

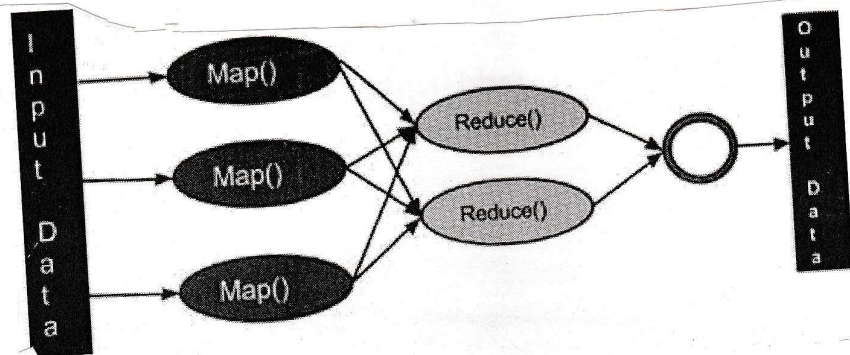
4.7 MAP REDUCE

"MapReduce is a framework using which we can write applications to process huge amount of data in parallel on large clusters of commodity hardware in a reliable manner."

4.7.1 Map Reduce :-

- MapReduce is a processing technique and a program model for distributed computing based on java.

- Two important task
- ① Map
- ② Reduce



① Map :-

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

② Reduce :-

- It takes the output from a map as an input and combines those data tuples into a smaller set of tuples.

- The reduce task is always performed after the map job.

- After processing, it produces a new set of output, which will be stored in the HDFS.

* During a MapReduce job, Hadoop sends the Map and Reduce tasks to the appropriate servers in the cluster.

* Most of the computing tasks place on nodes which data on local disks that reduces the network traffic.

* After completion of the given tasks, the cluster collects and reduces the data to form an appropriate result, and sends it back to the Hadoop server.

4.1.2 Map Reduce Input & Output:

- 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 as set of <key, value> pairs as the output of the job.

	<u>Input</u>	<u>Output</u>
map	$\langle k_1, v_1 \rangle$	list ($\langle k_2, v_2 \rangle$)
Reduce	$\langle k_2, \text{list}(v_2) \rangle$	list ($\langle k_3, v_3 \rangle$)

4.7.3 Input Splitting:-

- Hadoop divides the input into fixed-size pieces called input splits, or just splits.
- Each split is processed by a single map. Input Split represents the data to be processed by an individual Mapper.
- Each split is divided into records, and the map processes each record, which is a key value pair.
- Split is basically a number of rows and a record is that number.
- The input split does not contain input data but a reference to the data.

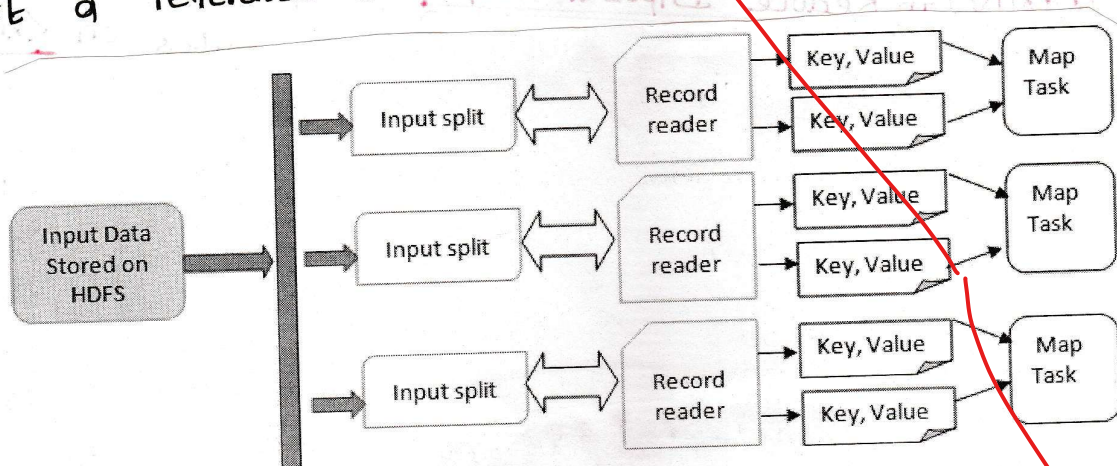


Fig: Input Splitting.

