

# 浙江大学



## Advanced Data Structures and Algorithm Analysis

Laboratory Projects 7

Map Reduce

Group 15

聂俊哲 张琦 张溢弛

Date: 2020-05-30

# Content

Chapter 1: Introduction to MapReduce

Chapter 2: Data Structure / Algorithm Specification

Chapter 3: Testing Results

Chapter 4: Analysis and Comments

Appendix: source code

References

Declaration

Signatures

# Chapter 1: Introduction to Map Reduce

## 1.1 Problem Description

### 1.1.1 Introduce

- In this project, we will introduce the framework of Map-Reduce and implement a MapReduce program to count the appearance of each word in a set of documents, which is a classic project in big-data called "word-count"

### 1.1.2 Input format

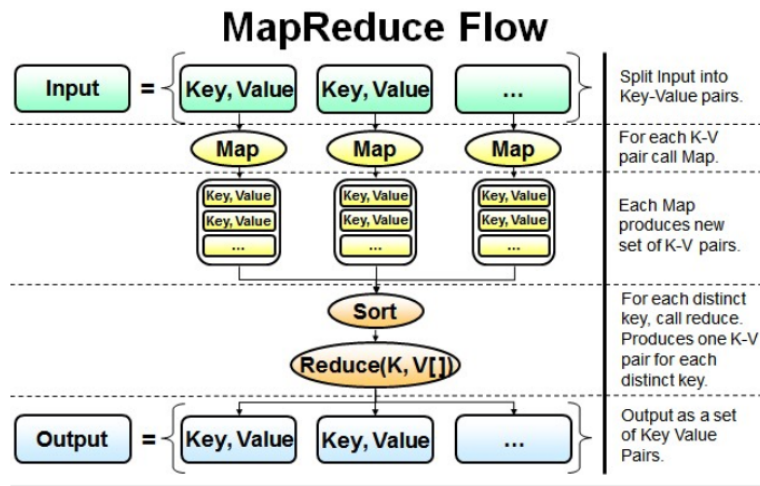
- A set of txt files in the folder "data", no need to input the data from keyboard and stdin.

### 1.1.3 Output format

- A map of words with their appearance time.
- The total number of the words.
- The running time of the program

## 1.2 Introduce to the framework of Map-Reduce

- MapReduce is a programming model and an associated implementation for processing and generating big data sets with a parallel, distributed algorithm on a cluster.
- A MapReduce Program includes a **map procedure** and a **reduce procedure**
  - map procedure performs filtering and sorting, which means making a pre-process of the big size data.
  - reduce procedure performs a summary operation which means counting the number of each word in this project.
- Another way to make a summary of MapReduce in 5 steps
  - **Prepare the Map() input** – the "MapReduce system" designates Map processors, assigns the input key  $K1$  that each processor would work on, and provides that processor with all the input data associated with that key.
  - **Run the user-provided Map() code** – Map() is run exactly once for each  $K1$  key, generating output organized by key  $K2$ .
  - **"Shuffle" the Map output to the Reduce processors** – the MapReduce system designates Reduce processors, assigns the  $K2$  key each processor should work on, and provides that processor with all the Map-generated data associated with that key.
  - **Run the user-provided Reduce() code** – Reduce() is run exactly once for each  $K2$  key produced by the Map step.
  - **Produce the final output** – the MapReduce system collects all the Reduce output, and sorts it by  $K2$  to produce the final outcome.
- The "MapReduce Framework" orchestrates the processing by marshalling the distributed servers, running the various tasks in parallel, managing all communications and data transfers between the various parts of the system, and providing for redundancy and fault tolerance.
- We can use such a graph to show the MapReduce algorithm.



- In our program code, we will use the **multiprocessing** library in Python3 to implement a parallel algorithm and a serial algorithm as comparasion.

