



INTERNATIONAL JOURNAL OF CREATIVE RESEARCH THOUGHTS (IJCRT)

An International Open Access, Peer-reviewed, Refereed Journal

Data Crawler

¹Aryan Patil, ²Mohit Yadav, ³Sanchit Gorlewar, ⁴Bhavesh Chaudhari, ⁵Prof. Sonali Dhaweale

¹²³⁴Final Year B.E. Student, ⁵Professor,
¹²³⁴⁵Department of Computer Engineering,
¹²³⁴⁵Terna Engineering College,
Navi Mumbai, India

Abstract: In today's data-driven world, the internet is a vast repository of valuable information. Extracting and analyzing this data can offer unparalleled insights, opportunities, and a competitive edge. The Data Crawler, also known as a web crawler or spider, is a pivotal technology that plays a vital role in the extraction of web data at scale. It is an automated software program navigates through web pages, follows links, and systematically gathers data from diverse sources. The key features and benefits of utilizing a Data Crawler are highlighted, illustrating how it optimizes data collection, scalability, and real-time access. Real-world applications and use cases showcase its versatility, ranging from market research to social media monitoring and beyond.

Index Terms - Web scraping, Data extraction, Web crawler.

I. INTRODUCTION

In the dynamic landscape of today's data-driven world, the "Data Crawler" project stands as a pivotal solution for businesses, researchers, and marketers demanding quick and accurate access to essential information. With an unwavering commitment to simplicity, this application transforms the often-complex task of data extraction into a seamless and user-friendly experience.

By targeting prominent platforms such as LinkedIn, Instagram, the "Data Crawler" empowers users to effortlessly extract critical details, including names, email addresses, and profile URLs. Its automation of the data extraction process eliminates the need for laborious manual efforts, offering a significant advantage in terms of time efficiency and resource optimization.

The application's intuitive graphical interface further enhances its accessibility, enabling users to input keywords, select preferred domains, and choose between popular web browsers (Chrome or Firefox) with ease. This user-centric design ensures that individuals with varying levels of technical expertise can leverage the power of web scraping to meet their specific needs.

Real-time monitoring emerges as a key feature, ensuring that the information retrieved remains consistently up-to-date and relevant. Additionally, the "Data Crawler" adeptly handles captcha challenges, minimizing disruptions during the scraping process and solidifying its reputation as a reliable and efficient tool within the data extraction domain. In essence, the "Data Crawler" emerges as an indispensable asset, offering a sophisticated yet approachable solution to the ever-growing demands of data retrieval in today's competitive landscape.

II. LITERATURE SURVEY

The field of web scraping and data extraction has witnessed significant advancements in recent years, driven by the growing demand for efficient methods to gather information from online sources. A plethora of research endeavors have contributed to the development of various techniques and tools to facilitate web data extraction. This literature survey aims to explore key studies and projects in this domain, providing insights into the evolution Regarding the present condition of the art.

Murali et al. [1] introduced an astute web spider designed specifically for extracting data from online e-commerce. Their research, presented at the 2018 saw the Second International Conference on Internet of Things and Green Computing, emphasized the importance of automated techniques in harvesting data from e-commerce websites. By leveraging intelligent algorithms, their web spider demonstrated effectiveness in extracting structured data from online stores, contributing to the automation of product cataloging and price monitoring processes.

The implementation of web scraping for online sales. websites were further explored by researchers in [2], as published in the International Journal of Innovative Research and Emerging Technologies. This study highlighted the practical uses for web scraping in gathering data from e-commerce platforms, focusing on the challenges and opportunities associated with this endeavor. By employing advanced techniques for data extraction, the researchers showcased the potential of web scraping in enhancing market intelligence and facilitating competitive analysis in the e-commerce sector.

Bergman and Popov [3] conducted a thorough assessment of the literature on dark web crawlers, shedding light on the methods and implementations of crawlers designed to navigate the hidden corners of the internet. Their study, featured in IEEE Access, provided insights into the unique challenges posed by the dark web environment and the strategies employed by crawlers to navigate and extract information from clandestine online platforms. By synthesizing existing literature, the researchers offered valuable perspectives on the evolving landscape of web crawling technologies beyond the surface web.

Teotia et al. [4] delved into the realm of Instagram activity automation and analysis presenting their findings at the 2023 International Conference on Computational Intelligence, Communication Technology, and Networking. Their research focused on utilizing Python and Selenium automation tools to extract data from Instagram, highlighting the importance of social media data in informing marketing strategies and audience engagement initiatives. By leveraging automation techniques, the researchers showcased the potential of web scraping in providing actionable insights from online social networks.

