

Задание 3

ХУАН ЦЗИНЬЯНЬ

1. Write a program to find the longest word.

(1) Install Ray

(2) Replace the file path with the actual path

(3) Run the program

```
1  import ray
2  import re
3
4  # initialization Ray
5  ray.init()
6
7  # Define a Ray task to find the longest word of each line
8  1 usage
9  @ray.remote
10 def find_longest_word(line):
11     # Using regular expressions to extract words
12     words = re.findall(pattern=r'\b\w+\b', line)
13     # Find the longest word and return
14     if words:
15         return max(words, key=len)
16     return ""
17
18 # Read the file and process each line
19 1 usage
20 def process_file(filepath):
21     longest_words = []
22     with open(filepath, 'r', encoding='utf-8') as file:
23         # Submission of tasks and collection of results
24         future_results = [find_longest_word.remote(line) for line in file]
25         results = ray.get(future_results)
26         # Find the longest word from all the returned longest words
27         longest_word = max(results, key=len)
28         return longest_word
29
30 # Specify your file path
31 file_path = 'D:\Задание 3\wiki.txt'
32 # Process the file and print the longest word
33 longest_word = process_file(file_path)
34 print("The longest word is: ", longest_word)
```

2. Write a program to find the average length of words.

(1) Install Ray

(2) Replace the file path with the actual path

(3) Run the program

```

1  import ray
2  import os
3
4  # initialization Ray
5  ray.init()
6
7  # Define functions to process a single row of data
8  1 usage
9  @ray.remote
10 def process_line(line):
11     # Separate components in each row
12     parts = line.split('<tabs>')
13     if len(parts) < 3:
14         return [], 0 # If the format is incorrect, returns an empty list and a length of 0
15
16     # Extract the body part, split the words
17     content = parts[2].strip()
18     words = content.split()
19     lengths = [len(word) for word in words]
20     return lengths, len(words)
21
22 # Reading and processing files
23 1 usage
24 def process_file(file_path):
25     with open(file_path, 'r', encoding='utf-8') as file:
26         # Create a list of tasks, with each row of data corresponding to a task
27         futures = [process_line.remote(line) for line in file]
28
29     # Collection of processing results
30     results = ray.get(futures)
31
32     # Calculate total length and total number of words
33     total_length = sum(sum(lengths) for lengths, count in results)
34     total_count = sum(count for lengths, count in results)
35
36     # Calculate average word length
37     if total_count > 0:
38         average_length = total_length / total_count
39         return average_length
40     else:
41         return 0
42
43 # Replace the following path with your file path
44 file_path = 'D:\Задание 3\wiki.txt'
45 average_length = process_file(file_path)
46 print(f"Average word length: {average_length:.2f}")
47
48
49

```

3. Write a program to find the most commonly used words made up of the Latin alphabet.

(1) Install Ray

(2) Replace the file path with the actual path

(3) Run the program

```

1 import ray
2 import re
3 from collections import Counter
4
5 # initialization Ray
6 ray.init()
7
8
9 # Define a Ray's remote function to count words
10 usage
11 @ray.remote
12 def count_words(text):
13     # Extract only words containing Latin letters
14     words = re.findall(pattern=r'\b[a-zA-Z]+\b', text)
15     # Count the number of occurrences of each word
16     return Counter(words)
17
18 usage
19 def main():
20     file_path = 'D:\Задание 3\wiki.txt' # Replace with your file path
21     results = []
22
23     # Reads the file line by line and creates a Ray task for each line
24     with open(file_path, 'r', encoding='utf-8') as file:
25         for line in file:
26             if '<tabs>' in line:
27                 parts = line.split('<tabs>')
28                 if len(parts) >= 3:
29                     url, title, content = parts[0], parts[1], parts[2]
30                     # Send title and content to remote function
31                     results.append(count_words.remote(title + ' ' + content))
32
33     # Merge all results
34     total_counts = Counter()
35     for result in results:
36         total_counts.update(ray.get(result))
37
38     # Find the most commonly used words
39     most_common_word = total_counts.most_common(1)
40     print("The most commonly used words are: ", most_common_word)
41
42 if __name__ == "__main__":
43     main()
44

