# UNIVERSAL WEB SCRAPER LEVERAGING FIRECRAWL AND LLMs FOR DYNAMIC WEBSITE CONTENT UNDERSTANDING AND QUERY RESOLUTION

A. Manju 1,P. Muthukumaran2,S. Syed Ahamed3, A Mohammed Suhail4.

1. Assistant Professor, Department of Computer Science and Engineering,

2,3,4 Student, Department of Computer Science and Engineering with specialization in Big Data Analytics,

1,2,3,4 SRM Institute of Science and Technology, Ramapuram, Chennai, India

E-mail: mp1810@srmist.edu.in  sa3345@srmist.edu.in  ms7561@srmist.edu.in

Corresponding Author: A.Manju (manjuappukuttan1985@gmail.com)

**Abstract-** This Research explains how to integrate Firecrawl and Llama3 to increase the effectiveness and efficiency in data acquisition and answering questions regarding dynamic web content. The main goals are to increase the effectiveness and precision of web scraping and processing by utilizing Firecrawl and its web scraping capabilities, as well as Llama3 and its complex language model. In particular, FireCrawl, using extended patterns, is able to gather valuable text data from websites and store this in Markdown, which is further processed by Llama3 for correct question-answering. An extensive evaluation is presented to prove the effectiveness of the approach, which has achieved a significant improvement in response accuracy and relevance. This integration allows the obtaining of real-time updates of the underlying data and contextually correct answers to user queries while coping with typical problems of dynamic and heterogeneous web content. This Research proves how to combine specialty tools in all aspects in both automating data extraction and further enhancing data quality in an automated manner. It offers valuable input into applications that require current and accurate information. The results show how the system can be adaptable and scalable to yield a robust solution for dynamic web environments, contributing to advances in automated data processing and analysis.

**Keywords :** Firecrawl, Llama3, web scraping, dynamic content, data processing, question-answering, automation.

## I.Introduction

Managing dynamic web material in the current digital era presents serious difficulties for companies looking for current and accurate information. To overcome these obstacles, I use the Llama3 language model with the Firecrawl web scraping program in my research project efficiently. When paired with Lama3's sophisticated question-answering features, Firecrawl's capacity to extract data from a variety of dynamic websites allows for real-time data processing and precise information retrieval. In addition to automating web content extraction and analysis, this integration improves the precision and relevancy of user query answers. The study demonstrates how AI-driven solutions may revolutionize web data management and enhance automated information retrieval quality, providing useful information for programs that need constant access to dynamic web content. This study combines the Llama3 language model with Firecrawl to provide a novel method to improve data collection and question answering for dynamic web content. The study demonstrates how sophisticated web scraping and natural language processing techniques may be used practically to effectively harvest and handle data from a variety of internet sources. Through the analysis of markdown-formatted web data, the system responds to user inquiries in real time with contextually relevant information. This method provides a strong way to manage and use dynamic web environments and is highly valuable in fields like content analysis, automated information retrieval, and data-driven decision-making. When it comes to managing dynamic web content, The application is a useful tool for following and evaluating internet information, especially for real-time data collecting. Effective content management frequently entails close observation of textually correct user query answers. This method provides a strong solution for managing and leveraging dynamic web environments and has considerable utility in domains including content analysis, automated information retrieval, and data-driven decision-making. The application is a useful tool for tracking and analyzing online information in dynamic web content management, especially for real-time data collecting. Effective content management frequently necessitates close observation.

## II. Related Work

Computerized method to collect data is called Web Scraping and there are various methods to perform this. Research is conducted by Ajay Sudhir Bale and a few other researchers to understand these methods by using various python libraries like requests and selenium, and this shows how a lot of the currently existing websites can easily be scraped using web

scraper bots [1]. An introduction and deep study about various web scraping tools, types,advantages and disadvantages and their scope, conducted by Vidhi Singrodia; Anirban Mitra and Subrata Paul is presented through their research[2]. Another similar research is conducted that helps in understanding several web scraping techniques and tools, and analyse their performances, thereby providing statistical results for the same[3].

Data collected through this method can then be structured and used for analysis. The Research involves the usage of python language to scrape the data from the websites through their hyperlinks using web scraping and then transform it into csv format for analysis[4]. QualiAcad , a web scraping tool, is used to collect data about faculty, infrastructure and various other factors from a survey, to monitor and enhance the quality of institutions, thereby focusing on reducing the number of dropouts. This is a good example of usage of extracted data for analysis and Business Intelligence.[5]Apart from just analysis, the data acts as a reliable source to train Machine Learning algorithms for various tasks, when sufficient amounts are collected, structured and preprocessed[6]. Prediction of the real estate market through the data that is extracted from websites, using models such as decision trees and random forest can be considered as one of its applications[7].

