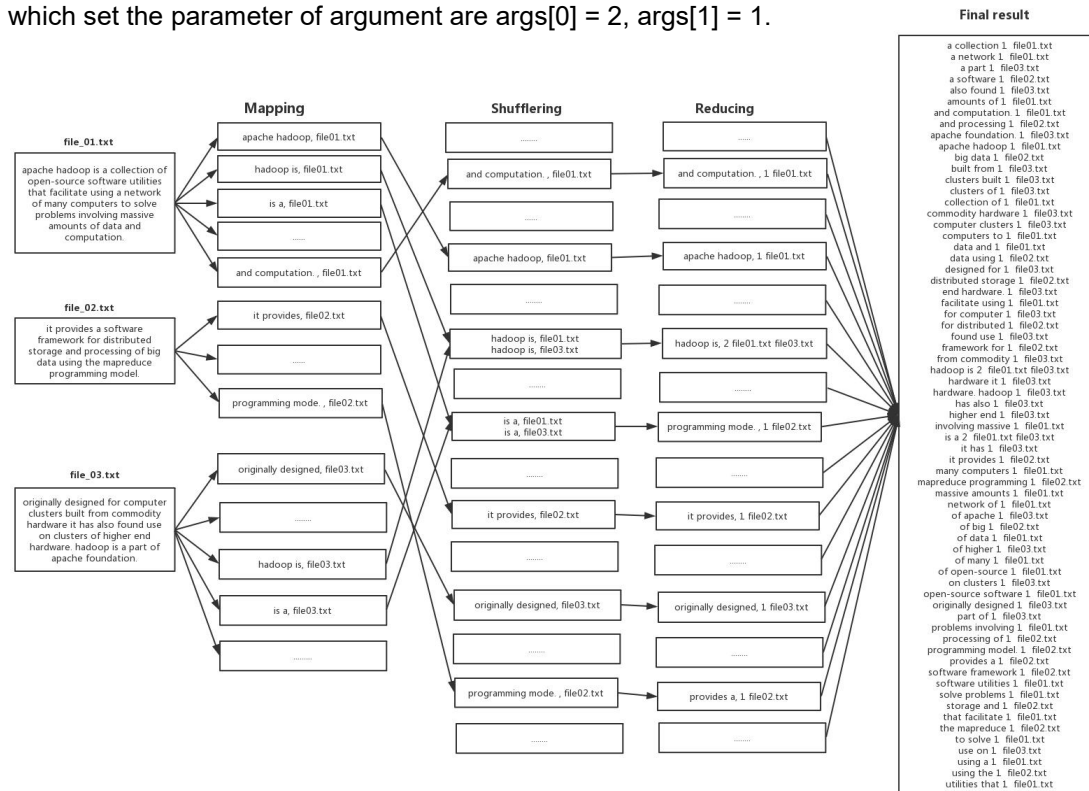# COMP9313 Assignment1 Report

## 1. Map and Reduce Process

As the graph show below, you can see the design process of map and reduce in assignment 1 which set the parameter of argument are args[0] = 2, args[1] = 1.



In Mapping function, I design the <Input key, Input value, output key, output value> is <Object, Text, Text, Text>. Firstly, Divide the contents of the file into words by space, and then use the iterator of the ngrams algorithm to sequentially obtain phrases consisting of n consecutive words. Next, write the key-value pairs(ngrams word, file name which word belong to) to context.

In Reducing function, The <Input key, Input value, output key, output value> is <Text, Text, Text, Text>. In the beginning, get the shufflering and sorting result <ngarms word, list of file name>. After that calculate the length of the list, we know the occurrences number of ngrams word. Next, I connect the length and the file name that appears in the file as the output value. The output key is ngrams word which is also input key. Finally, I write the output key-value pair into context.

## 2. Main(Driver)

In the main function, I creating a Configuration object and a Job object, assigning a job name for identification purposes. Then setting some attributes of job. After that, the HDFS input and output directory to be fetched from the command line, and two parameters for program. Finally, Submit the job to the cluster and wait for it to finish.