

Parallel Job Miner

*Note: Sub-titles are not captured in Xplore and should not be used

Jason Saini CECS University of Central Florida Orlando, FL jason.saini@ucf.edu	(of Affiliation) line 3: <i>name of organization (of Affiliation)</i> line 4: City, Country line 5: email address or ORCID	(of Affiliation) line 3: <i>name of organization (of Affiliation)</i> line 4: City, Country line 5: email address or ORCID	(of Affiliation) line 3: <i>name of organization (of Affiliation)</i> line 4: City, Country line 5: email address or ORCID
line 1: 2 nd Given Name Surname line 2: <i>dept. name of organization</i>	line 1: 3 rd Given Name Surname line 2: <i>dept. name of organization</i>	line 1: 4 th Given Name Surname line 2: <i>dept. name of organization</i>	
	line 1: 5 th Given Name Surname line 2: <i>dept. name of organization (of Affiliation)</i> line 3: <i>name of organization (of Affiliation)</i> line 4: City, Country line 5: email address or ORCID	line 1: 6 th Given Name Surname line 2: <i>dept. name of organization (of Affiliation)</i> line 3: <i>name of organization (of Affiliation)</i> line 4: City, Country line 5: email address or ORCID	

Abstract— The technology industry encompasses a diverse and continually evolving landscape, presenting a wealth of opportunities for computer science students on the cusp of graduation. For computer science students, the transition from academic study to professional practice is a critical phase of their careers. The job search process for students can be time-consuming and often inefficient, with relevant positions scattered across numerous platforms. Our team decided to make it easier for students to find job applications online with a parallel “job miner”, which will consolidate job application data from multiple popular job sites for computer science: LinkedIn, GitHub/Simplify, USAJobs (for those interested in government positions), Monster, and Indeed. We accomplished this by using a process called “web scraping” to gather job opportunity data from the web. In other words, our program gets necessary information about jobs from the sites listed above and provides that information to our users in the form of an xlsx (Excel) file. With this Excel file, users will be able to adopt this strategy of tracking positions they are interested in.

I. INTRODUCTION

Transitioning from studying in college to the ever-changing workforce requires finding jobs and internships within one’s field. This process often involves sending many applications and receiving mostly rejections. To make the process easier for students, particularly those majoring in Computer Science and Software Engineering, we have created a web scraping application.

The project is open source, and therefore free, to all that would like to use it. Being open source enables all those who wish to contribute and expand/improve upon it to do so. To get started, users simply need to follow the download and setup instructions on the [GitHub page](#), allowing them to run it

Identify applicable funding agency here. If none, delete this text box.

as needed. We were able to achieve all of this through using pre-existing Application Programming Interfaces (APIs), creating our own web scraping functions, and storing collected info into a DataFrame from Python’s Pandas library.

II. PROBLEM STATEMENT

A. Problem

Students encounter difficulty when applying to many different positions or companies trying to find relevant opportunities across various websites. After a website with relevant listings and the desired opportunities, they then must go into each card to read the details and click the actual application link if they are interested. Navigating through a multitude of cards on multiple sites for jobs that may not be relevant or applicable to you can quickly become overwhelming.

B. Proposed Solution

To help combat this, our project stores much of the needed information into a single file that is easy to read and contains the application links directly to all jobs found. From this spreadsheet, applicants can track their progress of individual applications using color coding, as demonstrated in the aforementioned commonplace strategy for tracking applications.

III. RELATED WORK

There are many resources available that scrape the sites mentioned above, along with many others, to collect job data. Where our project differs is that we are the only resource which collects job data from multiple sites at once, providing users with a comprehensive list of opportunities without needing much technical knowledge.

To achieve this, we use existing APIs to access content and data. APIs allow our project to interact with the platforms

efficiently and collect the needed information in a straightforward manner. This approach allows us to streamline data retrieval and adapt to anti-scraping measures used by major job finding sites.