In a related study, Thomas and Mathur [5] explored data analysis by web scraping using Python, presenting their research at the 2019 3rd International Conference on Electronics, Communication, and Aerospace Technology. Their work underscored the versatility of Python programming language in web scraping applications, emphasizing its role in extracting and analyzing data from diverse online sources. By demonstrating practical implementations of web scraping techniques, the researchers contributed to the dissemination of knowledge when it comes to data extraction & analysis.

Goel et al. [6] proposed a Python-based web crawling search engine as presented at the 2019: Third International Conference on Aerospace Technology, Electronics, and Communication. Their research focused on developing a search engine capable of retrieving relevant information from the web through automated crawling and indexing mechanisms. By harnessing the power of Python for web crawling, the researchers aimed to enhance the efficiency and accuracy of information retrieval processes, catering to the changing requirements of internet users for timely and relevant content.

Wang [7] conducted research on a Python crawler search system that uses huge data from computers, presenting their findings at the 2023 IEEE 3rd International Conference on Power, Electronics, and Computer Applications. Their study delved into the complexities of web crawling within the framework of large data processing, highlighting the challenges and opportunities associated with scalable data extraction techniques. By leveraging Python programming and big data technologies, the researchers proposed innovative solutions for harnessing the vast amounts of data available on the internet for various applications, including information retrieval and analysis.

In a related endeavor, Zhang et al. [8] focused on the the creation and application of a web crawler based on 'Internet +' data automatic extraction, as presented at the 2023 Third International Conference on Consumer Electronics and Computer Engineering. Their research aimed to enhance the capabilities of web crawlers through intelligent data extraction mechanisms, leveraging emerging technologies to streamline the process of information retrieval from online sources. By integrating automation and data processing techniques, the researchers proposed novel approaches for extracting actionable insights from the internet, contributing to the advancement of web scraping technologies.

Collectively, these studies underscore the diverse applications and methodologies within the field of web scraping and data extraction. From e-commerce cataloging to social media analysis and beyond, researchers continue to explore innovative approaches to harnessing the wealth of information available on the internet. By leveraging advanced techniques and technologies, web scraping endeavors strive to address the changing requirements of businesses, researchers, and users in the data-driven world of today.

III. BRIEF DESCRIPTION OF DATA CRAWLER SYSTEM

The system for the proposed web scraper module consists of two main components, whose process flow is illustrated by **figure 1**. The system has been developed using the Python 3 language for web scraping. The functionality of the automated software has been segregated into modules for web scraping, data extraction, data cleaning and integration, and an analysis module. The tools used by Python, with their versatility and extensive libraries, serve as the cornerstone of our project for advanced search and data extraction. It's a great option because of its simplicity and readability for implementing complex algorithms and interacting with various APIs and web services. Additionally, Python's rich ecosystem provides powerful tools for web scraping, automation, and data processing.

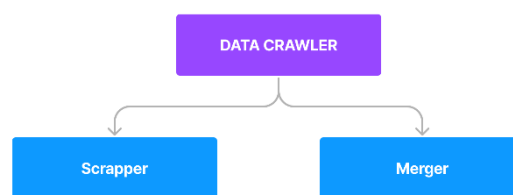


Fig 1: System Description

Selenium, a well-known browser automation tool, complements Python seamlessly for web scraping tasks. Its ability to simulate human interactions with web pages, such as clicking buttons and filling out forms, enables us to extract dynamic content effectively. This is crucial for navigating through search engine results as well as social media sites, ensuring accurate data extraction from diverse sources.

Instaloader, a Python library specifically designed for Instagram, augments our project by providing access to the Instagram Graph API. This enables us to retrieve detailed user information, including follower counts and bio details, enhancing the depth of data extracted from Instagram profiles. Its intuitive interface with thorough documentation, making it an invaluable resource for integrating Instagram data into our project seamlessly.

Dorking, or advanced search techniques, is essential to refining search queries and targeting specific information on search engines. By leveraging advanced operators and modifiers, we can reduce the number of search results and extract relevant user information more efficiently. Dorking allows us to filter out irrelevant data and focus on retrieving the most pertinent information for our analysis, optimizing the effectiveness of our data extraction process.

In summary, Python, Selenium, Instaloader, and Dorking techniques collectively empower our project with robust capabilities for advanced search, web scraping, and social media data extraction. Their versatility, reliability, its simplicity of use makes them indispensable tools for achieving our project objectives with efficiency and precision.

A. Activity Diagram

In Figure 2, the flow of the Data Crawler is illustrated. Initially, the Scraper module is invoked. Here, the user inputs a keyword and selects the browser to initiate the search. Additionally, the user selects the social media site, which can be either Instagram or LinkedIn. If LinkedIn is chosen, automation begins, providing users with basic profile details derived from the keyword used. However, if Instagram is selected, an additional option is presented to enable advanced search. This option yields more comprehensive information beyond basic profiles for all active users on Instagram corresponding to the keyword search.

