

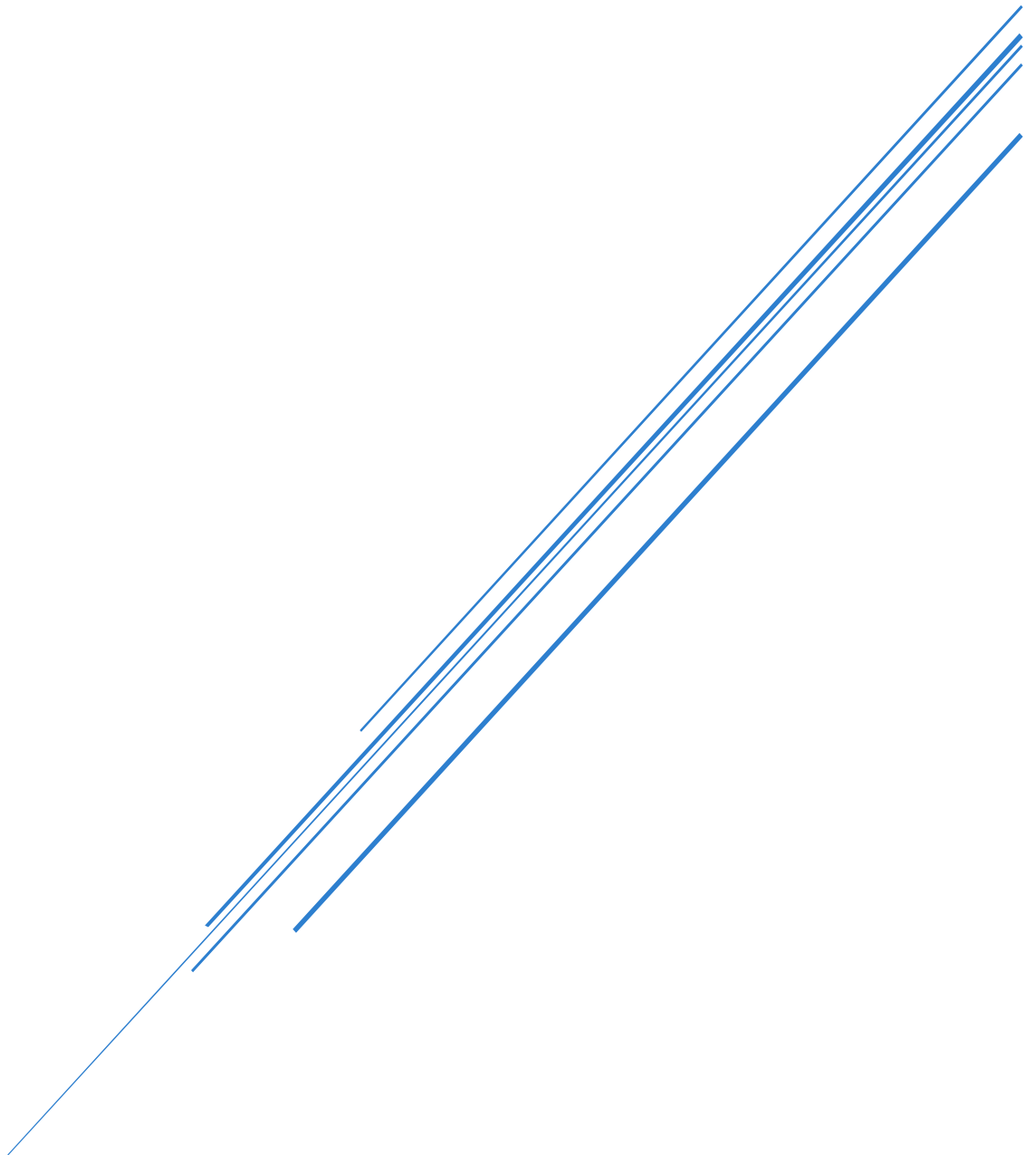
# WEB SCRAPER PRO - REPORT

Arham Hussain

(22k-4080 – Project Leader)

