# Code Explanation, Outputs And Screenshots :

## 1. Web Crawler :

### Code Explanation :

Scrapping in Simply hired and GlassDoor. The crawlWebPage method is part of a web scraping application that extracts job-related information from a website based on specified search terms. It utilizes Selenium WebDriver with a browser to navigate through the website, search for jobs corresponding to each term, and scrape relevant details from individual job pages. The method iterates through the provided search terms, accesses the search results page for each term, extracts job links, and continues to subsequent result pages until at least min job links are collected. It then proceeds to visit each job page, extracts job data using the scrapeJobData method, performs data validation, and stores the validated job information in a collection. The scrapJobLinks method uses Jsoup to parse the HTML source of the search results page, extracts job links, and adds them to a queue. The uniqueJobs set ensures that duplicate job entries are not included in the final collection. The extracted job data is then saved and appended to a JSON file using the saveAndAppendToJson method of the ScraperBot class. The WebDriver is closed at the end of the process.
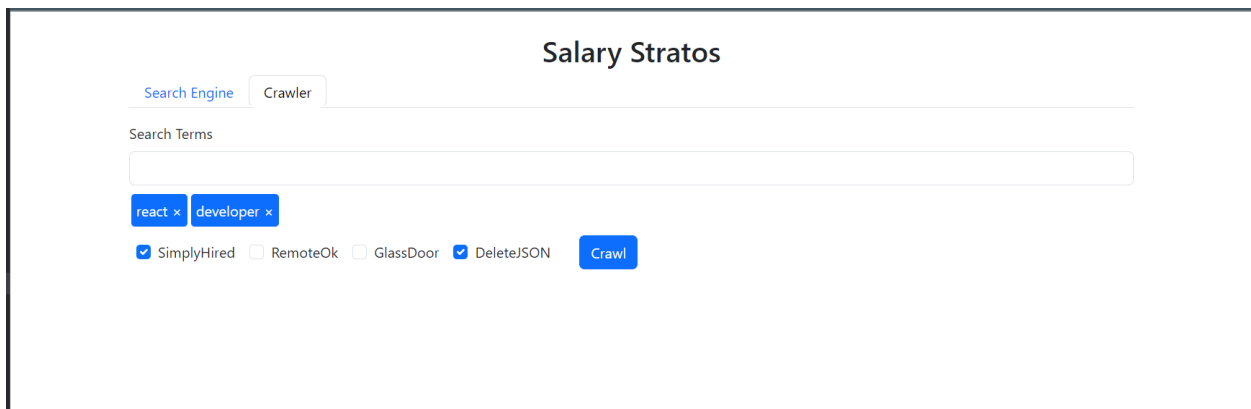
The only difference with scrapping Remote Ok is that the data we need is already existing on the search page and we don't need to open each and every job to get all the data.

### Output:

Users can enter the search terms, which website they want to crawl and whether they want to delete the json or append to the same database.json. The driver opens the chrome browser and starts crawling the website with the search terms to get the page source.

### Output 1: User wants to delete the database.json and crawl only Simply Hired for search teams "react", "developer"

User clicks on crawl



File deleted successfully

2023-12-04T06:53:23.655-05:00  INFO 32836 --- [nio-8080-exec-1] o.s.web.servlet.DispatcherServlet
JSON file deleted successfully.
SimplyHired Crawling Started

