# LAB 3: MAPREDUCE

Q1: How many times does the word Sherlock appear in the file?

The propose of part 1 is to statistic how many times a word occurs in the text. When we are using Hadoop to do the distributed calculation, what we should consider are the MAP job, REDUCE job and a MAIN function which invokes the program and configure the nodes. MAP is responsible to spilt all the word in the text and emit each word with a key-value pair, like ("word", 1) and send word-iterationList like ("word", 1, 1, 1…) to REDUCE nodes. REDUCE is responsible for aggregating the iteration list for each word and generating a ("word", count) pair for each unique word. Notice that a word is only sent to a REDUCE node. So the result of the REDUCE is the word-wordCount pair. To find the times that "Sherlock" appears, we only need to check the file and find the count, which is 345.

In *TokenizerMapper.java*,

> *StringTokenizer itr = new StringTokenizer(value.toString(), "-- \t\n\r\f,.:;?'\"");*

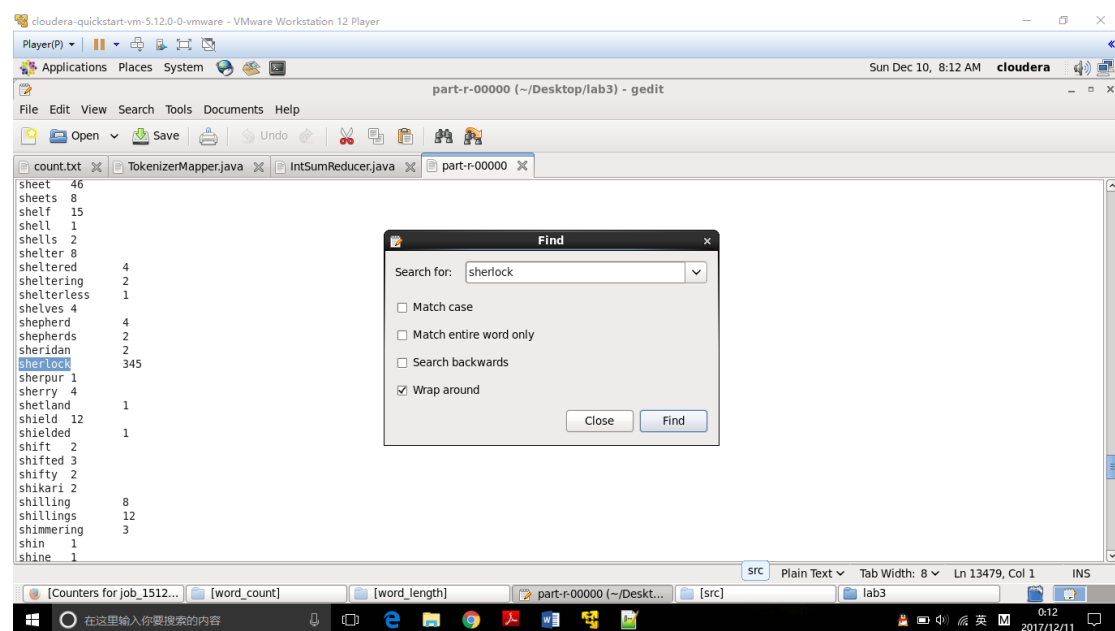means spilt the text into a list if it meets the following syntax: "-- \t\n\r\f,.:;?![]()'\"".

In *IntSumReducer.java*,

$$sum += value.get();$$

value.get() is to turn IntWritable type into int type; add the element of iterable list to sum.
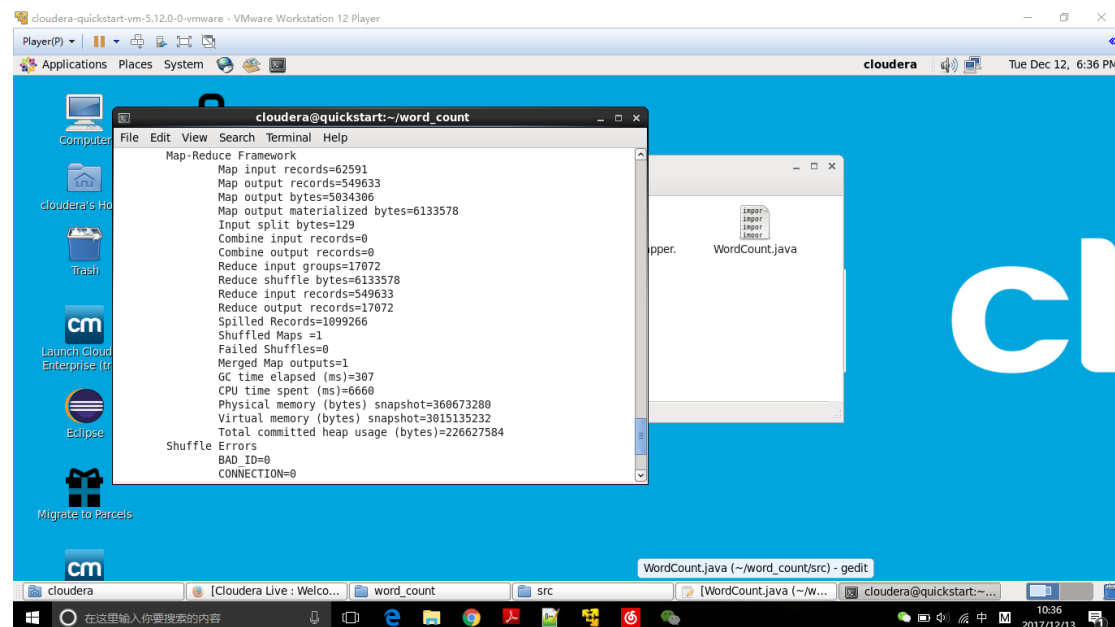
$$context.write(key, result);$$

emit the (word, count) pair.