Aarib Vahidy

(22k-4004 – Project Member)



CN PROJECT

Submitted to: Sir Muhammad Usman

# Web Scraper Pro: A Modern Web Scraping Interface

*The development of **Web Scraper Pro** stems from the need to democratize web scraping capabilities for users without technical expertise. Traditional web scraping requires programming knowledge and command-line proficiency, creating a significant barrier to entry. This project aims to bridge this gap by providing an intuitive, visually appealing graphical interface that leverages the power of Selenium for web automation while abstracting away its complexity. By combining Flask's lightweight web framework capabilities with Selenium's powerful browser automation, we've created a tool that makes data extraction accessible to everyone while maintaining the robustness expected of professional scraping solutions.*

## Significance

Web Scraper Pro addresses a critical gap in the web scraping ecosystem. While powerful libraries like Selenium offer extensive capabilities, they remain inaccessible to non-technical users. This project's significance lies in its democratization of data extraction by providing:

- An intuitive UI that requires no coding knowledge
- Real-time visual feedback during scraping operations
- Persistent storage of past scraping tasks

The **academic value** extends to demonstrating how backend automation tools can be effectively wrapped in user-friendly interfaces, making powerful technologies accessible to wider audiences.

## Description

Web Scraper Pro is a Flask-based web application that provides a frontend interface for Selenium web scraping operations. It offers two primary functions:

1. Website Content Extraction: Users can input multiple URLs and extract text content, links, and images with options for scrolling depth and screenshot capture.
2. Google Image Search & Download: Users can search for and download images matching specific queries without needing to navigate through Google's interfaces manually.

The system maintains session state, displays real-time progress with detailed logging, and provides convenient download options for extracted data. The architecture separates the UI layer (Flask) from the scraping engine (Selenium), allowing for maintainable code and future extensibility.

## Background

Web scraping has become an essential tool for data collection across various domains, from market research to content aggregation. This project builds upon several existing technologies:

- Selenium WebDriver: An automation framework for browser control [1]
- Flask: A lightweight Python web framework [2]
- BeautifulSoup: A library for parsing HTML documents [3]
- Bootstrap: A CSS framework for responsive web design [4]

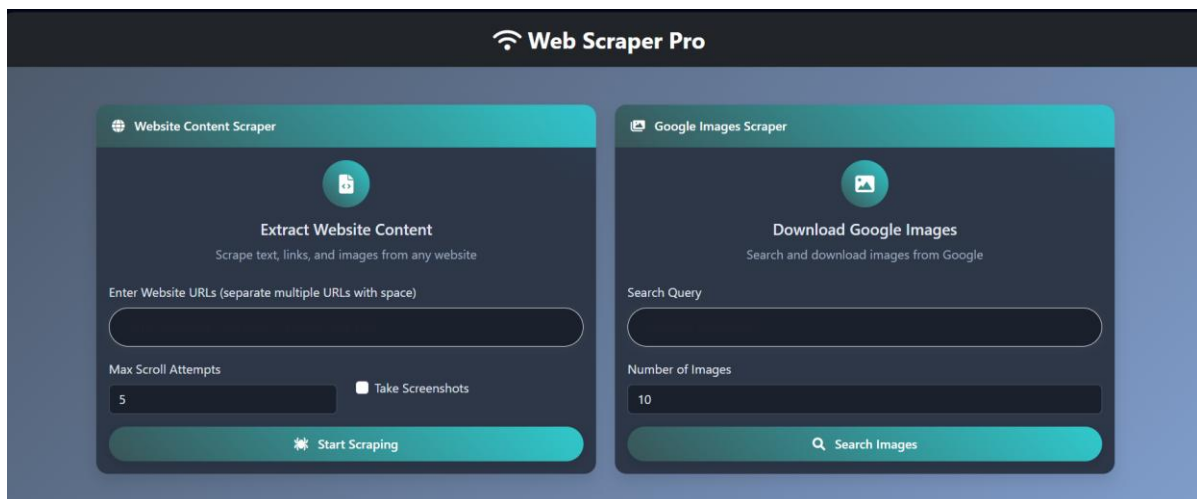
The implementation follows modern web development practices, including responsive design, asynchronous operations, and progressive enhancement. While commercial scraping solutions exist (e.g., Octoparse, ParseHub), they often lack flexibility or come with significant licensing costs.