It searches for react and gets all the links on that page

It opens each link and gets all the job data from that link

**Output 2: <u>User wants to crawl only Remote Ok for search teams "react", "developer" and wants to append to the database.json</u>**
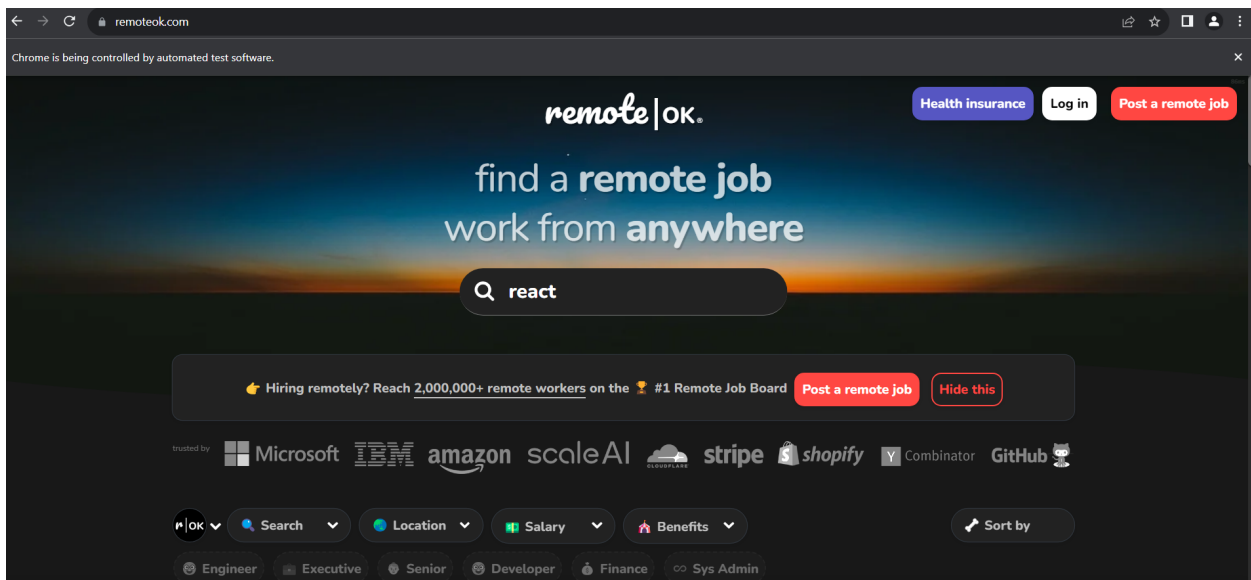
User clicks on crawl

Search for "react"



It gets all the jobs on this page

**Output 3 : If file does not exist when starting the application it will crawl the data for default search terms and create the database.json file.**

```
String[] searchTermsList = new String[]{
        "Engineer", "Exec", "Senior", "Developer", "Finance", "Sys
Admin", "JavaScript", "Backend", "Golang", "Cloud", "Front End"
};
```

File was not found

```
2023-12-04T07:04:59.747-05:00  INFO 41512 --- [  restartedMain] c.a.S.S
File does not exist re-crawl data for default search terms
Bot loaded for first time
```

Re-crawl data for the default search terms

## 2. HTML Parser :

**Code Explanation :**

The scrapeJobData method is part of a web scraping application designed to extract job-related information from a webpage. It uses the Jsoup library to parse the HTML source of the page and retrieve details such as job ID, title, company name, location, website name, salary, and job description. The method applies regular expressions to extract numerical salary information and converts it into minimum and maximum salary values. Various regular expression patterns are used to handle different salary formats, including daily, weekly, monthly, yearly, and hourly rates. The extracted data is then used to create a Job object, which encapsulates the job-related information, and this object is returned by the method. Overall, the method provides a structured way to extract and organize job data from the HTML source of a job posting webpage.

**Output 1 :**

The jobs with the below jobIds got scrapped

```
Job with id - 516430 scraped
Job with id - 514667 scraped
Job with id - 517540 scraped
Job with id - 517538 scraped
Job with id - 515328 scraped
Job with id - 504129 scraped
Job with id - 512297 scraped
Job with id - 508757 scraped
Job with id - 507482 scraped
Job with id - 506395 scraped
Job with id - 503595 scraped
Job with id - 502061 scraped
Size after inserting jobs :81
```

## 3. Data validation using regular expressions :

**Code Explanation :**

The code implements data validation for job-related information extracted during web scraping. It defines a set of regular expressions for validating different fields, such as job title, company name, website link, salary, location, and job description. The validateField method is a generic function that checks whether a given field matches the specified regular expression. The validateDataForOneObject method utilizes this generic validator for each field of a Job object, setting validation results for individual fields and an overall validation status. If any field fails validation, the allValid flag is set to false. During job data processing, only objects with valid information, as determined by the data validation process, are added to the jobsCollection. This

15

ensures that only well-formed and compliant job data is included in the final collection, enhancing the overall data quality.

**Output 1 :**

Validate data and only add to the database.json if the data is valid

```
Is Data Valid : true
Is Data Valid : true
Is Data Valid : true
Is Data Valid : true
Is Data Valid : true
Is Data Valid : false
Is Data Valid : true
Is Data Valid : true
Is Data Valid : true
Is Data Valid : false
Is Data Valid : true
```

## 4. Inverted Indexing and Frequency count :

**Code Explanation :**

```java
public TrieNode() {
    this.children = new HashMap<>();
    this.jobIds = new HashSet<>();
    this.wordFrequency = new TreeMap<>();
    this.isEndOfWord = false;
}
```

The code implements inverted indexing for efficient search operations in a large dataset of job descriptions. In the initialization of the trie (TrieDS), each job's description is processed to extract relevant tokens, excluding common English stop words. These processed tokens are then inserted into the trie data structure using the insertIntoTrie method. The trie enables quick retrieval of job IDs associated with specific words, forming an inverted index. This facilitates the search process by providing immediate access to job entries containing the queried words. The TrieNode class represents each node in the trie and includes a set of job IDs, establishing a direct link between words and their occurrences in job descriptions.

