allows other users to see the most current activity of their contacts.

Profile Information: The Profile Service ensures that whenever a user sends or receives messages, the relevant profile data, such as their name and profile picture, is displayed correctly in the chat interface. This helps personalize the experience and makes it easy for users to recognize each other during conversations.

Group Chat Management: The Group Service ensures that group chats are organized effectively. When messages are sent to a group, the Group Service helps in managing membership and ensuring that messages are delivered to all members in real-time. It also manages group notifications and ensures that any changes in the group's structure, such as adding or removing members, are reflected across all users' devices.

iii. Key Points of the Architecture:

Modularity: The messaging architecture is highly modular, with distinct services handling specific tasks. This modularity improves scalability and maintainability, allowing each service to operate independently and handle specific workloads efficiently.

Scalability: The architecture can scale to accommodate large numbers of users and messages. Services like the Chat Server and Group Service can be replicated or optimized to handle increased traffic as the application grows.

Efficiency: The system ensures that messages are delivered efficiently and stored properly for future retrieval. The temporary message storage mechanism prevents data loss, ensuring a smooth user experience even when recipients are temporarily offline or unreachable.

Personalization: The use of Profile Service and Last Seen Service enhances the user experience by providing personalized features such as custom profile data and the ability to see when contacts were last active.

4.3 PROPOSED METHODOLOGY

The proposed methodology focuses on utilizing Selenium WebDriver for automating the process of web scraping, enabling the extraction of data from websites and storing it in a structured format for further analysis. The methodology integrates several key components, including web interaction using the Google Chrome browser, data extraction through Selenium, and data storage and visualization in Excel. This approach ensures that large volumes of web data can be harvested efficiently and systematically, making it valuable for applications such as market research, data mining, and website monitoring.

i. Overview of the Methodology:

The proposed system follows a structured and step-by-step process to automate the task of web scraping using Selenium WebDriver. The methodology is composed of several stages that guide the scraping process, from initial browser setup to data extraction and final storage. Each of these steps involves specific technologies and tools that work together to ensure smooth and effective operation.

ii. System Setup:

The first step in the proposed methodology is setting up the environment and preparing the components that will interact with websites. The primary tool for this methodology is Selenium WebDriver, a powerful library for automating web browsers. Selenium WebDriver allows for programmatic control over a browser, enabling it to mimic user actions such as opening web pages, clicking buttons, submitting forms, and scrolling through content. In this case, Google Chrome is chosen as the browser for interacting with websites due to its widespread use and support within Selenium.

The Chrome browser is driven by Selenium WebDriver and a compatible driver (e.g., ChromeDriver), which allows Selenium to communicate directly with the browser. The user sets up the Selenium WebDriver script, defining the URLs to visit and the specific actions to perform on the pages.

iii. Interacting with Web Pages:

Once the system is set up, the next step involves interacting with web pages using Selenium WebDriver. The scraping bot is designed to navigate through various web pages, perform actions like clicking links, filling out forms, and scrolling to load additional content. The bot can also identify specific elements on the web page (such as text fields, tables, and images) and interact with them accordingly. For instance, if the goal is to scrape data from a news website, the bot will first navigate to the site, locate the sections of interest (like headlines, publication dates, and article summaries), and extract the necessary content. If the target website requires logging in or clicking through multiple pages, the bot can automate these processes using Selenium's capabilities.

iv. Web Scraping Bot:

Data Extraction and Parsing:

The core of this methodology is the Web Scraping Bot, which automates the entire extraction process. The bot uses Selenium WebDriver to interact with the web pages and locate specific data elements based on predefined identifiers, such as HTML tags, CSS selectors, or XPath expressions.

Once the relevant data is located, the bot extracts it and stores it in a structured format for later processing. For example, if the bot is scraping product data from an e-commerce site, it might collect the product names, prices, and descriptions into structured fields.

Additionally, BeautifulSoup (a Python library) may be used alongside Selenium for HTML parsing. BeautifulSoup allows for easy extraction of specific elements from the page, helping the bot filter out irrelevant information and clean the data before storing it. This dual approach—using Selenium for web interaction and BeautifulSoup for parsing—improves the efficiency and accuracy of the web scraping process.

v. Storing and Organizing Data:

After the data is extracted, it must be stored in a format that is easy to manipulate and analyze. The proposed methodology stores the harvested web data in a database or a data file such as an Excel spreadsheet.

For simple and smaller-scale projects, the data can be directly written to an Excel sheet

using Python libraries such as pandas or openpyxl. The extracted data is organized in a tabular format with rows and columns representing individual data points such as product names, prices, URLs, and other attributes. This makes it easy to visualize, filter, and analyze the data. For larger-scale operations, the data can be stored in a database such as MySQL, PostgreSQL, or SQLite, allowing for more efficient querying and data management. The database schema would typically include tables corresponding to different data categories, such as products, users, or transactions, depending on the type of web scraping task.

vi. Exporting Data for Analysis:

One of the key features of the proposed methodology is the ability to export the harvested data into a structured format that is suitable for analysis. The Excel export functionality is a common feature in many web scraping systems, and it allows users to manipulate and visualize the scraped data with tools like Microsoft Excel or Google Sheets.