## Chapter 2: Data Structure / Algorithm Specification

### 2.1 Data Structure and Function

- We implement the MapReduce algorithm in Python3 so that the data structure in the program is some simple linear list and **dictionaries**(which is similar to the STL map in C++)
- Some important data structures are
  - **pool:** it is used to make a parallel algorithm and create a process pool to do the word-count parallely. It is defined by the Python Parallel Compute Library **multiprocessing**

```
1 import multiprocessing
2 pool = multiprocessing.Pool(processes=len(file))
3 map_result = pool.map(word_count, file)
```

- **path\_set:** it is a list which stores all the txt files' path in the folder "data", in the program we use such an algorithm to get all the file paths

```
1 file_path = ".\\data\\"
2 path_set = [ ]
3 for file_name in os.listdir(file_path):
4     path_set.append(file_path+file_name)
5 return path_set
```

- **words:** it is a line of words using file reading operation in python and we use regular expression to get rid of **all the punctuation and Escape character**(like '\n') in the text. The algorithm of those operations is as following:

```
1 for line in fp:
2     # reduce the formats in the text.
3     words = re.sub(r"[%s]+" % string.punctuation, "", line)
```

```
1 words = re.sub("\n", " ", words).split(' ')
```

```
words = [x.lower() for x in words]
for single_word in words:
    if single_word in count.keys():
        count[single_word] += 1
    else:
        count[single_word] = 1
```

```
1
2 - result: the result after the **reduce operation** which merge the result of map operation.
```

```

3
4
5
6 ##### 2.2 Algorithm Specification
7
8 ##### 2.2.1 Overall Pseudocode
9
10 - The pseudocode of our program is
11
12 ```pseudocode
13 Algorithm: Map-Reduce
14 Input: none
15 Output: The map of word appearance
16
17 get all the file path in the folder "data";
18 # the main operation of map
19 for i in file_path:
20     single_map = word_count(i) #map the appearance of words
21     add single_map to the all_maps
22
23 # the main operation of reduce
24 result_map = { }
25 for i in all_maps:
26     for j in i.keys:
27         result_map[j]+=1
28 show the result of the program

```

## 2.2.2 Important functions

- There are some importance functions in the program

### 2.2.2.1 word\_count function

- We implement a parallel algorithm to count the appearance of each file.
- The function to count the appearance of the word in a single txt file, though its algorithm is simple using Python. The pseudocode of the function is:

```

1 Function: word_count
2 Input: the path of a single txt file
3 Return value: a map of word appearance
4
5 open the file using the path;
6 while read a line from the file:
7     simplify the words with reduction to punctuation and Escape character;
8     change all the words into lowercase;
9     split the word by ' ';
10    count the appearance of each words using a dictionary.
11 return the dictionary as result

```

### 2.2.2.2 Reduce function

- The reduce function will merge the single maps into a large one using the following algorithm
  - This is a serial version, we also have a parallel version which is a little different from the following algorithm

```

1 Function: Reduce
2 Input: the single maps collected in a list
3 Return value: Show the result of the total appearance
4
5 for i in map_list:
6     for j in i(i is a map now):
7         if j is not a digit and j is not a " ":
8             result_map[j]+=i[j]
9         else:
10            result_map[j]=i[j]
11 print the result_map in alphabetical order

```

## Chapter 3: Testing Results

- In this project, the test data is difficult to create by a test program. And this time, we created some text files or find some of them online to be the test data, which is in all sizes from KB to GB.

- What's more, we write a **test program based on the source code** to test the **running time** of our MapReduce system, we will draw a table and a statistical graph in the final of the part3.
- You can **find such test data in the folder "test data"**, but when running the code, the data will still be selected from the folder "data". In the **following parts we will not show the data in the report but just introduce them to avoid being a dirty report.**

### 3.1 Test data 1

- A single txt file with **only 1 word** in it
- It is a smallest condition of the project
- We do such a map-reduce operation 1000 times per test to calculate the average running time.
- The result is (for **100 times run**, just part of the final input)
  - In the serial program

```
Total Number of different words: 1
Total Number of different words: 1
Total Number of different words: 1
Total Number of different words: 1
Running time: 0.550546407699585
```