For frequency counting, the code utilizes the TrieNode's wordFrequency TreeMap within the insertIntoTrie method. This TreeMap stores the frequencies of each word within the job descriptions, organized by job ID. As words are inserted into the trie, their frequencies are updated in the corresponding TrieNode. This frequency count mechanism allows for efficient analysis of word occurrences across the dataset. During subsequent searches or data analysis, the Trie structure provides immediate access to both the job IDs containing specific words and the frequencies of those words within each job description. Overall, the combination of inverted indexing and frequency counting in the trie enhances the speed and efficiency of search operations in large-scale job datasets.

**Output 1 :**

The trie gets initialized with database.json every time the application starts



The words from database.json is used to initialize the trie



**Output 2 :**

The trie gets reinitialized whenever we scrape the data from the ui

Size after inserting jobs :296
RemoteOk Crawling Ended
Job data loaded from json
Trie init with jobs

## 5. Page Ranking :

### Code Explanation :

The PageRanking class implements page ranking by searching inverted indexed data stored in a Trie data structure. It takes an array of search terms, retrieves their frequencies from the Trie, and ranks jobs based on the cumulative word frequencies of the search terms within each job's description. For each term, it iterates over the corresponding Trie node to obtain job IDs and their word frequencies. Then, for each job ID, it retrieves the corresponding job details from the stored job data, updating the job's accumulated word frequency and word field. The jobs are then inserted into a SortedArray, sorted based on their word frequencies. If a job already has a non-zero frequency, it is temporarily removed from the sorted structure for updating before being reinserted. This is done because if we don't remove the word it will create duplication of jobs in the response array if we search for multiple search terms. The resulting SortedArray contains jobs ranked according to the cumulative frequency of the search terms within their descriptions, providing a page ranking reflecting the relevance of jobs to the search criteria.

### Output 1 :

Result on the ui for one search term

18

**Output 2 :**

Result on the ui for two search terms

Page Ranking

**Rank 1**
Senior Software Engineer Fullstack

Search term occurrences: 30
react   developer

Location: United States,Canada                                    **Salary:** 110000

| Description                                              ⌄ |

Remote Ok

**Rank 2**
React Native Developer

Search term occurrences: 27
react   developer

Location: Probably worldwide                                      **Salary:** 110000

| Description                                              ⌄ |

Remote Ok

**Rank 3**
React Developer

Search term occurrences: 18
react   developer

Location: Remote                                                 **Salary:** 124800

| Description                                              ⌄ |

SimplyHired

**Rank 4**
Senior Ruby Engineer Core

Search term occurrences: 12
react   developer

Location: Ottawa                                                 **Salary:** 110000

Description                                                ⌄

**Output 3 :**

Result on the ui for multiple search terms

## Salary Stratos

| Search Engine | Crawler |

| react software developer UI |          | Search |

Recent Searches With Frequency Count

| Frequency: 3 | Frequency: 2 | Frequency: 1 | Frequency: 1 |
| developer | react | UI | software |

20

## 6. Spell checking :

**Code Explanation :**

The SpellChecker class is designed to perform spell checking by maintaining a dictionary of valid words and employing the EditDistanceAlgo for suggesting similar words in case of misspelled input. The initialization process, handled by the initializeSpellChecker method, reads words from an external dictionary file (dictionaryOfWords.txt) and populates both a Trie data structure (TrieDS) and a frequency map (dict). The dict also contains all the words and the frequencies of the word from the database.json. The TrieDS efficiently stores words for quick retrieval and supports autocomplete features.