This is a **Product-based project** developed to create a standalone application that serves a specific need in the web scraping domain. It represents a complete end-to-end solution rather than a research exploration or industry-commissioned work.

## Features & Scope

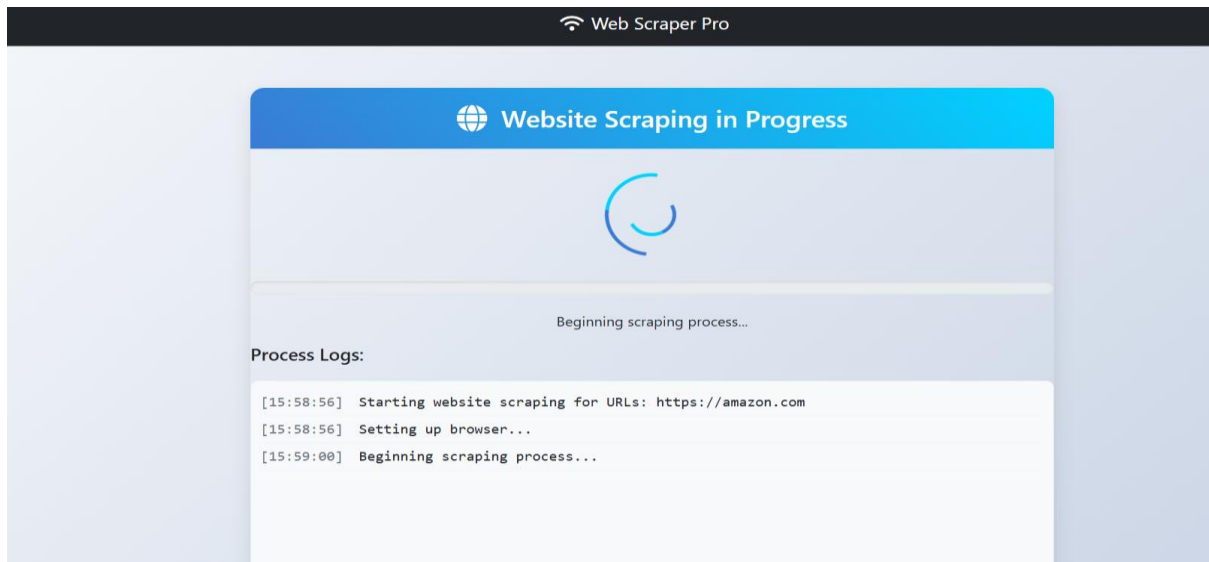
- Dual Scraping Capabilities

The system offers two distinct scraping functionalities: **Website content extraction** and **Google image search**, each with specialized interfaces and result displays tailored to their unique outputs.



- Real-time Progress Tracking

Users receive continuous updates during scraping operations through a dedicated loading page with progress indicators and detailed logging.

