

Project Summary: Rufus - AI-Powered Web Scraper for RAG Systems

- Kunal Shah
- kunaljshah03@gmail.com

GitHub: <https://github.com/kunalshah03/rufus>

Overview

As a recent Master's graduate with experience in web scraping and NLP from my Omdena internship, I took on the challenge of building Rufus, an AI-powered web scraping tool specifically designed for RAG systems. The project allowed me to combine my experience with BeautifulSoup and Scrapy with newer technologies like OpenAI's GPT models.

Technical Approach

Building on my experience with web scraping during my Omdena internship, I designed Rufus with a modular architecture focusing on three core components:

- Asynchronous web crawler
- AI-powered content processor
- RAG-optimized document structuring

Challenges and Solutions

1. Dynamic Content Handling

Coming from traditional scraping tools like BeautifulSoup, I initially struggled with dynamic JavaScript-rendered content. Given my background in distributed systems from my Search Engine project, I implemented the following:

- Asynchronous crawling using aiohttp
- Playwright integration for JavaScript rendering
- Smart content detection mechanisms

2. AI Model Integration

While I had experience with NLP from my Omdena internship, integrating OpenAI's API presented new challenges:

- Initially faced rate-limiting issues
- Spent \$10 on API credits for testing
- Implemented intelligent rate limiting and error handling
- Optimized prompt engineering for better content extraction

3. RAG-Specific Optimization

Drawing from my experience with the FinanceGPT project where I worked with RAG systems:

- Designed document structure specifically for RAG pipelines
- Implemented efficient chunking mechanisms
- Added metadata enrichment for better retrieval
- Created JSONL export functionality

4. Scalability Concerns

Leveraging my experience with microservices from my Carikture internship:

- Implemented concurrent request handling
- Added configurable crawling limits
- Created efficient memory management
- Designed for horizontal scalability

Technical Decisions

Architecture

Based on my distributed systems experience from the Search Engine project:

- Modular design for maintainability
- Clear separation of concerns
- Robust error handling
- Comprehensive logging

Performance Optimization

Drawing from my work on the WolfMedia project:

- Implemented async operations
- Added batch processing
- Created configurable limits

- Efficient resource management

Results and Impact

- Successfully crawled and processed various web structures
- Generated RAG-ready documents with rich metadata
- Achieved stable performance with rate-limiting
- Created comprehensive documentation for easy integration

Learning Outcomes

- Deepened understanding of RAG systems
- Gained practical experience with OpenAI's API
- Improved async programming skills
- Enhanced error-handling strategies

This project demonstrates my ability to build scalable, production-ready software while combining my academic knowledge with practical experience from internships. The challenges were turned into learning opportunities, resulting in a robust tool that meets real-world requirements.