In the suggestSimilarWord method, given an input word, the code first checks if it's empty or belongs to a predefined list of invalid words. If the word is valid, it searches for an exact match in the Trie. If no match is found, it utilizes the EditDistanceAlgo to calculate the edit distance between the input word and all words in the dictionary. The results are organized in a nested TreeMap structure, where the outer TreeMap sorts suggestions based on edit distance, and the inner TreeMap further organizes them by word frequency. This prioritization ensures that suggestions are not only linguistically similar but also ranked by their prevalence in the dictionary, contributing to more accurate and contextually relevant spell-checking suggestions.
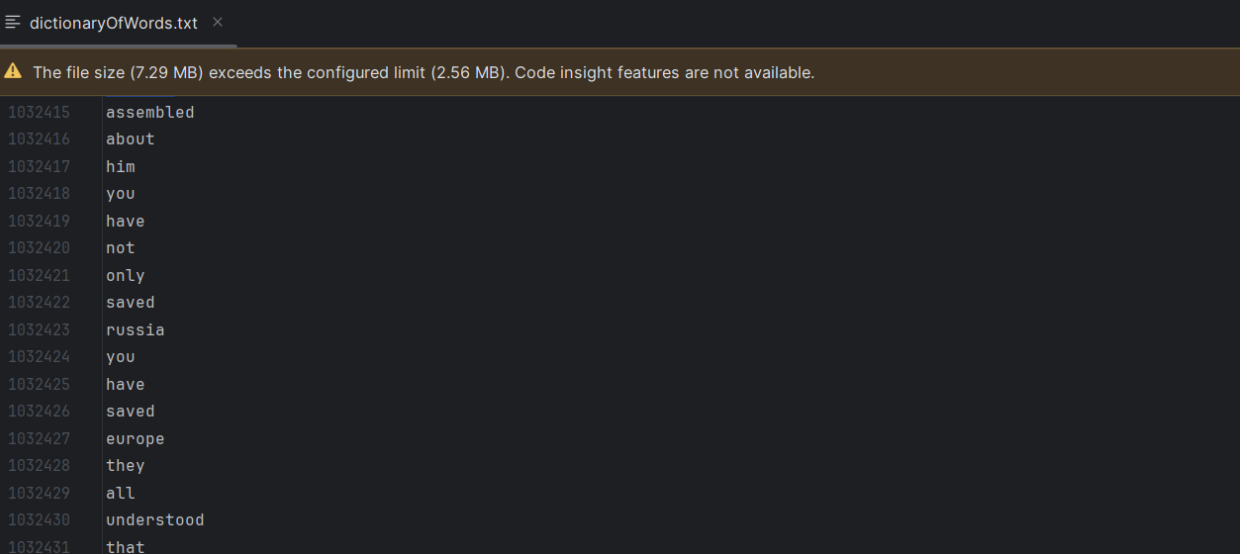
21

In summary, the spell-checking functionality leverages a Trie data structure, an external dictionary, and the EditDistanceAlgo to offer efficient and accurate suggestions for potential corrections based on word similarity and frequency. The data structures and algorithms employed contribute to a robust spell-checking mechanism within the provided code.

**Output 1 :**

When the program starts the spell checker dictionary is initialized.

```
Job data loaded from json
Trie init with jobs
Spell Checker Initialized
```

Data gets initialized from the dictionaryOfWords.txt

```
≡ dictionaryOfWords.txt ×

⚠ The file size (7.29 MB) exceeds the configured limit (2.56 MB). Code insight features are not available.

1032415    assembled
1032416    about
1032417    him
1032418    you
1032419    have
1032420    not
1032421    only
1032422    saved
1032423    russia
1032424    you
1032425    have
1032426    saved
1032427    europe
1032428    they
1032429    all
1032430    understood
1032431    that
```

**Output 2 :**

When the Scrapper is called, the spell checker dictionary is initialized.

```
RemoteOk Crawling Ended
Job data loaded from json
Trie init with jobs
Spell Checker Initialized
```

**Output 3 :**

Spell checking for one search term when the user hits enter. The search term native is correct.

It will search for native



## Output 4 :

Spell checking for two search terms when the user hits enter. The second search term nativp is not correct but the closest word to this is native so it suggested that based on edit distance

Search Engine   Crawler

react nativp                                                              [ Search ]

**Recent Searches With Frequency Count**

| Frequency: 10 | Frequency: 7 | Frequency: 4 | Frequency: 2 | Frequency: 1 | Frequency: 1 | Frequency: 1 | Frequency: 1 | Frequency: 1 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| react | operations | developer | native | I | System | am | UI | software |

**Word Suggestions**

react :

| Frequency: 217 | Frequency: 96 | Frequency: 24 | Frequency: 10 | Frequency: 10 |
| --- | --- | --- | --- | --- |
| react | reaction | reactjs | reactionary | reactive |

**Spell Checker**

react:
Search term is correct

nativp:

| EditCost: 1 | EditCost: 2 | EditCost: 3 | EditCost: 4 | EditCost: 5 | EditCost: 6 |
| --- | --- | --- | --- | --- | --- |
| Frequency: 117 | Frequency: 339 | Frequency: 354 | Frequency: 13587 | Frequency: 160060 | Frequency: 80060 |
| native | nation | nature | at | the | of |

It will only search for react in this case as the second term is wrong.