As noted above, many job finding companies implement constantly evolving anti-scraping features. These features include bot detection, dynamic loading, and including unique job ids into the URL. Having these measures in place makes basic web scraping challenging, but APIs help combat that challenge and enable our data collection from various sites.

Additionally, our program is easy to use and designed to be efficient. The user only needs to download the files, follow the quick setup on the GitHub page, and run the file with a title to receive 100+ job opportunities from 4 sites related to what they are looking for. The program is made to be efficient by utilizing threads to scrape all the websites at once, instead of one by one. While existing resources do exist, our project offers a unique version by scraping and consolidating information from multiple sites at once.

IV. TECHNIQUE

Our approach to a multi-threaded web scraper was using BeautifulSoup for the manual web scraping, RapidAPI to get our APIs, and Python's multiprocessing library to run each scraper on a different thread.

BeautifulSoup is a popular Python library used for web scraping. We use it in conjunction with the Requests library, which sends a get request to a provided link and returns the contents. After getting the content, you pass that through to BeautifulSoup to parse the content it was given, we choose to use the HTML parser.

Now that you have the HTML content, you can use BeautifulSoup to then traverse the page via its html tags. RapidAPI is a service that allows users to publish their own API and use others' published APIs, either free or for a cost. Most published APIs tend to have multiple price plans which allow them a certain number of requests.

Python's multiprocessing library enables a program to run processes on separate threads concurrently. It allows one to assign a 'Process' to a function with any arguments it may need. When the thread is started, it executes that function with the given arguments.

We use these in conjunction with each other to create our current application. BeautifulSoup allows us to find the needed content on the pages that we manually scrape,

RapidAPI returns us the data of jobs that it gets, and multiprocessing allows us to run the scrapers concurrently instead of sequentially.

To ensure the scalability of our program we have a clean and organized file structure that is modular. For example, you can add a scraper to another site by simply adding a file for it and adding the function to the main file.

V. EVALUATION

Initially, our team hypothesized that the sequential and concurrent methods of web scraping result in very similar execution times, with the sequential method having the faster runtime around 70% of the time. This, however, was mainly due to the need to finish adding scraping/querying functionalities for multiple job sites such as Monster, Glassdoor, and Indeed.

After adding multiple sources for jobs and more job listings in general, we observed that the execution times started to deviate away from each other, with the concurrent method having a faster runtime.

```
Projects\ParallelJobMiner> &
Sequential search complete in 2.19 seconds
Projects\ParallelJobMiner> &
Concurrent Search complete in 1.18 seconds
```

After running our concurrent and sequential algorithms separately, we observed a nearly 86% decrease in runtime from end to end, from processing the user specified job title to generating the spreadsheet. While a noticeably small increase in the quantifiable runtime, a decrease of this magnitude on a much larger scale will be incredibly beneficial to those compiling job applications in larger quantities

VI. LIMITATIONS AND CHALLENGES

In the early stages of the project, we discovered that websites (i.e. LinkedIn) are weary against web scrapers for reasons including protecting intellectual property, preserving bandwidth and server resources, maintaining data accuracy and quality, protecting user privacy, preventing competitive advantage, avoiding legal complications, etc.

Our quick workaround was to use the official APIs, but of course these come with their own set of limitations. For example, the free LinkedIn API we are using allows 25 calls per month, which will not be serviceable for a full production launch. This is another challenge for scalability and maintenance, and a significant consideration when incorporating a data source into our parallel job scraper.

Another considerable challenge is accessing our pandas dataframe through threads, which requires a thread-safe extended class of the dataframe. A final challenge is finding ways to elevate the user experience from a heavily manual setup process to a more direct way of accessing our application.

VII. FIRST MILESTONE

On March 8th, 2024, we reached a significant milestone in our project at this point, using a modular approach to thread different scrapers and API requests, we can quickly compile relevant job application information for prospective applicants.

[illegible]

IEEE conference templates contain guidance text for composing and formatting conference papers. Please ensure that all template text is removed from your conference paper prior to submission to the conference. Failure to remove template text from your paper may result in your paper not being published.