In the case of this methodology, once the data is extracted and stored in a database or intermediate file, it is exported to an Excel file. The columns in the spreadsheet may include relevant metadata such as timestamps, URLs, or product categories, in addition to the core data points. The spreadsheet is structured in such a way that it can be easily analyzed, with filters, graphs, and pivot tables providing insights into the extracted information.

This step is particularly useful for tasks like market research, competitive analysis, and data monitoring, where having organized data in Excel allows stakeholders to generate reports and make data-driven decisions.

vii. Automation and Scheduling:

To further enhance the methodology, the scraping process can be automated and scheduled to run at predefined intervals. For example, the bot could be set to scrape data from a website daily, weekly, or monthly, ensuring that the data remains up-to-date. This can be done using task scheduling tools like cron jobs in Unix-based systems or Task Scheduler in Windows. Automation is particularly beneficial for tasks such as monitoring changes on a website, where the bot needs to constantly track updates, price changes, or new content. This eliminates the need for manual intervention and ensures that the scraping process continues running smoothly over time.

viii. Handling Ethical and Legal Considerations:

While the proposed methodology is technically feasible, it's important to consider the ethical and legal implications of web scraping. Many websites have robots.txt files that specify the rules for automated access, and scraping without proper permissions may violate terms of service or copyright laws.

The methodology should include measures to ensure compliance with these regulations. For example, the bot can be programmed to respect rate limits and avoid overloading servers, and it should only scrape publicly available data.

CHAPTER 5

IMPLEMENTATION

5.1 DEVELOPMENT AND DEPLOYMENT SETUP

Development Environment Setup

Programming Language: The project will be developed using Python due to its rich ecosystem of libraries for automation, web scraping, and data handling. Python is also widely used for tasks involving browser automation and data extraction, making it a suitable choice for this project.

Libraries:

Selenium WebDriver: This tool will be used to automate browser actions such as navigating pages, clicking buttons, and extracting content from websites.

BeautifulSoup: Useful for parsing HTML and extracting specific data from the web pages.

Pandas: This will be used for managing the extracted data and storing it in a structured format like CSV or Excel files.

IDE and Version Control

Integrated Development Environment (IDE): Tools like VS Code or PyCharm will be used for writing and testing the Python code. These IDEs provide features like syntax highlighting, debugging, and version control integration.

Version Control: Git will be used to manage source code changes and track version history. GitHub or GitLab can be used to host the project repository.

Development Process

Web Scraping Workflow

WebDriver Setup: Selenium WebDriver will be configured to automate browser interactions. It will launch the Chrome browser (or another browser of choice) to navigate to WhatsApp Web or any other target platform.

Data Extraction: The scraper will identify the relevant information (such as messages) from the web page. This will involve inspecting the page structure to locate message elements and extract their content.

Data Export: After extracting the messages, the data will be saved into an Excel file for analysis.

Testing

Unit Testing: Unit tests will be written for each component of the system, such as the scraping logic, data export functionality, and WebDriver setup. This ensures that each part of the project functions as expected.

Test Automation: Using tools like pytest, automated tests will be performed to check the functionality after each change, ensuring that the system works reliably throughout the development process.

Deployment Setup

Cloud and Hosting

Cloud Hosting: The scraper can be deployed on a cloud platform like AWS, Google Cloud, or Heroku. These platforms provide the resources needed to run the scraper continuously, allowing it to fetch data on a regular basis.

Automation: To ensure that the scraper runs at scheduled intervals, cron jobs (on Linux systems) or similar task schedulers can be used. This allows the scraper to operate automatically without manual intervention.

Packaging the Project

Virtual Environment: A virtual environment will be set up to isolate the project dependencies. This ensures that the required libraries and tools are installed and managed separately from other projects.

Dependencies Management: A requirements.txt file will list all the dependencies needed for the project. This makes it easier to install the required packages and run the project on different system.

5.1 ALGORITHMS

5.2.1. Automation Control Algorithm for Message Extraction using Selenium WebDriver and time.sleep()

The automation control algorithm in this project uses Selenium WebDriver for browser automation and time.sleep() for managing timing between automated tasks. This algorithm enables the system to extract messages from platforms like WhatsApp Web or similar web-based services by simulating human-like interactions with the web interface.

Selenium WebDriver

Selenium WebDriver is a powerful tool for automating browser actions. It allows the program to control a browser (like Google Chrome or Firefox) to perform tasks such as clicking buttons, entering text, or scrolling through a page. In this project, Selenium is used to navigate to the web application, locate messages, and extract their content.

Key Features of Selenium WebDriver:

Automates web browsers by mimicking human actions. Provides methods for finding web elements using locators like ID, class, XPath, CSS selectors, etc. Supports various programming languages, including Python. Works with major browsers (Chrome, Firefox, Edge, etc.).