**Page Ranking**

**Rank: 1**
**React Native Developer**                    Search term occurrences: 21
                                              react

**Location:** Probably worldwide                          **Salary:** 110000

Description ⌄

Remote Ok

**Rank: 2**
**Senior Frontend Engineer**                  Search term occurrences: 12
                                              react

**Location:** Worldwide                                    **Salary:** 105000

Description ⌄

**Output 5 :**

Spell checking for multiple search terms when the user hits enter.

If we observe below we entered 3 words "react developer tedting". We can see that react and operations are correct words but tesdting is a wrong word. It also gave the closest word to tesdtion i.e. testing based on Edit Distance.

The search page ranking is done only for two words in this case "react developer"



## 7. Word Completion :

### Code Explanation :

The WordCompletion class provides functionality to generate word suggestions based on validated search terms using a Trie data structure. The getWordSuggestions method takes a list of validated search terms, a Trie structure (JobDataTrie), and the desired count of suggestions. It iterates through each search term, retrieves word suggestions from the Trie using the searchInTrieWithPrefix method, and constructs a response containing the original search term and a SortedArray of suggested words along with their frequencies. The searchInTrieWithPrefix method traverses the Trie to find words with the specified prefix, and the collectWords recursive function gathers word frequencies from the Trie nodes, constructing a SortedArray of WordFrequency objects sorted by frequency. The resulting WordSuggestionResponse

encapsulates the original search terms, whether the response is valid, and a list of suggested words with their frequencies.

**Output 1 :**
Word Suggestion for one search term when the user enters a partial word. It's not necessary click on enter to see the word suggestions



**Output 2 :**
Word Suggestion for two search terms when the user enters a partial word



**Output 3 :**
Word Suggestion for multiple search terms when the user enters a partial word

26

## 8. Search frequency :

**Code Explanation :**

The SearchFrequency class manages search term frequencies using an LRUCache, maintaining a limited size to track the most recent and frequently used search terms. The LRUCache, implemented as a custom class extending LinkedHashMap, enforces a Least Recently Used behavior. It stores search terms and their frequencies, automatically removing the least recently used entry when the cache exceeds the specified size. The displaySearchFrequencies method retrieves the last 10 search term frequencies in descending order using a custom SortedArray data structure, providing a dynamically updated view of recent search term activities. The updateSearchFrequency method ensures accurate tracking of user interactions by updating the search term frequencies based on an array of search terms. Overall, this component efficiently combines LRUCache and sorted arrays to offer insights into recent user search patterns.

LRU Cache size is 50

The return array size to ui is 10

If the Cache size becomes bigger than 50 it will remove the least recently used word from the LRU Cache

**Output 1 :**

First word entered into the LRU Cache "react"

27

### Output 2 :

Second word entered into the LRU cache "developer" which is searched 3 times



### Output 3 :

After inserting 10 words we see that "react" is not shown in the ui anymore because it is not used for a while, but it still remains in the cache till the cache size becomes 50, when this happens if "react" is the least recently used word it will get removed from the cache.



## 9. Finding patterns using regular expressions :

### Code Explanation :

The scrapeWebPage method in classes such as SimplyHiredScrapper, RemoteOk, and GlassDoorScrapper utilizes regular expressions (Regex) to identify various salary formats and convert them into a consistent yearly format. The code initializes multiple regular expression

patterns (regexYearlyWithK, regexYearlyFrom, etc.) designed to match different salary structures, such as hourly, weekly, monthly, or yearly, with or without the use of 'K' (thousands) and other variations.