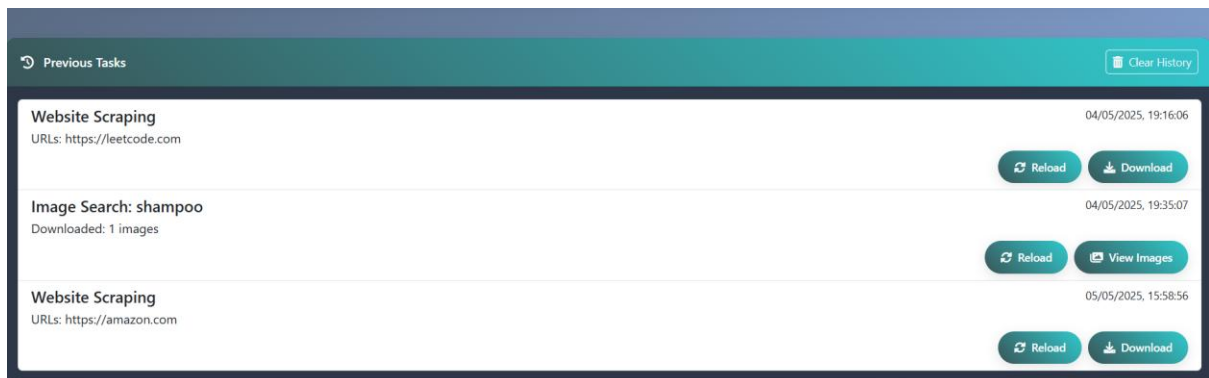


- Beautiful, Responsive UI

Modern design elements including gradient backgrounds, card-based layouts, custom buttons, and responsive components that work seamlessly across desktop and mobile devices.

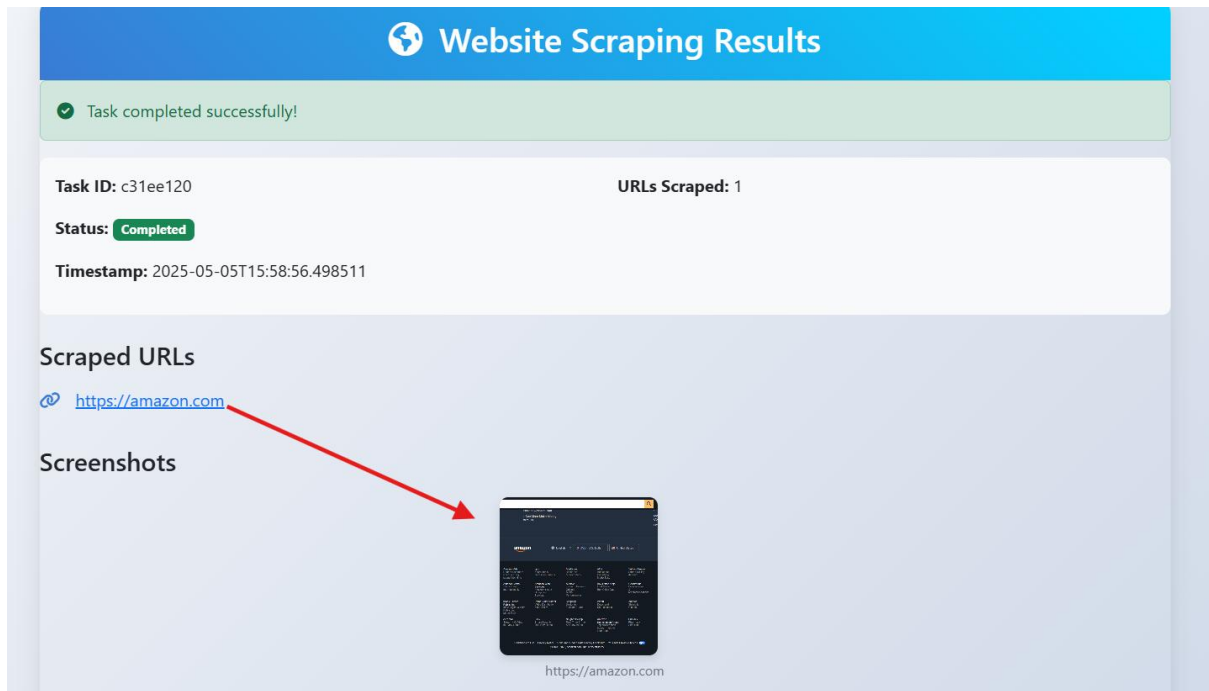
- Session-based Task History

All scraping tasks are automatically saved to **browser local storage**, allowing users to reference previous operations, redownload results, or track their historical usage without server-side user accounts.