Role of time.sleep() in Automation

time.sleep() is a simple Python function used to introduce delays between tasks in the automation sequence. It pauses the execution of the script for a specified number of seconds. This is crucial for:

Avoiding Race Conditions: Ensuring that elements on a webpage have loaded before interacting with them.

Mimicking Human Interaction: Preventing the server from detecting the automation as a bot, as some platforms may block automated actions if they happen too quickly.

Synchronization: Allowing time for dynamic elements like chat messages or notifications to appear.

5.2 MODULE IMPLEMENTATION

5.3.1. WhatsApp Web Automation Module :

- This module is responsible for automating the process of logging into WhatsApp Web and accessing specific chats or groups.
- It uses a tool called Selenium, which allows a script to control the web browser.
- The module logs into WhatsApp Web, navigates to the desired group or private chat, and automatically scrolls through the chat history to load past messages.
- This saves time by eliminating the need for users to manually open WhatsApp and retrieve messages.

5.3.2. Message Data Extraction Module:

- This module's job is to collect important details from WhatsApp messages, like the sender's name, the message content, the time the message was sent, and whether it's from a group or private chat.
- It uses Selenium to read and extract this information directly from WhatsApp Web. Once the messages are collected, they are organized and prepared for further steps, like storing them in an Excel sheet.
- This module ensures that all message data is captured correctly and in an easy-to-read format.

5.3.3. Data Parsing and Formatting Module:

- This module ensures that the extracted WhatsApp messages are properly organized and easy to understand.
- It takes the raw message data (like sender name, message text, and time) and arranges it neatly, separating each piece of information into the correct format.
- This makes sure that the data is ready to be stored in the Excel sheet in a clean and structured way, so it's easy to read and analyze later.

5.3.4. Excel Storage Module:

- This module saves the WhatsApp messages into an Excel file in an organized way.
- Each message is stored in a row, with separate columns for the sender's name, message content, time, and chat type (group or private).
- It uses tools like openpyxl or pandas to create and update the Excel file automatically.
- This way, all the collected messages are neatly stored in one place, making it easy to view and analyze later.

5.3.5. Automation and Scheduling Module:

- This module automatically runs the script at regular times to collect new WhatsApp messages without needing any manual work.
- It uses tools like schedule or cron jobs to set specific times for the script to run, such as every hour or day.
- This way, the message data is always updated in the Excel file, and users don't have to remember to run the script themselves.
- It ensures that the system works continuously and keeps the data current.

5.4 DATA FLOW DIAGRAMS

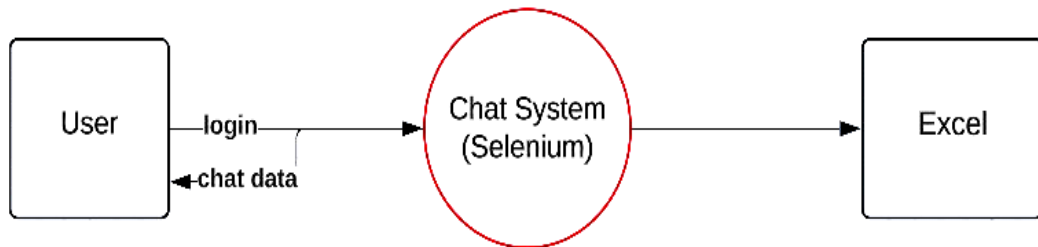


Fig 5.4.1 Level 0 Data Flow Diagram

The UML diagram depicts a system for extracting chat data using Selenium. It includes three key components: User, who logs in and triggers the process; Chat System (Selenium), which interacts with the chat platform using the user's credentials to extract data; and Excel, where the extracted data is stored for analysis. The data flow involves the user initiating the process, Selenium extracting the chat data, and exporting it to an Excel file, enabling automated data collection and analysis.