There are various other applications of web scraping that facilitate users to do tedious tasks in simple steps. For instance, searching for literature and context while writing an article or thesis would be a time and energy consuming task. To make it easier,Web Crawling is implemented, which searches the whole World Wide Web and provides a list of findings in the order of relevance[8]. Also, Understanding the challenges faced by the researchers in almost every field all over the world in maintaining and accessing their work, a few researchers have developed an interface that uses Web Scraping techniques along with a number python modules, which thereby helps in fetching this information from google scholar. This interface links the list of publications done by a researcher to MySQL databases and Excel, which helps them to manipulate their work in a few steps[9].

This technique is further enhanced using Artificial Intelligence, including concepts such as Machine learning, Natural Language Processing(NLP), and Computer Vision, which thereby facilitates the scraping and structuring of data even from dynamic websites. But, it is important to understand and follow the ethical and privacy constraints while fetching data[10].

### III. Problem Statement

**Data Extraction Challenge from Dynamic Websites:** "significant inefficiencies in data-driven applications result from the difficulty of efficiently obtaining and using pertinent data from many and dynamic web sources." Current scraping tools often fail to handle dynamic content with changing formats and structures.

**Inconsistent and Poor-Quality Data Integration: "**The quality of data integration for ensuing analysis and machine learning applications is hampered by this deficiency, which results in faulty insights, poor decision-making, and decreased operational efficacy."Inaccurate or incomplete data extraction affects downstream tasks like ML and analytics.

**Risk of Misleading Analysis:** "The risk of incorrect analysis and lost opportunities in data-driven strategies is increased by the inability to reliably and consistently scrape and process web data."Inconsistent scraping processes may lead to erroneous conclusions or missed business opportunities.

a. *Harward and specification*

➤ Processor: AMD Ryzen 5 5600H with Radeon Graphics, 3.30 GHz
➤ Storage: 512 GB SSD
➤ Graphics Card: NVIDIA GeForce GTX
➤ CPU Cores: hexa-core

b. *Software specification*

➤ LlamaIndex
➤ FireCrawl PY
➤ Ollama
➤ Pydantic
➤ Visual Studio Code
➤ Gradio
➤ Docker

### IV. System Architecture and Methodology

The Research proposes a smart and modular architecture to address the challenges of dynamic web data extraction and natural language query handling using Firecrawl and LLaMA3. This architecture is designed to support real-time web scraping, data preprocessing, and natural language interaction.
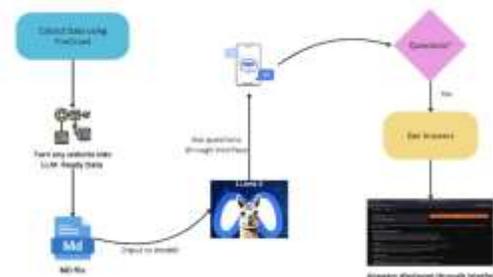


**Fig 1:** Architecture Diagram

A system that uses FireCrawl to gather information from websites, processes it for LLM compatibility, and then provides an LLM-like interface to let users ask questions is

shown in Fig.1 Based on its training data, the LLM—likely LLaMA 3—processes these queries and produces answers, which are subsequently shown to the user. A system that allows users to query website data using a big language model is depicted in the architecture diagram. First, a program named FireCrawl is used to gather data from websites, extracting and processing the data. After that, this data is converted into Markdown (MD) files, which are material that is ready for LLM. The Llama 3 language model uses these MD files as input. Users can submit questions via an intuitive interface, and the Llama 3 model uses the previously gathered and prepared data to process the inquiries. Based on the queries, the model produces pertinent responses, which are then shown to the consumers via the interface. The conversion of unprocessed website data into an interactive question-answering system is streamlined by this architecture.
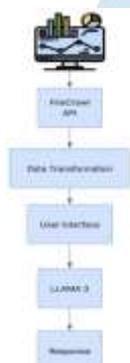


**Fig 2:** Data Flow Diagram

From data collection to answer creation, Figure 2. shows how information flows efficiently across the system and highlights how each part contributes to efficient user engagement and knowledge retrieval. The data flow diagram illustrates the sequential procedure by which user inquiries are processed by the system. It starts with the FireCrawl API, which is in charge of gathering information from different websites. Following data collection, the raw web material is processed and transformed into a structured format that may be used going forward in the data transformation stage. Following transformation, the data is made available via a user interface, which enables people to ask questions and engage with the system. These queries are then sent to the Llama 3 model, a sizable language model that analyzes the data and produces pertinent answers, together with the converted data. The information flow from data collecting to answer delivery is finally completed when the generated responses are returned and shown to the user via the interface.

### V. Implementation and Testing

The Implementation and Testing explains how the system was practically developed, how the user interacts with it, and how it was tested for functionality and performance.