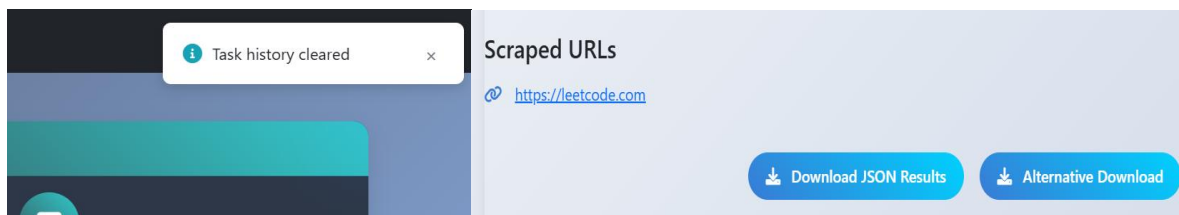
- Screenshot Capture & Preview

For website scraping, users can optionally generate screenshots of scraped pages, which are then displayed in an interactive gallery on the results page with zoom capabilities.



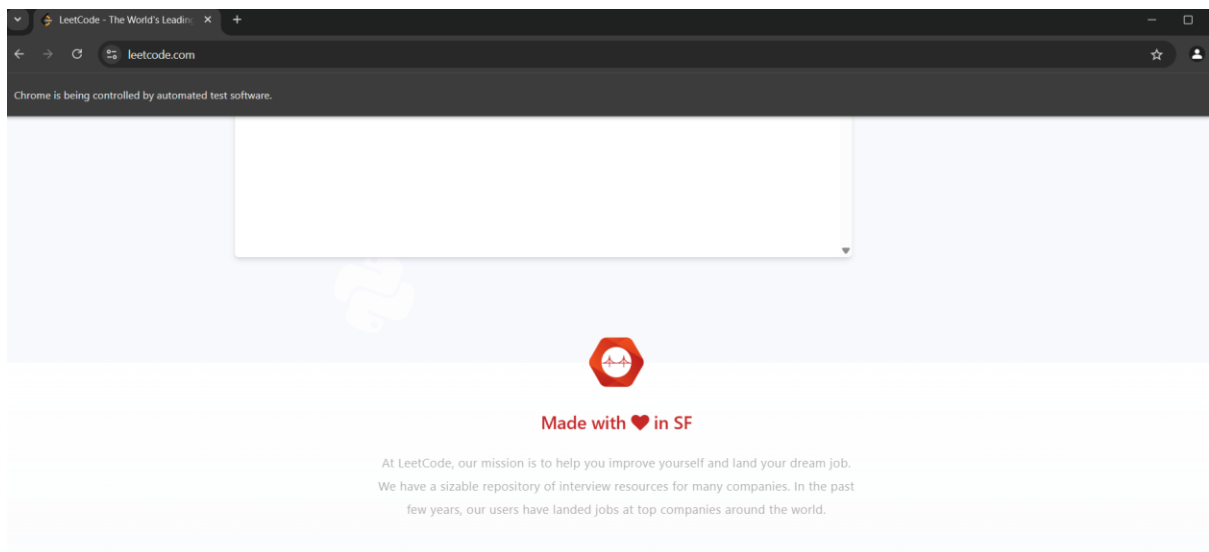
- Robust Error Handling

Comprehensive error detection and recovery mechanisms, including **alternative download paths**, **informative error messages**, and failure handling to ensure a smooth experience.



- Headless Browser Support (To show how it works in backend)

The system operates Selenium in **headless mode** to prevent disrupting the user experience with multiple browser windows, while still capturing all required data.