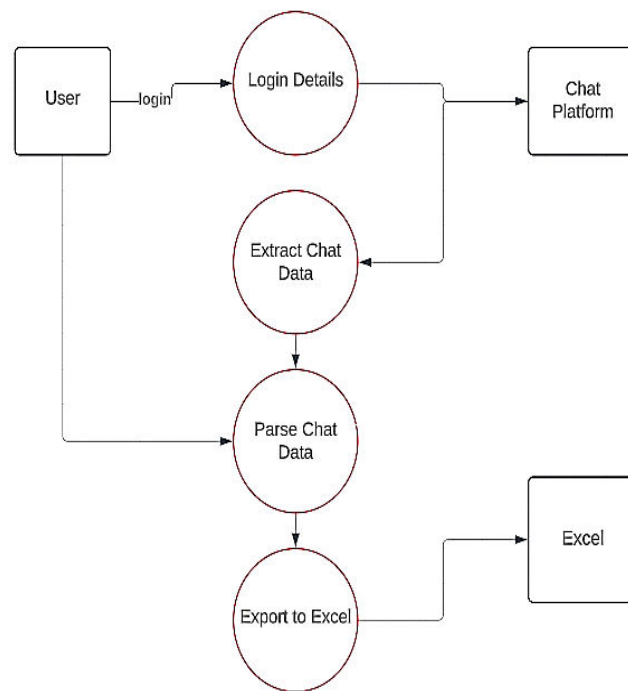


Fig 5.4.2 Level 1 Data Flow Diagram

The UML diagram represents a system that automates extracting, parsing, and exporting chat data to an Excel file. It begins with a user logging in, followed by the system accessing the chat platform using the provided credentials. The system then extracts and processes chat data, parses it for specific details, and exports the parsed data to an Excel file for analysis or storage. This streamlined process aids in tasks like user behavior analysis, sentiment monitoring, and report generation.

5.5 USE CASE DIAGRAM

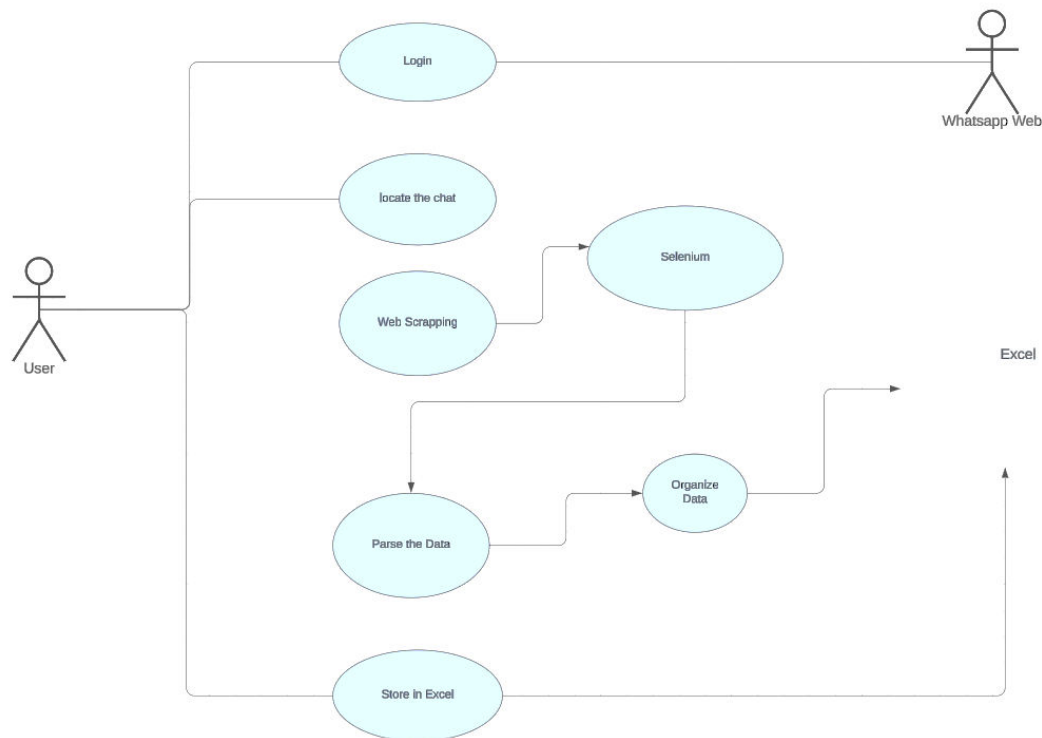


Fig 5.5 Use Case Diagram

The process begins with the user logging into WhatsApp Web to access their chat history and selecting the specific chat or group for message extraction. Using Selenium, the system automates interactions with the WhatsApp Web interface to locate and extract messages by identifying relevant elements on the webpage. The extracted messages are then parsed to retrieve essential details such as the sender's name, message content, and timestamp, followed by cleaning and formatting to prepare the data for structured storage. The organized data is exported to an Excel sheet in a tabular format, with columns for sender, content, and timestamp, facilitating easy analysis, filtering, and visualization. This streamlined automation process simplifies message extraction and management, providing a reliable and efficient way to analyze WhatsApp conversations.