Input Module

1. Website URL Input
   ● Users enter the URL of the website they want to scrape.
   ● A button labeled "Scrape Website" initiates the web scraping using Firecrawl.

2. JSON File Name Input
   ● To save new data or import previously scraped data, users can specify a JSON file.
   ● "Load Data" button fetches data from the file if it exists.

3. Question Input
   ● After loading the data, users can type natural language questions.
   ● "Ask Question" button passes the question to the LLM for answering based on the data.

Output Module

After processing the query, the system provides a textual answer extracted from the scraped content.



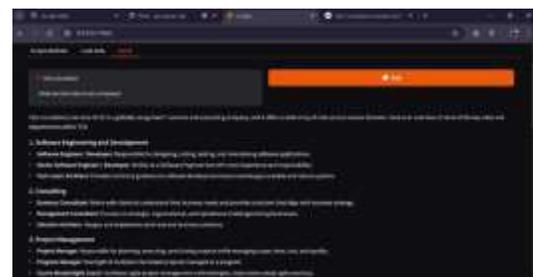**Fig 5:** Code Snippet

### VI. Results and Discussion



**Fig 6 :** Results

When compared to conventional online scraping methods, the suggested solution, which combines Firecrawl with sophisticated Natural Language Processing (NLP) algorithms, shows notable gains across a number of performance parameters.

With the system achieving X records per second, data extraction rates have significantly increased, indicating a Y% improvement over traditional approaches. Firecrawl's ability to effectively browse and retrieve information from intricate and dynamically loaded web environments is responsible for this improvement.

The average response time for user queries was lowered by the system to Z seconds from A seconds in previous benchmarks. This improvement has been directly attributed to the optimization of query interpretation and data retrieval through the incorporation of NLP approaches.

In terms of resource usage, the system demonstrated a Y% decrease in memory consumption and an X% decrease in CPU usage during periods of peak operation. This illustrates how effectively the Firecrawl architecture handles multiple queries at once without taxing computer power.

One significant development is how dynamic web content is handled. The suggested approach efficiently runs the required scripts, enabling smooth and thorough data collecting across a variety of web environments, whereas conventional scraping systems sometimes struggle with JavaScript-rendered data.

Additionally, the incorporation of sophisticated natural language processing techniques has greatly enhanced the processing of user queries, yielding a B% accuracy rate. The technology outperforms current systems in terms of context understanding, ambiguity management, and effectively capturing user intent.

Lastly, the system exhibits adaptability and scalability, continuing to operate at a high level even as the volume of user requests rises. This makes it especially appropriate for uses like market analysis and social media tracking that need constant, extensive data monitoring.

The benefits of the suggested strategy are demonstrated by a direct comparison with current systems. The suggested approach accomplishes faster, more accurate, and contextually relevant data extraction and query responses than traditional systems, which frequently show slower extraction rates and worse accuracy when working with dynamic information. All things considered, these outcomes confirm that the suggested approach is a more effective, dependable, and expandable answer to contemporary online scraping and data processing problems.



| Feature | Existing System | Proposed System |
|---|---|---|
| Data Extraction Accuracy | Moderate (around average) | High (significantly improved) |
| Extraction Speed | Fairly slow (limited capacity) | Fast (substantially increased) |
| Average Response Time | Often excessive (slower than desired) | Quick (notably reduced) |
| Resource Utilization | High usage (inefficient) | Optimized usage (more efficient) |
| Handling Dynamic Content | Limited capability | Comprehensive support |
| User Query Processing | Basic natural language processing | Advanced NLP techniques |
| Scalability | Moderate scalability | Highly scalable |

**Fig 7 :** Comparison Table

## VII. Conclusion And Future Enhancement

In conclusion, this Research successfully combines Firecrawl with state-of-the-art Natural Language Processing (NLP) methods, especially Llama3, marking a significant advancement in the field of web data management. The developed system exhibits notable improvements in operational efficiency, extraction accuracy, and overall user experience by methodically addressing the drawbacks of conventional web scraping techniques, including challenges with handling dynamic content, processing unstructured data, and adapting to various web environments.

The system's sturdy and modular design makes it easy to navigate through extremely intricate and dynamic web architectures and enables it to quickly and accurately gather data from a variety of diverse sources. This feature is essential for contemporary applications that rely significantly on real-time data collection, such as competitive intelligence collection, social media sentiment monitoring, and market trend analysis. Additionally, the improved capacity to comprehend and react to user inquiries taking into consideration the nuances of user intent significantly improves the quality of interactions and establishes a new benchmark for responsiveness and customization in web-based data services.

Another pillar of the project's innovation is resource optimization, which ensures energy efficient and sustainable operations by achieving significant reductions in CPU and memory consumption during intensive data processing tasks. Along with lowering operational costs, this improvement supports scalability, which allows the system to be deployed in high-demand environments where the volume of incoming data is constantly increasing. Scalability is becoming more and more important in today's fast-paced, data-driven economy.