- In parallel program

```
Total Number of different words: 1
Running time: 1.062551498413086
```

### 3.2 Test data 2

- A single txt file with only Repeated words in a medium size.
- We copy the word "apple" for 200 times to be the test data.
- The result of **100 times of running** is:
  - In the serial program

```
Total Number of different words: 2
The number of total words: 39900
Running time: 0.07810568809509277
```

- In parallel program

```
apple 200
Total Number of different words: 2
Running time: 1.1764190196990967
```

### 3.3 Test data 3

- some txt files with KB-level sizes
- We do the operation 100 times per test to calculate the average running time.

类型	大小
文本文档	3 KB
文本文档	1 KB
文本文档	18 KB

- The running result of **100 times** of those files is
  - In serial program

```
Total Number of different words: 139
The number of total words: 400900
Running time: 0.9059920310974121
```

- In parallel program

```
Total Number of different words: 139
Running time: 1.2047772407531738
```

### 3.4 Test data 4

- Some txt files with MB-level sizes but smaller than test data 5

类型	大小
文本文档	2 KB
文本文档	62 KB
文本文档	18 KB

- We do the operation 100 times per test to caculate the average running time.
  - In the serial program

```
Total Number of different words: 2150
The number of total words: 1591000
Running time: 3.124260902404785
```

- In the parallel program

```
Total Number of different words: 2150
Running time: 3.4205074310302734
```

### 3.5 Test data 5

- some txt files with MB-level sizes

文本文档	9,583 KB
文本文档	28,749 KB
文本文档	28,722 KB

- We do the operation 10 times per test to caculate the average running time.
  - In the serial program

```
Total Number of different words: 2099
The number of total words: 132374610
Running time: 82.8415732383728
```

- In parallel program

```
Total Number of different words: 2099
Running time: 42.72328853607178
```

### 3.6 Test data 6

- large amount of files with MB-level sizes(25 files in total)

大小
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB
9,583 KB

- We run this program in **one time**, the result is
  - In the serial program

```
Total Number of different words: 2097
The number of total words: 47296200
Running time: 29.439714193344116
```

- In parallel program

```
Total Number of different words: 2097
Running time: 13.217740297317505
```

### 3.7 Test data 7

- A large size file with nealy 112MB data (I tried to run GB size data but my computer is out of work)
  - In the serial program



```
Total Number of different words: 2097
The number of total words: 22702192
Running time: 14.038269996643066
```

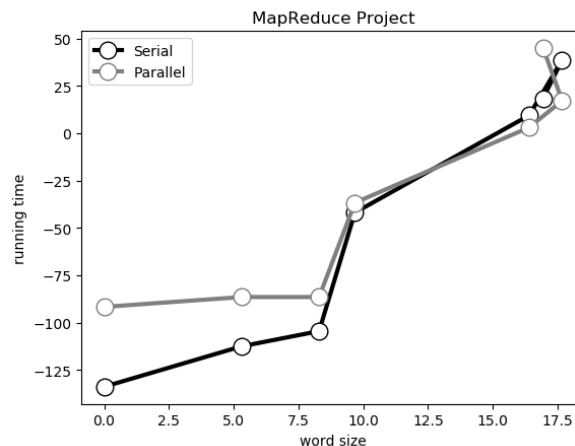
- In parallel program

```
Total Number of different words: 2097
Running time: 35.09289836883545
```

### 3.7 The static data

	Total Word Size	File size	running time(per time)(s) serial	running time(per time)(s) parallel
1	1	1	0.006	0.011
2	200	1	0.008	0.012
3	4000	3	0.009	0.012
4	15910	3	0.031	0.036
4	13237461	3	8.284	4.272
5	47296200	25	29.440	13.218
6	22702192	1	14.038	35.0928

- We could draw a static graph to show the result more directly
  - To make the graph a more contrast one, we use the **matplotlib** in python and do some **pre-process** to the data that is replace the origin data with its log and other operations, you can see the specific operation in our code **draw.py**
  - The analysis of the project is in the chapter 4