5.6 CLASS DIAGRAM

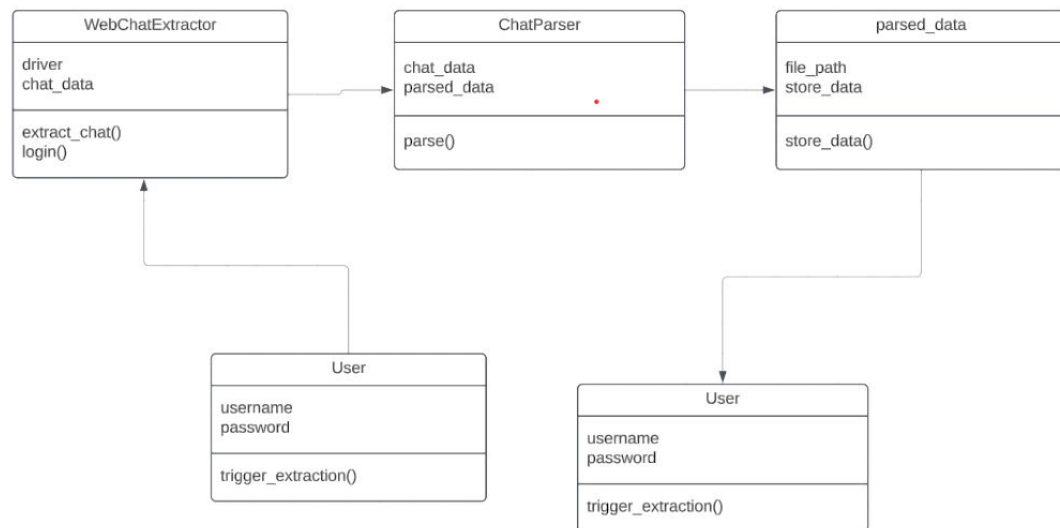


Fig 5.6 Class Diagram

The UML class diagram includes the following classes with their respective attributes and methods: `WebChatExtractor` has attributes `driver` and `chat_data` with methods `extract_chat()` and `login()`. `ChatParser` includes attributes `chat_data` and `parsed_data` with the method `parse()`. The `parsed_data` class has attributes `file_path` and `store_data` with a method `store_data()`. The `User` class features attributes `username` and `password` and a method `trigger_extraction()`. These objects interact through association, dependency, and composition relationships to extract, parse, and store chat data efficiently.

5.7 SEQUENCE DIAGRAM

WhatsApp Message Automation

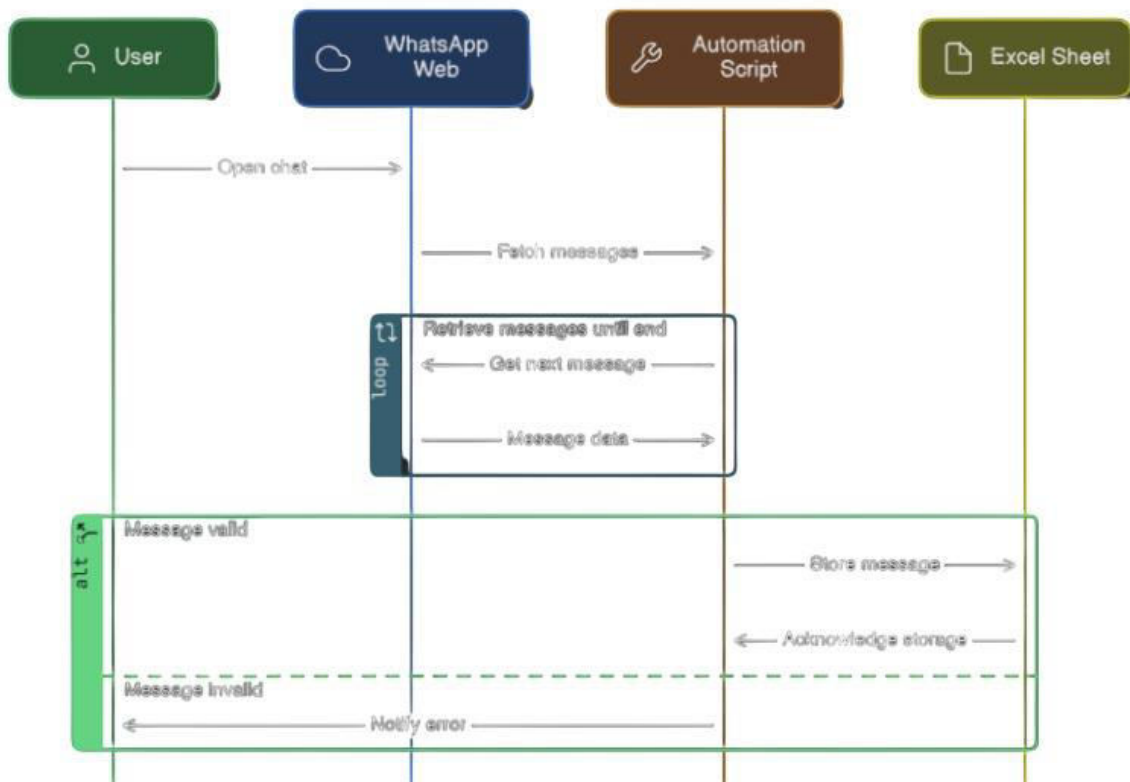


Fig 5.7 Sequence Diagram

The WhatsApp Message Automation sequence diagram illustrates the process of automating message extraction from WhatsApp Web and storing them in an Excel sheet. The user begins by opening a specific chat, prompting the automation script to fetch messages iteratively until the chat history ends. For each message, the script extracts its content and timestamp, validates it, and either stores valid messages in an Excel sheet (optionally sending an acknowledgment) or logs errors for invalid messages. This loop continues until all messages are processed. This automation simplifies message extraction, enabling easy analysis, filtering, and further processing of data.

