# Implementing MapReduce

### Siddhant Dilip Godshalwar

### February 25, 2023

## 1  Introduction

MapReduce is a programming model that is designed for processing and analyzing large datasets. It was originally introduced by Google in 2004 and has since become a widely used approach for large-scale data processing. MapReduce allows for efficient distributed processing of large datasets across a cluster of computers, enabling users to process data that would otherwise be too large to process on a single machine.

## 2  Implementation

The MapReduce model is a programming model that allows for distributed processing of large datasets across a cluster of computers. The MapReduce model consists of two primary phases: the Map phase and the Reduce phase.

### 2.1  Master Node

In a MapReduce framework, the MasterNode is the central coordinator that manages the entire processing of the data. It is responsible for dividing the data into smaller chunks and distributing these chunks to different worker nodes or mappers. The MasterNode also oversees the progress of the Map and Reduce tasks and monitors the health and status of the worker nodes.

Via TCP/IP sockets, typically, the MasterNode communicates with the worker nodes over a network. By keeping track of the health and condition of the worker nodes and reassigning work to other nodes in the event of failure, it also offers fault tolerance.

The MapReduce framework's MasterNode is a crucial part since it offers a centralized control method for handling massive volumes of data. The MasterNode may greatly increase the speed and effectiveness of data processing by dividing the task across several worker nodes and coordinating their efforts.

### 2.2  Data Partitioning

Data Partitioning is done by Master Node. The Master decides which mapper gets how many files and then assigns accordingly. The Data Partitioning in this implementation has three rules:-

- CASE 1 If the number of files to be Map Reduced is less than the number of mappers then each mapper gets 1 file and the remaining mappers get nothing.

- CASE 2 If the number of files is divisible by the number of mappers then each mapper has equal distribution of files

- CASE 3 If the number of files is greater than the number of mapper but is not divisible then all mappers get equal distribution except the last mapper which gets the remainder of files as well.
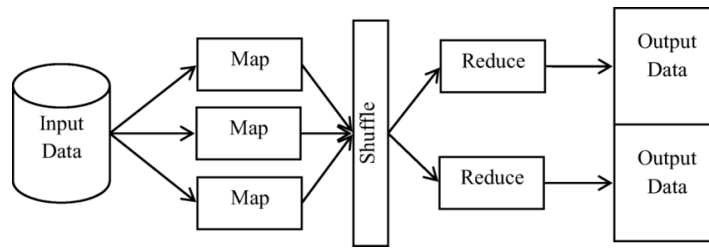
Figure 1: Architecture of MapReduce.

## 2.3  Mapper Phase

In the Map phase, the input data is processed and transformed into intermediate key-value pairs. The input data is typically a large dataset that is partitioned across the nodes in the cluster. Each node in the cluster processes a portion of the data independently using a user-defined mapping function. The mapping function takes in a single record of the input data and generates a set of intermediate key-value pairs.

The intermediate key-value pairs generated by the mapping function are then sorted, partitioned, and shuffled so that all the intermediate key-value pairs with the same key are sent to the same node for processing in the Reduce phase.

## 2.4  GroupBy Phase

In MapReduce, the GroupBy phase is the process of grouping together the intermediate key-value pairs produced by the mappers, based on their keys. The purpose of this phase is to ensure that all key-value pairs with the same key are processed together by the reducer, as the reducer receives input in the form of !key, list of values! pairs.

The GroupBy phase takes place after the Map phase and before the Reduce phase, and it involves sorting and shuffling the intermediate key-value pairs. Specifically, the intermediate key-value pairs generated by the mappers are first sorted by key. Then, pairs with the same key are grouped together and sent to the same reducer.

## 2.5  Reducer Phase

In the Reduce phase, the intermediate key-value pairs are combined and aggregated to produce the final output. The Reduce phase also consists of a user-defined function that is applied to each set of intermediate key-value pairs with the same key. This function takes in a key and a list of values and returns a set of output key-value pairs.

The output of the Reduce phase is typically written to a distributed file system, such as Hadoop Distributed File System (HDFS) or Amazon S3. The final output can then be used for further analysis or processing.

# 3  Workflow

## 3.1  Steps to run the files:

- 1 Insert all the files you want to MapReduce in the Input File directory in .txt format

- 2 Decide the number of Mappers you want and Reducers you want and change the values in the config.py file accordingly

- 3 Run the command MapReduce.py file

## 3.2  Output:

There are two types of output files generated:

- 1 ReducerInvertedOutput:- this contains the output in an inverted index format for that particular ReducerID

- 2 ReducerOutput:- this contains the output in word-count format for that particular ReducerID

# 4    Testing

The testing files in this implementation have been added to the 'inputFiles' directory available in the code. Please feel free to add any number of files you are interested in adding.

## 4.1    TestCase1:- Base Case

**File Count:-**1
**Mapper:-**1
**Reducer:-**1



Figure 2: TestCase 1:- Word Count



Figure 3: TestCase 1:- Inverted Index

## 4.2   TestCase2:- No. of Files = No. of Mappers = No. of Reducers

**File Count:-**6
**Mapper:-**6
**Reducer:**6



Figure 4: TestCase 2: Word Count



Figure 5: TestCase 2:- Inverted Index

## 4.3 TestCase3:- No. of Files is greater than No. of Mappers and is divisible

**File Count:-**6
**Mapper:-**3
**Reducer:**3



Figure 6: TestCase 3: Word Count



Figure 7: TestCase 3:- Inverted Index

## 4.4 TestCase4:- No. of Files is greater than No. of Mappers and is not divisible

**File Count:-**7
**Mapper:-**3
**Reducer:**4



Figure 8: TestCase 4: Word Count



Figure 9: TestCase 4:- Inverted Index

## 4.5 TestCase5:- No. of Files is less than No. of Mappers and No. of Reducers

**File Count:-**6
**Mapper:-**9
**Reducer:**9



Figure 10: TestCase 5: Word Count



Figure 11: TestCase 5:- Inverted Index

## 4.6 TestCase6:- No. of Reducers is less than No. of Files is less than No. of Mappers

**File Count:-**6
**Mapper:-**9
**Reducer:**4



Figure 12: TestCase 6: Word Count



Figure 13: TestCase 6:- Inverted Index

# 5 Limitations and Assumptions

## 5.1 Limitation

The Limitations of this implementation of MapReduce is:-

- 1 This doesn't account for the Fault Tolerance of any one of the mappers/reducers. So if a Mapper/Reducer fails, there will be some loss of data

- 2 This MapReducer also can only consume files that are in '.txt' format and are using 'UTF-8' codec. Other files are not parseable

## 5.2 Assumptions

The Assumptions of this implementation of MapReduce is:-

- 1 The communication between the Master and the Mappers and Reducers is secure.

- 2 The Masters and Reducers will work without failure.

- 3 The input files are all plain text files and use 'UTF-8' codecs.