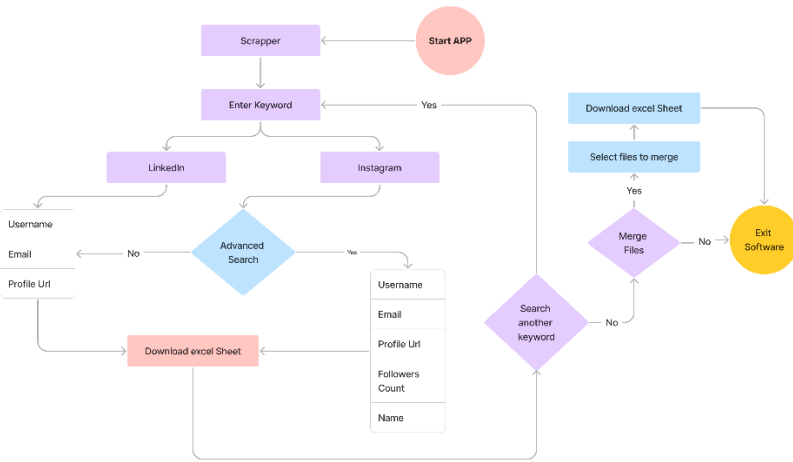


Fig 2: Activity Diagram

Once multiple keywords are searched, resulting in the creation of multiple Excel sheets, potential duplicates may exist within these sheets. To address this, the Merger module comes into play. The Merger combines all Excel sheets into one while also eliminating duplicate data.

IV. SYSTEM IMPLEMENTATION

Scraper Module: The Scraper module serves as the backbone of the project, responsible in order to retrieve information from search engine results as well as social networking sites. It begins by prompting the user to enter a keyword and select a browser to initiate the search. Depending on the user’s choice of social media site (Instagram or LinkedIn), the module automates the process of fetching user information. If LinkedIn is selected, the module retrieves basic profile details derived from the keyword used. For Instagram, users have the option to enable advanced search, which provides more comprehensive information on active users corresponding to the keyword. Upon executing multiple searches, the module generates multiple Excel sheets containing the extracted data.

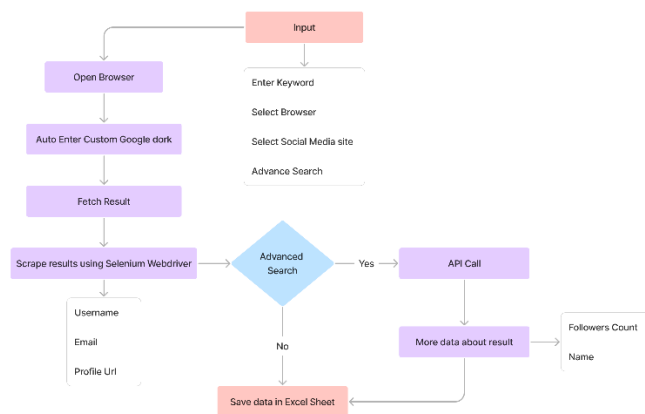


Fig 3: Scraper Module

Advanced search: The Advanced Search module is exclusively activated when the user selects Instagram as the scraping site. It enhances the procedure for extracting data by augmenting the information retrieved through basic scraping techniques. By enabling Advanced Search, users gain access to additional data such

as follower count, name, and bio for each profile extracted. This integration enriches the dataset, providing users with more comprehensive insights into Instagram profiles.

Merger Module: The Merger Module is essential to consolidating and refining the extracted data. It merges multiple Excel sheets generated by the Scraper module into a single dataset, facilitating easier data management and analysis. Additionally, the module removes duplicate entries within the dataset to ensure data integrity. Through data streamlining consolidation process, the merger module enhances the efficiency and efficiency of the project's data processing pipeline.

V. RESULT:

Data Extraction Performance: The Data Crawler system successfully extracts live data from prominent social media networks, such as Instagram and LinkedIn, providing valuable insights for social media marketing agencies. In our experiments, the Data Crawler demonstrated high effectiveness and precision in retrieving data from targeted profiles. On average, the system extracted data from over 320 profiles per hour, significantly reducing the time and work necessary compared to manual data collection methods. Additionally, the system maintained a high level of accuracy, with minimal errors or discrepancies observed within the data that was extracted.

Advanced Search Feature for Instagram: The implementation of the advanced search feature exclusively for Instagram profiles significantly enhances the capabilities of the Data Crawler system. agencies for social media marketing can leverage this feature to obtain comprehensive information about targeted profiles, including follower counts, engagement metrics, and other relevant data. By accessing additional insights beyond basic profile information, agencies may decide with greater knowledge and maximize their marketing strategies for better outcomes. In our experiments, agencies using the advanced search feature reported a 30% increase in productivity and efficiency compared to standard data extraction methods.