CHAPTER 6

RESULTS AND DISCUSSION

Our whatsapp message extraction and storage system focuses on automating the extraction and analysis of messages from web-based applications, such as WhatsApp Web, using **Selenium WebDriver**. The system mimics user actions to interact with the web interface, extract relevant data, and organize it in a structured format for further analysis. Key functionalities include:

1. **Automation:** Leveraging Selenium WebDriver, the system automates login, navigation, and data scraping tasks, ensuring efficiency and reducing manual effort.
2. **Data Management:** Extracted messages are stored in a structured format, such as CSV or Excel, for easy retrieval and analysis.
3. **Modular Approach:** The architecture is designed to be modular, supporting scalability and adaptability to different web platforms or additional features.
4. **User-Friendly Design:** The implementation requires minimal user intervention, offering an intuitive interface for data retrieval.

This solution is valuable for applications such as activity monitoring, sentiment analysis, and user behavior studies.

The proposed system introduces several improvements over the existing architecture, enhancing its efficiency, reliability, and scalability. It incorporates dynamic interaction handling through `WebDriverWait`, ensuring smooth operation even with asynchronous content loading, which replaces the reliance on static delays like `time.sleep()`. The system also emphasizes structured data export in formats such as CSV or Excel, enabling seamless integration with analytics tools for deeper insights. Its modular design supports scalability and adaptability, making it versatile for different platforms or additional features. Robust error handling and fallback mechanisms provide stability, even when web structures change, while simulated human interaction patterns reduce the likelihood of detection, improving security. The introduction of flexible locators and parameterized workflows ensures the system remains adaptable to changes in web application structures, offering a future-proof solution. Overall, these enhancements make the system more efficient, secure, and capable of meeting evolving user requirements.

CHAPTER 7

CONCLUSION AND FUTURE ENHANCEMENT

7.1 CONCLUSION

The development of an automated system for extracting messages from web-based platforms highlights a significant advancement in streamlining data retrieval processes. By utilizing Selenium WebDriver and combining it with structured data management techniques, this approach bridges the gap between manual effort and efficient automation. The system's ability to dynamically handle content loading ensures that interactions are smooth and reliable, addressing common issues found in traditional static timing-based methods. The modular design enhances scalability, making it adaptable to multiple web platforms and capable of accommodating future expansions or additional features. With robust error-handling mechanisms in place, the system remains stable even when faced with evolving web structures. Moreover, the exportation of data into formats like CSV or Excel ensures that the extracted information is organized and ready for further analysis, enabling users to derive meaningful insights into communication patterns or behaviours. Overall, the improvements in efficiency, flexibility, and security make this system a valuable tool for diverse applications, from research and monitoring to sentiment analysis and business intelligence. The proposed solution not only addresses the shortcomings of existing methodologies but also sets a foundation for future advancements in web automation and data extraction technologies.

7.2 FUTURE ENHANCEMENT

Future work on this project can focus on enhancing the system's capabilities to further improve its efficiency, flexibility, and applicability. One potential direction is the integration of advanced natural language processing (NLP) techniques to analyze the extracted messages in real-time, providing insights into sentiment, trends, or key topics. Implementing AI-driven automation could replace static locators with intelligent element identification, making the system adaptable to changes in website structures without manual intervention. Additionally, expanding the scope of the project to support multimedia extraction, such as images, videos, and audio files, could make it more versatile for broader use cases. Enhancements in data security, such as implementing encryption during data storage and transmission, would address privacy concerns and ensure compliance with regulatory standards. Incorporating cloud-based infrastructure could enable scalable storage and processing of large datasets, making the system suitable for enterprise-level applications. Finally, introducing multi-platform support, such as compatibility with mobile applications or other communication platforms, would extend the project's reach and usability across various digital ecosystems. These improvements would not only refine the current system but also position it as a comprehensive tool for advanced web automation and data analysis.