Q2: The Hadoop job outputs a set of statistics after the completion. What was the number of key-value pairs generated by the Mapper, the number of unique inputs to the Reducer, and the final number of records emitted by the Reducer? Based on these values, answer the following two questions: How many words has the document in total? How many unique words does the document have?

According to our MAP method, which emits ("words", 1) for each word, the number of key-value

pairs generated by the Mapper is actually the total number of words in the text. From the screenshot, we can see the value is 549633. The number of unique inputs to the Reducer is the Reduce input group which is 17072. The final number of records emitted by the Reducer is 17072 since the reducer adds the occurrence times of each word, so it is also equal to the number of unique words. So base on the analysis above, the document has 549633 words in total. There are 17072 unique words in the document.



Q3: How many Map and Reduce tasks does your MapReduce job have? Can you explain the difference with the Sherlock job in Part 1?

Since we added *"job.setNumReduceTasks(2);"* in our "main class", which means we set two Reduce nodes to do the job. So we have 1 Map task and 2 Reduce tasks. Notice that the task in 2 Reduce nodes is executed parallel and each node is allocated about half of the total output of Map job. So compared with Sherlock job which only has 1 Reduce node, the Reduce job is faster in time in the case of the same Reduce input.

Q4: Can you see a clear pattern on how Hadoop partitions the keys among multiple reducers? Does it make it easier or harder the problem of manually retrieving information about a specific key? E.g. Answering questions such as "How many times the word Sherlock appears in the text?"

I can't see a clear pattern on how Hadoop partitions the keys. All pairs of a key are allocated to one reducer. It is the internal mechanism of Hadoop to handle which key is distributed to which Reduce node and it should be balance for each node.
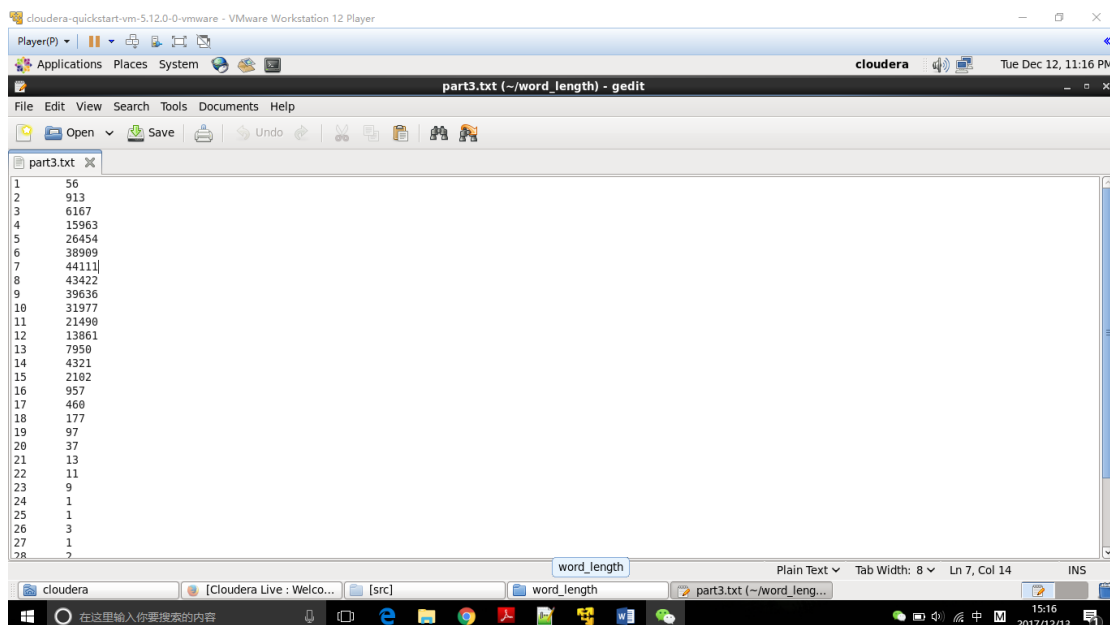
Since there are two Reduce nodes to handle the aggregation process and the Reduce tasks is parallel, I think it make it easier to handle the problem.

Sherlock appears 4 times in the text.

Q5: What is the most common word length for the onehundredM.txt dataset? Briefly explain the Mapper and Reducer code that you have implemented.

According to the screenshot below, we find that the most common word length for the onehundredM.txt dataset is 7.



The Mapper code is shown below:

The input of Map is the document: count.txt, which generated in part 2. So the input data will be *<Object, Text>*. Then we want to emit pair for each unique word, which is *<word_length, 1>*. So the output key and value are both *IntWritable*. The format of the file has each pair in a line and the key (unique word) and value (how many times the word occurs) is divided by "\t". After we scan the whole file, the value is the whole text. Then we spilt each line through *StringTokenizer* and for each line, we spilt it to a string array. In this case, components[0] represents the word, components[1] represents the times of the words occur. Then we acquire the length of the word and emit the pair (*word_length, 1*).

The Reduce code is shown below:



Since the type that Map emit is <IntWritable, IntWritable>, the input key should be IntWritable and the value is a IntWritable list. The output of Reduce is also pairs < IntWritable, IntWritable >. Since

Reduce method is called for each key, what we need to do is to add the 1 of the list and emit the pair (word_length, sum_of _uniqueWord_has_the_same_length).