After extracting the raw salary string from the web page, the code checks which regex pattern matches the format and proceeds with appropriate processing. For instance, it replaces the matched pattern with specific placeholders and then uses replaceAll to extract numeric values. It further adjusts these values based on the specific time unit (hours, weeks, etc.) to convert them into a standardized yearly representation. The resulting minimum and maximum salary values are then utilized in the application for analysis and comparison. This approach ensures accurate extraction and uniform representation of salary data, providing consistency in further processing and analysis

**Output 1 :**

If salary in Simply Hired is anything other than the below it will print the failure message

```
String regexYearlyWithK = "\\$([\\d.]+)K - \\$([\\d.]+)K a year";
String regexYearlyFrom = "\\$([\\d,]+) a year";
String regexWeeklyWithoutK = "\\$([\\d,]+) - \\$([\\d,]+) a week";
String regexMonthlyWithoutK = "\\$([\\d,]+) - \\$([\\d,]+) a month";
String regexYearlyWithoutK = "\\$([\\d,]+) - \\$([\\d,]+) a year";
String regexHourly = "\\$([\\d.]+) - \\$([\\d.]+) an hour";
String regexHourlyNonDecimal = "\\$([\\d]+) - \\$([\\d]+) an hour";
String regexHourlyFrom = "\\$([\\d.]+) an hour";
String regexDaily = "\\$([\\d.]+) a day";
String regexWeekMax = "\\$([\\d,]+) a week";
String regexYearFrom = "\\$([\\d,]+) a year";
```

## 10. Compare sorting algorithms for Page ranking :

**Code Explanation :**

The CompareRunTimesData class is designed to store and manage runtimes for different sorting algorithms, particularly Merge Sort, Binary Search, and QuickSort, within the context of inverted indexing and page ranking. Each sorting algorithm's runtime is tracked through dedicated variables (mergeSort, binarySearch, quickSort) along with respective getter and setter methods. This implementation is intended to facilitate runtime comparison between these algorithms when applied to the specific task of sorting jobs retrieved from the inverted index. By setting and retrieving the runtimes for each sorting operation, this class provides a streamlined way to evaluate and contrast the efficiency of these sorting methods in the context of page ranking tasks.

In the context of maintaining a sorted array for inverted indexing and page ranking, binary search holds advantages over Merge Sort and QuickSort. Binary search excels in scenarios where the data is already partially sorted, which aligns with the situation encountered when maintaining a sorted array. Given its time complexity of O(log n) for finding the correct insertion position in a

sorted array, binary search proves notably faster than both QuickSort (O(n log n)) and MergeSort (O(n log n)) in this specific context. Binary search's efficiency in locating insertion positions within a sorted array directly complements the task of maintaining order in the context of inverted indexing, making it a more suitable choice in scenarios where incremental sorting is required while minimizing time complexity.

**Output 1 :**



**Output 2 :**



## 10. Recommend top paid jobs based on search results :

### Code Explanation :

The searchInvertedIndexedDataBySalary method in the PageRanking.java class efficiently ranks jobs based on a combination of their maximum salary and word frequency. Leveraging a Trie data structure (TrieDS) for inverted indexing of job-related terms, the function iterates through the provided search terms, retrieves relevant information from the trie, and calculates the cost for each job. The cost is determined by adding the maximum salary and the word frequency associated with each job. Jobs are then inserted into a SortedArray based on their calculated cost, ensuring that the array remains sorted in descending order of cost. This approach enables the

quick identification of top-ranking jobs, and in cases where multiple jobs have the same salary, they are further sorted based on their word frequency. The result is an efficient mechanism for users to explore and discover jobs that align with their preferences, particularly focusing on the combined factors of salary and word frequency.

**Output 1 :**

Top paid jobs for "React"



**Output 2 :**

Top paid jobs for "React Native Developer"

Top Ranking Jobs by Salary

**Rank: 5**
Associate Frontend Engineer

Search term occurrences: 2
React

**Location:** Probably worldwide                                    **Salary:** 125000

Description ⌄

Remote Ok

**Rank: 6**
Primum Full Stack Developer

Search term occurrences: 5
React  Native  Developer

**Location:** United States                                         **Salary:** 120000

Description ⌄

Remote Ok

**Rank: 7**
Senior Software Engineer

Search term occurrences: 2
Developer

**Location:** Probably worldwide                                    **Salary:** 120000

Description ⌄

Remote Ok

## Output 3 :

We also show the description of the job along with the location, title, salary etc. When we click on the description it shows the job description.

Top Ranking Jobs by Salary

Remote Ok

**Rank: 49**
Full Stack PHP Web Developer

Search term occurrences: 4
React  Developer

**Location:** Remote                                               **Salary:** 100000

Description ← ⌄

SimplyHired

32

**Output 4 :**

The job website name is a button which has the link to the job. When clicked on the job website name it takes to the job website where the job is posted.