## Chapter 4: Analysis and Comments

### 4.1 Analysis of Time and Space complexity

- It's a special project which is not a traditional algorithm or data structure. And it is hard to make a conclusion of its time and space complexity for the reason that the main cost depends on the total size of the words in several files.
- Actually, the topic request is to analysis the test results, which will be detailed in 4.2

### 4.2 Analysis of the test results

- According to the results of testing, we could have the following conclusions

- In the serial algorithm of MapReduce, the time cost has a **linear relationship** with the size of words (We assume that the size of words is  $M$ ). We think the reason is that the serial algorithm will traverse all the words once and calculate the sum of their appearance, so the time complexity is  $T = O(M)$  in the traditional algorithm. And the number of files doesn't make a great difference in the time cost.
- In the parallel algorithm of MapReduce, we could make a conclusion that both size of the word and file impact greatly on the time cost because in the program we will create a new process for each file to count the word appearance and finally merge them one by one. So the time complexity is  $T = f(M, N)$  where  $N$  is the size of files, but it's still  $O(M)$  for a certain  $N$ .
- When the test data is in a small size, the serial algorithm runs faster than the parallel algorithm. However, when the test data is in a very large size with only one file, the serial algorithm still runs faster. Only **when the size of words and files are both in a large size, the parallel algorithm will run faster than the serial algorithm**. We could infer that it's because the time cost of creating a new process for a new file to count the word appearance is large compared to the total running time. So when the size of files is just one, the parallel algorithm just waste more time than serial algorithm. But a parallel algorithm can make operations on the files together by making full use of the CPU and memory. So when the size of words and files are both suitably large, the parallel algorithm will do a good job.

## Appendix

### Source Code

- Source code (Parallel version)

```

1  import re
2  import os
3  import time
4  import string
5  import multiprocessing
6
7  number_of_word = 0
8
9
10 # count the appearance of words in the text file
11 def word_count(file_path):
12     """
13     :param file_path: the path of the text file
14     :return: a dictionary of words' appearance
15     """
16     fp = open(file_path)
17     count = {}
18     global number_of_word
19     number = 0
20     for line in fp:
21         # reduce the formats in the text.
22         words = re.sub(r"[%s]+" % string.punctuation, "", line)
23         words = re.sub("\n", " ", words).split(' ')
24         words = [x.lower() for x in words]
25         # words = re.sub(r"[%s]+" % string.punctuation, "", line).split(' ')
26         for single_word in words:
27             number += 1
28             if single_word in count.keys():
29                 count[single_word] += 1
30             else:
31                 count[single_word] = 1
32         # print(count)
33         number_of_word += number
34         # print(number)
35         # print(number_of_word)
36     return count
37
38
39 # find all the file path
40 def get_file_path():
41     """
42     :no param
43     :return: a list of all the file paths
44     """

```

```

45     file_path = ".\\test data\\test4\\"
46     path_set = []
47     for file_name in os.listdir(file_path):
48         path_set.append(file_path + file_name)
49     return path_set
50
51
52 # show the result of the map-reduce algorithm
53 def show_result(result_map):
54     """
55
56     :param result_map: the result map that will be shown to us
57     :return: none
58     """
59     # result_map.sort()
60     count = 0
61     for i in sorted(result_map.keys()):
62         if count != 0:
63             print(i, result_map[i])
64             count += 1
65     print("Total Number of different words:", count)
66     # global number_of_word
67     # print("The number of total words:", number_of_word)
68
69
70 # the main function of map operation
71 def Map_Operation():
72     file = get_file_path()
73     # print(file)
74     # parallel
75     pool = multiprocessing.Pool(processes=len(file))
76     map_result = pool.map(word_count, file)
77     # print(map_result)
78     Reduce_Operation(map_result)
79     # print(count_dict)
80
81
82 # the main function of reduce operation
83 def Reduce_Operation(count_dict):
84     # print(count_dict)
85     # print(len(count_dict))
86     result = {}
87     # i is a dict now
88     for i in count_dict:
89         # print(i)
90         for j in i.keys():
91             # print(j)
92             if str(j).isdigit() == 0 and str(j) != " ":
93                 if str(j) in result.keys():
94                     # print(1)
95                     result[str(j)] += int(i[j])
96             else:
97                 # print(0)
98                 result[str(j)] = int(i[j])
99     # print(result)
100    show_result(result)
101
102
103 if __name__ == '__main__':
104     # test the running time of the program
105     start = time.time()
106     Map_Operation()
107     # print(count_dict)
108     # Reduce_Operation()
109     end = time.time()
110
111     print("Running time:", (end - start))
112