4.7.4 Other Terminology:-

* NameNode :- Node that manages the Hadoop Distributed File System (HDFS).

* DataNode :- Node where data is presented in advance before any processing takes place.

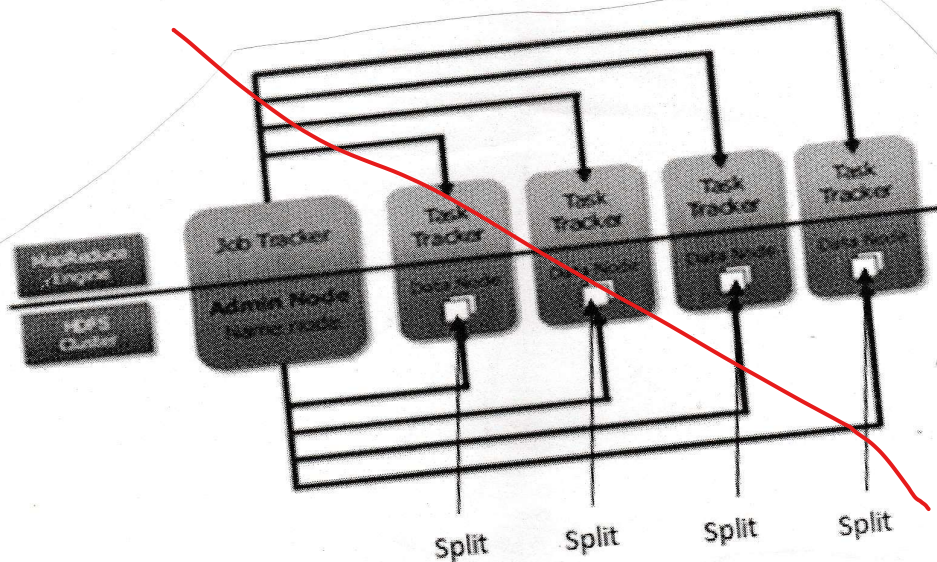


Fig: Job tracker Working Diagram.

* Master Node :- Node where JobTracker runs and which accepts job requests from clients.

* Job Tracker :- Schedules jobs and tracks the assign job to Task Tracker.

* Task Tracker :- Tracks the task and reports status to JobTracker.

Example Program :-

Mapper program:

```
package com.tom_e_white.drdobbs.mapreduce; // Package name

import org.apache.hadoop.io.LongWritable; // Import library packages
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException; // be
public class ProjectionMapper extends Mapper<LongWritable, Text, Text, LongWritable> { // Every Mapper class must extend from MapReduceBase class and it must implement Mapper Interface
    private Text word = new Text();
    private LongWritable count = new LongWritable(); // key, value pair

    @Override
    protected void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {
        // value is tab separated values: word, year, occurrences, #books, #pages
        // we project out (word, occurrences) so we can sum over all years
        String[] split = value.toString().split("\t+");
        word.set(split[0]);
        if (split.length > 2) {
            try {
                count.set(Long.parseLong(split[2]));
                context.write(word, count);
            } catch (NumberFormatException e) {
                // cannot parse - ignore
            }
        }
    }
}
```

Running through our tiny dataset, the map output looks like this:

```
("dobbs", 20)
("dobbs", 22)
("doctor", 545525)
("doctor", 668666)
```

In our abstract representation, the input to the reduce step looks like this:

```
("dobbs", [20, 22])
("doctor", [545525, 668666])
```

Reducer:

package org.apache.hadoop.mapreduce.lib.reduce;

Package name

import java.io.IOException;

import org.apache.hadoop.io.LongWritable;

import org.apache.hadoop.mapreduce.Reducer;

*Import library
Packages*

public class LongSumReducer<KEY> extends Reducer<KEY, LongWritable,
KEY, LongWritable> {

*Every reducer class must extend
mapreduce Base class and implement
Interface*

private LongWritable result = new LongWritable();

public void reduce(KEY key, Iterable<LongWritable> values,
Context context) throws IOException, InterruptedException {

*Every reducer
class must
provide begin
of reduce for*

long sum = 0;

for (LongWritable val : values) {

sum += val.get();

}

result.set(sum);

context.write(key, result);

}

The output of the reducer will be:

("dobbs", 42)

("doctor", 1214191)