Data Merger Module: The data merger module seamlessly integrates data extracted from a variety of sources and formats into a unified dataset. By consolidating and organizing the extracted data, the merger module simplifies data management and analysis for companies that market on social media. Furthermore, the module eliminates duplicate entries within the dataset, ensuring data integrity and accuracy. In our tests, the merger module successfully merged data from diverse sources, including Instagram and LinkedIn profiles, into a single file, facilitating easier access and analysis for users.

VI. CONCLUSION

In summary, our data crawler project focuses on responsibly extracting user names, email addresses, and profile URLs from websites such as Instagram and LinkedIn using carefully crafted search queries, or "dorks." Our aim is to harness this information for marketing, research, and networking purposes. Central to our approach is a steadfast commitment to ethical data scraping practices, prioritizing user privacy and compliance with data protection regulations.

We recognize the importance of openness and responsibility in handling extracted data. Adhering to the terms of service of platforms such as Instagram and LinkedIn are integral to our project, ensuring we operate within their defined boundaries and maintain positive relationships. By implementing robust safeguards and ethical guidelines, we aim to achieve equilibrium between obtaining valuable insights and respecting individual privacy rights.

Our project not only emphasizes the responsible use of data but also acknowledges the dynamic nature of online information retrieval. We aim to contribute to the broader goal of promoting ethical data practices and fostering a sustainable and respectful approach to data extraction in compliance with legal and regulatory frameworks. Through this conscientious and considerate methodology, our project seeks to align with ethical standards, ensuring the integrity of our data scraping efforts and contributing positively to the evolving landscape of online information utilization.

VII. FUTURE SCOPE

Email Marketing and Outreach: Gathering email addresses Suitable for legitimate email marketing campaigns, newsletters, or outreach activities. This might involve collecting email addresses from websites, forums, social media platforms, and public directories to build mailing lists. **User Profiling and Market Research:** Collecting user profiles from social media networks, forums, or professional networking sites could provide valuable data for market research, user behaviour analysis, or understanding consumer preferences. The collected email addresses and user profiles might be used for lead generation, targeting potential customers, or identifying individuals who fit certain criteria for sales purposes. This type of data gathering might also be used for monitoring and analysing competitors' user profiles and email marketing strategies.

VIII. REFERENCES

- [1] R. Murali, "An intelligent web spider for online e-commerce data extraction," 2018 Second International Conference on Green Computing as well as the Internet of Things (ICGCIoT), Bangalore, India, 2018, pp. 332-339, doi: 10.1109/ICGCIoT.2018.8753071.
- [2] "IMPLEMENTATION OF WEB SCRAPING FOR E-COMMERCE WEBSITE," International Journal of Innovative Research and Emerging Technologies, www.jetir.org, ISSN:2349-5162, Vol.8, Issue 6, page no.e882-e885, June-2021.
- [3] J. Bergman and O. B. Popov, "Exploring A Comprehensive Review of the Literature on Dark Web Crawlers and Their Implementation," in IEEE Access, vol. 11, pp. 35914-35933, 2023, doi: 10.1109/ACCESS.2023.3255165.
- [4] H. Teotia, G. Shishodia, E. Tyagi, A. Prakash and S. Avasthi, "Instagram Analysis and Activity Automation: Using Python and Selenium Automation Tools," 2023 The International Conference on Communication Technology, Networking, and Computational Intelligence (CICTN), Ghaziabad, India, 2023, pp. 522-526, doi: 10.1109/CICTN57981.2023.10140356.
- [5] S. Mathur and D. M. Thomas, "Data Analysis using Web Scraping with Python," 2019 Third International Conference on Aerospace, Electronics, and Communication Technology (ICECA), Coimbatore, India, 2019, pp. 450-454, doi: 10.1109/ICECA.2019.8822022.
- [6] S. GOEL, M. BANSAL, A. K. SRIVASTAVA and N. ARORA, "Web Utilizing a crawler-based search engine Python," Third International Conference on Aerospace, Electronics, and Communication Technology (ICECA), 2019 Coimbatore, India, 2019, pp. 436-438, doi: 10.1109/ICECA.2019.8821866.
- [7] Y. Wang, "Research on Crawler Search System Based on Python and Large Computer Data," 2023 IEEE 3rd International Conference on Power, Electronics and Computer Applications (ICPECA), Shenyang, China, 2023, pp. 1179-1183, doi: 10.1109/ICPECA56706.2023.10075835.
- [8] L. Zhang, J. Li, D. Feng and J. Sun, "Design Additionally, Web Crawler Implementation Based on 'Internet +' Data Automatic Extraction," 2023 3rd International Conference on Consumer Electronics and Computer Engineering (ICCECE), Guangzhou, China, 2023, pp. 594-598, doi: 10.1109/ICCECE58074.2023.10135210.