```

- Source code(Serial version)

```

1  import re
2  import os
3  import time
4  import string
5
6
7
8  count_dict = []
9  number_of_word = 0
10
11
12 # count the appearance of words in the text file
13 def word_count(file_path):
14     """

```

```

15 :param file_path: the path of the text file
16 :return: a dictionary of words' appearance
17 """
18 fp = open(file_path)
19 count = {}
20 global number_of_word
21 for line in fp:
22     # reduce the formats in the text.
23     words = re.sub(r"%s+" % string.punctuation, "", line)
24     words = re.sub("\n", " ", words).split(' ')
25     words = [x.lower() for x in words]
26     # words = re.sub(r"%s+" % string.punctuation, "", line).split(' ')
27     for single_word in words:
28         number_of_word += 1
29         if single_word in count.keys():
30             count[single_word] += 1
31         else:
32             count[single_word] = 1
33     # print(count)
34     global count_dict
35     count_dict.append(count)
36     # print(count)
37
38
39 # find all the file path
40 def get_file_path():
41     """
42     :no param
43     :return: a list of all the file paths
44     """
45     file_path = ".\\test data\\test4\\"
46     path_set = []
47     for file_name in os.listdir(file_path):
48         path_set.append(file_path+file_name)
49     return path_set
50
51
52 # show the result of the map-reduce algorithm
53 def show_result(result_map):
54     """
55
56     :param result_map: the result map that will be shown to us
57     :return: none
58     """
59     # result_map.sort()
60     count = 0
61     for i in sorted(result_map.keys()):
62         if count != 0:
63             print(i, result_map[i])
64         count += 1
65     print("Total Number of different words:", count)
66
67
68 # the main function of map operation
69 def Map_Operation():
70     file = get_file_path()
71     # print(file)
72     for i in file:
73         word_count(i)
74     # print(count_dict)
75
76
77 # the main function of reduce operation
78 def Reduce_Operation():
79     global count_dict
80     # print(count_dict)
81     # print(len(count_dict))
82     result = {}
83     # i is a dict now
84     for i in count_dict:
85         # print(i)
86         for j in i.keys():
87             # print(j)
88             if str(j).isdigit() == 0 and str(j) != " ":
89                 if str(j) in result.keys():
90                     # print(1)
91                     result[str(j)] += int(i[j])
92             else:
93                 # print(0)
94                 result[str(j)] = int(i[j])
95     # print(result)
96     show_result(result)
97
98
99 if __name__ == '__main__':
100     # test the running time of the program
101     start = time.time()

```

```
102     Map_Operation()  
103     # print(count_dict)  
104     Reduce_Operation()  
105     end = time.time()  
106     print("The number of total words:", number_of_word)  
107     print("Running time:", end-start)
```

- Test Program is adapted from the source code.

## References

List all the references here in the following format:

- [1] 《算法导论》第三版
- [2] PTA website <https://pintia.cn/>
- [3] Course Slides of ADS

## Author List

- Zhang Yichi 3180103772
- Zhang Qi 3180103162
- Nie Junzhe 3180103501

## Declaration

We hereby declare that all the work done in this project titled "Huffman Codes" is of our independent effort as a group.

## Signatures

夏俊哲

张琦

张溢邦