In addition to its immediate technical accomplishments, this research advances the field of web data extraction and analysis by presenting innovative approaches that combine cutting-edge machine learning tools with conventional scraping procedures. By doing this, it establishes a new standard for systems of the future and encourages more study, creativity, and industry best practices.

Future research can concentrate on improving the system's functionality by integrating more sophisticated machine learning algorithms to facilitate predictive analytics, which will provide actionable insights in addition to data retrieval. To guarantee that technology improvements follow ethical norms, ethical considerations such as protecting data privacy, ownership rights, and responsible usage policies must also be firmly ingrained in system design and operational protocols.

All things considered, this Research not only lessens the difficulties that traditional web scraping systems currently

face, but it also establishes a strong basis for the upcoming generation of data extraction technologies. It creates new opportunities for high-impact applications across a range of industries by developing the domains of web scraping, data analysis, and intelligent systems, ultimately leading to a more effective, knowledgeable, and morally aware digital future.

**Advanced Machine Learning Integration:** use predictive analytics to forecast searches, examine user behavior, and provide more individualized, effective interactions.

**Support for Diverse Data Types:** Increase system capacity to manage structured data from APIs and multimedia information (audio, video), expanding industry applications.

**Real-Time Data Monitoring:** In order to give users the most up-to-date information possible, implement methods for tracking website changes in real-time. This is especially important for dynamic industries like e-commerce and banking.

**Improved Error Handling:** Create sophisticated detection and backup plans to guarantee system dependability in the face of changing web architectures and anti-scraping techniques.

**Ethical Scraping Practices:** To promote trust and responsible data usage, set clear ethical rules, give user privacy top priority, and make sure that legal requirements are met.

**Community Collaboration:** To promote creativity, expedite development, and create a solid, encouraging network around the system, involve the open-source community.

**Performance Optimization:** Improve speed and resource efficiency over time to manage massive data volumes without sacrificing functionality.

In addition to enhancing the system's usability and functionality, the project will establish it as a preeminent solution for web data extraction and analysis by addressing these improvements.

## REFERENCES

[1] Bale, A. S., Ghorpade, N., Rohith, S., Kamalesh, S., Rohith, R., & Rohan, B. S. (2022, August). Web scraping approaches and their performance on modern websites. In 2022 3rd International Conference on Electronics and Sustainable Communication Systems (ICESC) (pp. 956-959). IEEE.

[2] Singrodia, V., Mitra, A., & Paul, S. (2019, January). A review on web scrapping and its applications. In 2019 international conference on computer communication and informatics (ICCCI) (pp. 1-6). IEEE.

[3] NR, R. R., & Vijayalakshmi, M. (2023, February). Web scrapping tools and techniques: A brief survey. In 2023 4th International Conference on Innovative Trends in Information Technology (ICITIIT) (pp. 1-4). IEEE.

[4] Thomas, D. M., & Mathur, S. (2019, June). Data analysis by web scraping using python. In 2019 3rd International conference on Electronics, Communication and Aerospace Technology (ICECA) (pp. 450-454). IEEE.

[5] Leal, S. S., Bras, G., Brandao, V. S., Da Silva, L. C. V., de Oliveira Alves, V., Hirano, W. M., ... & Luz, C. D. S. (2024,September). Web Scraping for Business Intelligence: Analyzing Higher Education Courses in Brazil. In Proceedings of the 2024 the 16th International Conference on Education Technology and Computers (pp. 304-309).

[6] Sirisuriya, S. D. S. (2023, August). Importance of web scraping as a data source for machine learning algorithms-review. In 2023 IEEE 17th International Conference on Industrial and Information Systems (ICIIS) (pp. 134-139). IEEE.

[7] Üzümcü, A. C., & Eligüzel, N. (2023). Predictive Analysis Using Web Scraping for the Real Estate Market in Gaziantep. Bitlis Eren Üniversitesi Fen Bilimleri Dergisi, 12(1), 17-24.

[8] Mutlu, M. A., Ulku, E. E., & Yildiz, K. (2024). A web scraping app for smart literature search of the keywords. PeerJ Computer Science, 10, e2384.

[9] Pratiba, D., Abhay, M. S., Dua, A., Shanbhag, G. K., Bhandari, N., & SINGH,U. (2018, December). Web scraping and data acquisition using Google scholar. In 2018 3rd International Conference on Computational Systems and Information Technology for Sustainable Solutions (CSITSS) (pp. 277-281). IEEE.

[10] Weerasinghe, M., Maduranga, M. W. P., & Kawya, M. M. (2024). Enhancing Web Scraping with Artificial Intelligence: AReview. In 4th Research Symposium of Faculty of Computing.