**Department of Computer Science and Engineering (Data Science)**

**Subject: Big Data Engineering (DJ19DSL604)**
**AY: 2022-23**

**Experiment 1**

**(Hadoop Ecosystem)**

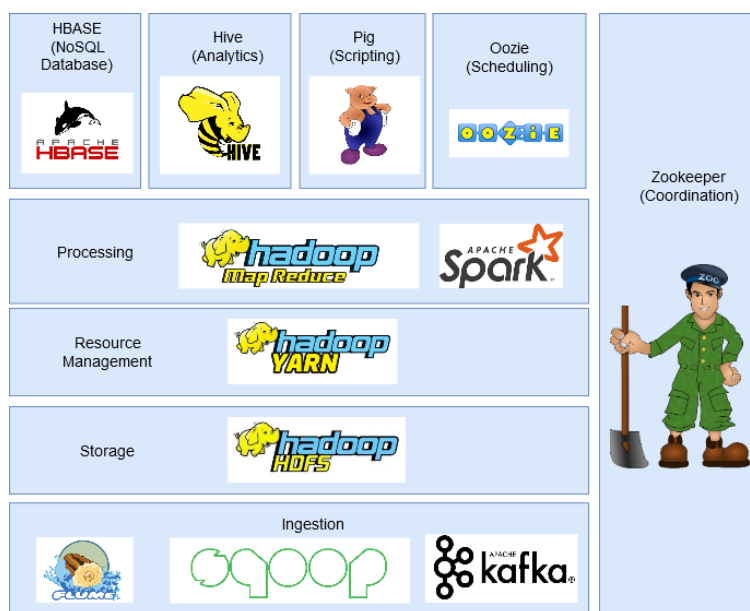Name : Sarvagya Singh
SAPID : 60009200030
BATCH : K1

**Aim:** Hands on experience on Hadoop Ecosystem and MapReduce.

**Theory:**

**1. Understanding inputs and outputs of MapReduce Concept of Hadoop:**

Apache Hadoop is an open-source framework based on Google's file system that can deal with big data in a distributed environment. This distributed environment is built up of a cluster of machines that work closely together to give an impression of a single working machine.

Components of Hadoop Ecosystem:



**2. The Map Tasks, grouping by Key, The Reduce Tasks, Combiners, Details of MapReduce Execution:**

To handle Big Data, Hadoop relies on the MapReduce algorithm introduced by Google and makes it easy to distribute a job and run it in parallel in a cluster. It essentially divides a single task into multiple tasks and processes them on different machines.

The MapReduce framework operates on key-value pairs, that is, the framework views the input to the job as a set of key-value pairs and produces a set of key-value pair as the output of the job,

**Department of Computer Science and Engineering (Data Science)**

conceivably of different types.

The key and value classes have to be serializable by the framework and hence, it is required to implement the Writable interface. Additionally, the key classes have to implement the WritableComparable interface to facilitate sorting by the framework.

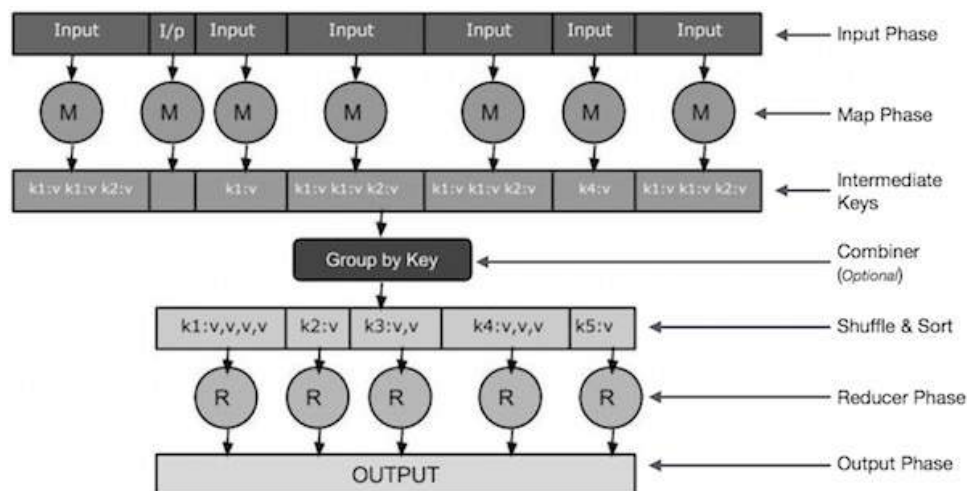Both the input and output format of a MapReduce job are in the form of key-value pairs −