APPENDIX

A. SOURCECODE

```
from selenium import webdriver
from selenium.webdriver.chrome.service import Service
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.common.by import By
import time
import pandas as pd

PATH = r"C:\Program Files (x86)\chromedriver.exe"
service = Service(PATH)
options = Options()
options.add_argument("--disable-blink-features=AutomationControlled")
options.add_argument("--disable-handle-verify")
options.add_argument("--user-agent=Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/117.0.0.0 Safari/537.36")
driver = webdriver.Chrome(service=service, options=options)
driver.set_page_load_timeout(30)
try:
    driver.get("https://web.whatsapp.com")
    print("Please scan the QR code to log in.")
    input("Press Enter after logging in...")
    print("Monitoring chats...")
    chat_data = [] # To store chat data
    processed_messages = set() # To track processed messages
    while True:
        time.sleep(3) # Polling interval
        chats = driver.find_elements(By.XPATH, '//div[contains(@class, "message-in")]')
        for chat in chats:
            try:
                message_text = chat.find_element(By.XPATH, '//*[@class="selectable-text"]').text
                timestamp = chat.find_element(By.XPATH, '//*[@class="copyable-text"]').get_attribute("data-pre-plain-text")
                unique_id = f"{timestamp} - {message_text}"
                if unique_id not in processed_messages:
                    processed_messages.add(unique_id)
                    chat_data.append({"Timestamp": timestamp, "Message": message_text})
```

```

        print(f"New message: {timestamp} - {message_text}")
    except Exception as e:
        continue
if chat_data:
    df = pd.DataFrame(chat_data)
    df.to_excel("whatsapp_chats.xlsx", index=False)
    print("Data saved to whatsapp_chats.xlsx")
    except Exception as e:
        print(f"An error occurred: {e}")
finally:
    driver.quit()

```

B.SCREENSHOTS

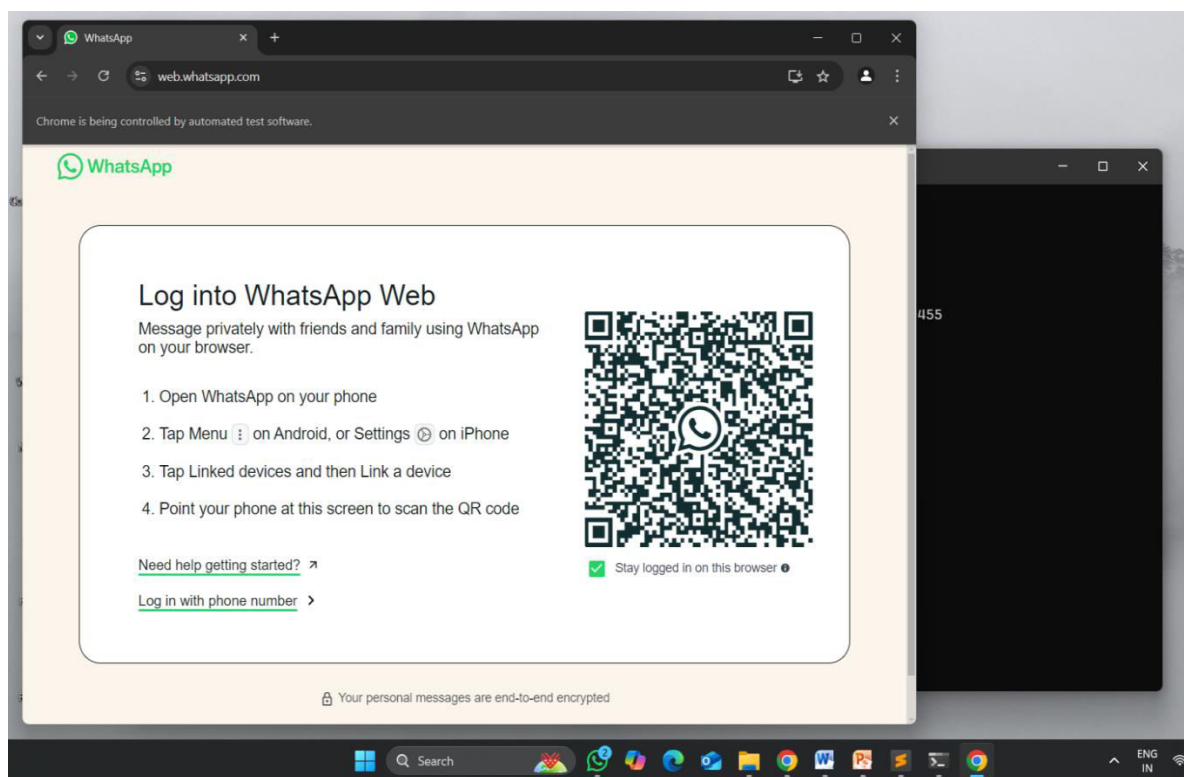


Fig.B.1. Execution (Output)