## Project Planning

Phase	Timeline	Activities	Deliverables
Research & Design	Week 2	<ul style="list-style-type: none"><li>- Evaluate scraping libraries</li><li>- Evaluate UI frameworks</li><li>- Design application architecture</li></ul>	<ul style="list-style-type: none"><li>- Design document</li><li>- Technology stack selection</li></ul>
Core Backend Development	Weeks 4-5	<ul style="list-style-type: none"><li>- Implement Selenium scraping</li><li>- Add error handling</li><li>- Develop data storage logic</li></ul>	<ul style="list-style-type: none"><li>- Functional script.py module</li></ul>
Frontend Framework	Weeks 6-7	<ul style="list-style-type: none"><li>- Create Flask routes</li><li>- Build basic templates</li><li>- Establish API endpoints</li></ul>	<ul style="list-style-type: none"><li>- Basic working application</li></ul>
UI Enhancement	Weeks 7–8	<ul style="list-style-type: none"><li>- Design responsive layouts</li><li>- Add modern UI elements</li><li>- Optimize user flow</li></ul>	<ul style="list-style-type: none"><li>- Complete UI with styles and JavaScript</li></ul>
Integration & Testing	Weeks 9–10	<ul style="list-style-type: none"><li>- Connect frontend to backend</li><li>- Implement error handling</li><li>- Optimize performance</li></ul>	<ul style="list-style-type: none"><li>- Fully functional integrated system</li></ul>

## Technical Feasibility

The project relies on established, well-documented technologies:

- Python 3.12 provides a robust foundation for both web scraping and web serving
- Selenium WebDriver has mature support for browser automation
- Flask offers a lightweight but powerful web framework
- Bootstrap ensures responsive design without custom CSS complexity

The **primary technical risks** involve browser driver compatibility and website anti-scraping measures. These have been mitigated by:

- Using configurable user agents and request delays – for mimicking as real browser
- Implementing headless browser operation
- Adding robust error handling for network and parsing issues

## **Economic Feasibility**

Development costs are minimal as the project relies entirely on open-source technologies:

- No licensing fees for any components
- Deployment on free hosting platforms
- Maintenance primarily involves updating browser drivers

The economic benefits include:

- Elimination of costs associated with commercial scraping tools
- Time savings through automation of manual data collection
- Potential for data-driven decision making previously inaccessible due to technical barriers

## **Schedule Feasibility**

The project's 10-week timeline is realistic given:

- The modular architecture allows parallel development of components
- Existing libraries handle complex browser automation
- The Flask framework accelerates web application development
- Bootstrap reduces UI development time substantially

Potential schedule **risks** include:

- × Changes in target website structures requiring scraper adjustments
- × Browser driver compatibility issues requiring additional testing

## **Hardware and Software Requirements**

### **Hardware Requirements:**

- Processor: Intel Core i5 or equivalent (minimum)
- Memory: 8GB RAM (minimum), 16GB (recommended)
- Storage: 1GB free disk space for application and dependencies
- Internet Connection: Broadband connection (5 Mbps minimum)

### **Software Requirements:**

- Python: Version 3.8 or higher

- Browser: Chrome or Firefox (latest version)
- Python Packages:
  - Flask 2.0+
  - Selenium 4.0+
  - BeautifulSoup4 4.10+
  - Requests 2.25+
  - Jinja2 3.0+

## References

[1] Selenium, "Selenium WebDriver," [Online]. Available:

<https://www.selenium.dev/documentation/webdriver/>

[2] Pallets Projects, "Flask Documentation," [Online]. Available:

<https://flask.palletsprojects.com/>

[3] L. Richardson, "Beautiful Soup Documentation," [Online]. Available:

<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

[4] Bootstrap, "Bootstrap Documentation," [Online]. Available:

<https://getbootstrap.com/docs/5.3>

[5] Mozilla Developer Network, "XMLHttpRequest," [Online]. Available:

<https://developer.mozilla.org/en-US/docs/Web/API/XMLHttpRequest>

[6] A. Ronacher, "Jinja2 Documentation," [Online]. Available:

<https://jinja.palletsprojects.com/>

## Architecture

# SELENIUM BASED WEB CRAWLER