**(Input) <k1, v1> -> map -> <k2, v2>-> reduce -> <k3, v3> (Output).**

|        | Input            | Output         |
|--------|------------------|----------------|
| Map    | <k1, v1>         | list (<k2, v2>) |
| Reduce | <k2, list(v2)>   | list (<k3, v3>) |

The MapReduce algorithm contains two important tasks, namely **Map** and **Reduce**.

- The Map task takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key-value pairs).

- The Reduce task takes the output from the Map as an input and combines those data tuples (key-value pairs) into a smaller set of tuples.

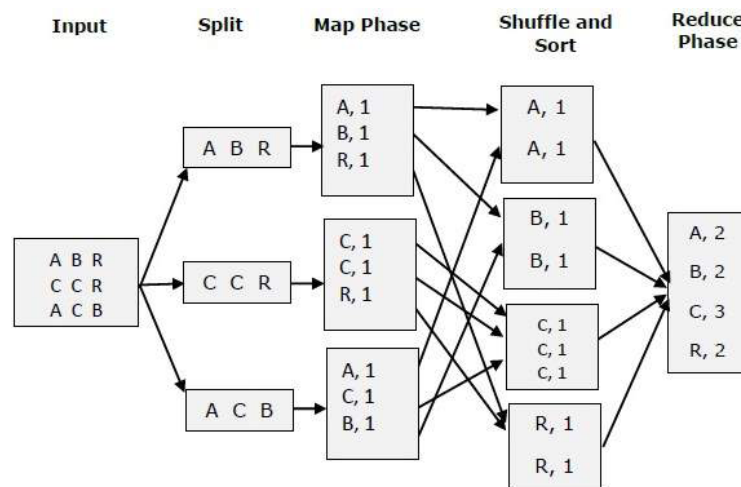The reduce task is always performed after the map job.



- **Input Phase** − Here we have a Record Reader that translates each record in an input file and sends the parsed data to the mapper in the form of key-value pairs.

- **Map** − Map is a user-defined function, which takes a series of key-value pairs and processes each one of them to generate zero or more key-value pairs.

- **Intermediate Keys** − They key-value pairs generated by the mapper are known as intermediate keys.

2

**Department of Computer Science and Engineering (Data Science)**

- **Combiner** − A combiner is a type of local Reducer that groups similar data from the map phase into identifiable sets. It takes the intermediate keys from the mapper as input and applies a user-defined code to aggregate the values in a small scope of one mapper. It is not a part of the main MapReduce algorithm; it is optional.

- **Shuffle and Sort** − The Reducer task starts with the Shuffle and Sort step. It downloads the grouped key-value pairs onto the local machine, where the Reducer is running. The individual key-value pairs are sorted by key into a larger data list. The data list groups the equivalent keys together so that their values can be iterated easily in the Reducer task.

- **Reducer** − The Reducer takes the grouped key-value paired data as input and runs a Reducer function on each one of them. Here, the data can be aggregated, filtered, and combined in a number of ways, and it requires a wide range of processing. Once the execution is over, it gives zero or more key-value pairs to the final step.

- **Output Phase** − In the output phase, we have an output formatter that translates the final key-value pairs from the Reducer function and writes them onto a file using a record writer.