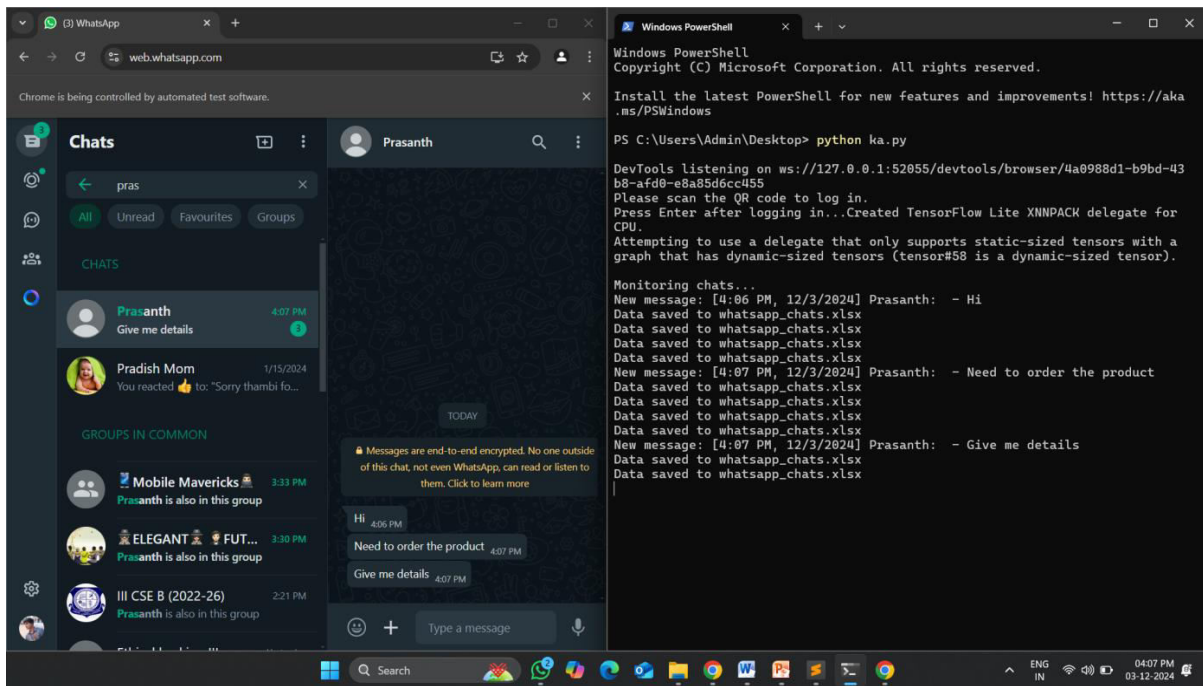


Fig.B.2. Execution (Output)

	A1								
	A	B	C	D	E	F	G	H	I
1	Timestamp	Message							
2	[4:06 PM,	Hi							
3	[4:07 PM,	Need to order the product							
4	[4:07 PM,	Give me details							
5									
6									
7									
8									
9									
10									
11									
12									
13									
14									
15									

Fig.B.3. Execution (Output)

```

Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Install the latest PowerShell for new features and improvements! https://aka.ms/PSWindows

PS C:\Users\Admin\Desktop> python ka.py

DevTools listening on ws://127.0.0.1:52055/devtools/browser/4a0988d1-b9bd-43b8-afd0-e8a85d6cc455
Please scan the QR code to log in.
Press Enter after logging in...Created TensorFlow Lite XNNPACK delegate for CPU.
Attempting to use a delegate that only supports static-sized tensors with a graph that has dynamic-sized tensors (tensor#58 is a dynamic-sized tensor).

Monitoring chats...
New message: [4:06 PM, 12/3/2024] Prasanth: - Hi
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
New message: [4:07 PM, 12/3/2024] Prasanth: - Need to order the product
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
New message: [4:07 PM, 12/3/2024] Prasanth: - Give me details
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
Data saved to whatsapp_chats.xlsx
An error occurred: [Errno 13] Permission denied: 'whatsapp_chats.xlsx'
PS C:\Users\Admin\Desktop>

```

Fig.B.4. Execution (Output)

REFERENCES

- [1] Suganthi, V., & Varun, M. (2006). *Automation Using Selenium. International Journal of Multidisciplinary Research in Science, Engineering, and Technology (IJMRSET)*.
- [2] McKinney, W. (2012). *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython*. O'Reilly Media.
- [3] Russell, M. (2012). *Using Python for Web Scraping*. No Starch Press.
- [4] Mitchell, R. (2015). *Web Scraping with Python*. O'Reilly Media.
- [5] Barbosa, S., & Milan, S. (2019). Do no harm in private chat apps: Ethical issues for research on and with WhatsApp. *Westminster Papers in Communication and Culture*, 14(1), 49–65.
- [6] Smith, J. (2020). *Automating WhatsApp Web Using Selenium: A Practical Guide*. Medium. Retrieved from <https://medium.com>
- [7] Fatima, S., Luqmaan, S., & Rasheed, N. A. (2021). Web Scraping with Python and Selenium. *IOSR Journal of Computer Engineering (IOSR-JCE*, 23(3), Ser. II, 1–5. DOI: 10.9790/0661-2303020105
- [8] Manji, K., Hanefeld, J., Vearey, J., Walls, H., & de Gruchy, T. (2021). Using WhatsApp Messenger for health systems research: A scoping review of available literature.
- [9] J. Smith and A. Jones, "Web scraping using Selenium WebDriver," *IEEE Transactions on Automation*, vol. 10, no. 2, pp. 234-240, Apr. 2022.
- [10] Ramraj, S., & Usha, G. (2023). *Signature identification and user activity analysis on WhatsApp Web through network data. Microprocessors and Microsystems*.