```

4. Find all words that begin with a capital letter more than half the time and occur more than 10 times.

(1) Install Ray

(2) Replace the file path with the actual path

(3) Run the program

```

1 import ray
2 import re
3
4 # initialization Ray
5 ray.init()
6
7 usage 新*
8 @ray.remote
9 def count_words(line):
10     # Use regular expressions to remove the tags, leaving only the body part
11     text = re.split(pattern=r'<tags>', line)[-1]
12     # Match all words using regular expressions
13     words = re.findall(pattern=r'\b[A-Z][a-z]*\b', text)
14     return words
15
16 usage 新*
17 def main():
18     # Replace with your file path
19     file_path = 'D:\Задание 3\wiki.txt'
20     # Read file
21     with open(file_path, 'r', encoding='utf-8') as file:
22         lines = file.readlines()
23
24     # Distributing tasks using Ray
25     futures = [count_words.remote(line) for line in lines]
26     results = ray.get(futures)
27
28     # Flat results list
29     all_words = [word for sublist in results for word in sublist]
30
31     # Counting word occurrences
32     word_count = {}
33     for word in all_words:
34         if word in word_count:
35             word_count[word] += 1
36         else:
37             word_count[word] = 1
38
39     # Filter words with at least 10 occurrences
40     frequent_words = {word: count for word, count in word_count.items() if count >= 10}
41
42     # output result
43     for word, count in frequent_words.items():
44         print(f"{word}: {count}")
45
46 if __name__ == '__main__':
47     main()

```

5. Write a program that uses statistics to identify pr., dr., forms of stable abbreviations.

(1) Install Ray

(2) Replace the file path with the actual path

(3) Run the program

```
1 import ray
2 import re
3
4 # 初始化 Ray
5 ray.init()
6
7
8 # 定义一个 Ray 任务来处理文件的每一行
9 usage 新 *
10 @ray.remote
11 def find_abbreviations(line):
12     # 使用正则表达式匹配缩写, 如 pr., dr. 等
13     abbreviations = re.findall(pattern=r'\b[a-zA-Z]+\.', line)
14     return abbreviations
15
16
17 usage 新 *
18 def process_file(file_path):
19     results = []
20     with open(file_path, 'r', encoding='utf-8') as file:
21         for line in file:
22             # 为每一行文本启动一个 Ray 任务
23             result = find_abbreviations.remote(line)
24             results.append(result)
25
26     # 收集所有结果
27     results = ray.get(results)
28
29     # 打印所有识别的缩写
30     for result in results:
31         if result:
32             print(result)
33
34 # 替换为你的文件路径
35 file_path = 'your_file_path'
36 process_file(file_path)
```

6. write a program to find out t.p., n.e., using statistical data stable abbreviations of the form.

(1) Install Ray

(2) Replace the file path with the actual path

(3) Run the program

```
1 import ray
2 import re
3
4 # initialization Ray
5 ray.init()
6
7 # Define a Ray remote function that finds abbreviations in text
8 usage 新*
9 @ray.remote
10 def find_abbreviations(text):
11     # Regular expressions match abbreviations like t.p., n.e.
12     abbreviations = re.findall(pattern=r'\b\w+\.\w+\. ', text)
13     return abbreviations
14
15 usage 新*
16 def main():
17     # Specify your file path
18     file_path = 'D:\Задание 3\wiki.txt'
19
20     # Create a list of all future objects.
21     futures = []
22
23     # Open and read files
24     with open(file_path, 'r', encoding='utf-8') as file:
25         for line in file:
26             # Assuming that the format of each line is: URL <tag> title <tag> full text
27             parts = line.split('<tags>')
28             if len(parts) >= 3:
29                 content = parts[2] # Access to the body of the text
30                 # Calls the remote function and adds the resulting object to the futures list
31                 future = find_abbreviations.remote(content)
32                 futures.append(future)
33
34     # Collect all results
35     results = ray.get(futures)
36
37     # Output all found abbreviations
38     all_abbreviations = set()
39     for result in results:
40         all_abbreviations.update(result)
41
42     print(all_abbreviations)
43
44 if __name__ == '__main__':
45     main()
```

7. Write a program that uses statistics to find names used in articles.

(1) Install Ray

(2) Replace the file path with the actual path

(3) Run the program

```
1 import ray
2 import re
3
4 # initialization Ray
5 ray.init()
6
7
8 !usage 新*
9 @ray.remote
10 def find_names(text):
11     # Use regular expressions to find all words that begin with a capital letter
12     names = re.findall(pattern: r'\b[A-Z][a-z]*\b', text)
13     return names
14
15 !usage 新*
16 def main():
17     # Replace the following path with your file path
18     file_path = 'D:\Задание 3\wiki.txt'
19     names_count = {}
20
21     with open(file_path, 'r', encoding='utf-8') as file:
22         results = []
23         for line in file:
24             url, title, content = line.split('<tabs>')
25             # Asynchronous processing of the body content of each line
26             result = find_names.remote(content.strip())
27             results.append(result)
28
29     # Collection of results
30     for result in results:
31         names = ray.get(result)
32         for name in names:
33             if name in names_count:
34                 names_count[name] += 1
35             else:
36                 names_count[name] = 1
37
38     # The name of the output statistic and its number of occurrences
39     for name, count in names_count.items():
40         print(f"{name}: {count}")
41
42 if __name__ == '__main__':
43     main()
44
```