Below is the simple structure of mapper, reducer and driver java files: (the highlighted text needs to be changed as per the problem statement)

**Mapper Code:**
```
// Importing all libraries
public class MapperClassName extends MapReduceBase implements Mapper<LongWritable,
Text, Text, IntWritable> {

        // Map function
        public void map(LongWritable key, Text value, OutputCollector<Text,
                        IntWritable> output, Reporter rep) throws IOException
```

3

```
        {

            String line = value.toString();

            // Write your Code here
            output.collect(new Text(finalvariablethatcontainsdata), new IntWritable(1));
}
```

**Reducer Code:**

```
// Importing libraries
public class ReducerClassName extends MapReduceBase implements Reducer<Text, IntWritable,
Text, IntWritable> {

        // Reduce function
        public void reduce(Text key, Iterator<IntWritable> value, OutputCollector<Text,
IntWritable> output, Reporter rep) throws IOException
        {

            // Write your code here

            output.collect(key, new IntWritable(finalvariablethatcontainsdata));
        }
}
```

**Driver Code:**

```
// Importing libraries

public class DriverClassName {
    public static void main(String[] args) {
        JobClient my_client = new JobClient();
        // Create a configuration object for the job
        JobConf job_conf = new JobConf(DriverClassName.class);

        // Set a name of the Job
        job_conf.setJobName("Nameofthejob");

        // Specify data type of output key and value
        job_conf.setOutputKeyClass(Text.class);
        job_conf.setOutputValueClass(IntWritable.class);

        // Specify names of Mapper and Reducer Class
        job_conf.setMapperClass(foldernamecontainingclassfiles.MapperClassName.class);
        job_conf.setReducerClass(foldernamecontainingclassfiles.ReducerClassName.class);

        // Specify formats of the data type of Input and output
```

**Department of Computer Science and Engineering (Data Science)**

```
job_conf.setInputFormat(TextInputFormat.class);
job_conf.setOutputFormat(TextOutputFormat.class);

// Set input and output directories using command line arguments,
//arg[0] = name of input directory on HDFS, and arg[1] =  name of output directory to be
created to store the output file.

FileInputFormat.setInputPaths(job_conf, new Path(args[0]));
FileOutputFormat.setOutputPath(job_conf, new Path(args[1]));

my_client.setConf(job_conf);
try {
  // Run the job
  JobClient.runJob(job_conf);
} catch (Exception e) {
  e.printStackTrace();
}
  }
}
```

## 3. Hadoop Distributed File System (HDFS) - Command line:

It is the storage component of Hadoop that stores data in the form of files.

Each file is divided into blocks of 128MB (configurable) and stores them on different machines in the cluster.

It has a master-slave architecture with two main components: Name Node and Data Node.

- **Name node** is the master node and there is only one per cluster. Its task is to know where each block belonging to a file is lying in the cluster

- **Data node** is the slave node that stores the blocks of data and there are more than one per cluster. Its task is to retrieve the data as and when required. It keeps in constant touch with the Name node through heartbeats

## 4. Overview of resource management – YARN:

YARN or Yet Another Resource Negotiator manages resources in the cluster and manages the applications over Hadoop. It allows data stored in HDFS to be processed and run by various data processing engines such as batch processing, stream processing, interactive processing, graph processing, and many more. This increases efficiency with the use of YARN.

The Resource Manager is the core component of YARN – Yet Another Resource Negotiator. In analogy, it occupies the place of JobTracker of MRV1. Hadoop YARN is designed to provide a generic and flexible framework to administer the computing resources in the Hadoop cluster.

In this direction, the YARN Resource Manager Service (RM) is the central controlling authority

**Department of Computer Science and Engineering (Data Science)**

for resource management and makes allocation decisions ResourceManager has two main components: Scheduler and ApplicationsManager.



**Lab Assignments:**

1. Implement the following for the dataset given below:

   https://www.kaggle.com/datasets/narayan63/netflix-popular-movies-dataset

   a. Apply various HDFS commands on command line to get the appropriate file structure to run the MapReduce code.

   b. Write a Java MapReduce program to find out the list of movies who have been given ratings >= 8.

2. Implement a Java MapReduce program to display the occurrences of each character from a file.

# BIG DATA ENGINEERING

# EXPERIMENT 1.a

Name : Sarvagya Singh
SAPID : 60009200030
BATCH : K1

**Aim**: Hands on experience on Hadoop Ecosystem and MapReduce.

**Implementation**:

**Conclusion**: We have performed the various commands to achieve the given task.