

Map - Reduce Flow Chart

Word-Count Job

\$ hadoop jar test.jar DriverCode <outputdir>

file.txt (roomB)

input file (file.txt (roomB))

hi how are you
how is your job

how is your family
how is your winter

what is the time now
what is the time now

how is your brother
how is your brother

how is the hadoop
how is the hadoop

input split

input split

input split

input split

Record reader

Record reader

R.R

Record reader

Mapper

Mapper

Mapped

Mapper

Intermediate data
(data generated
by mapper & reducer)

(hi, 1) (how, 1)
(how, 1) (is, 1)
(is, 1) (your, 1)
(your, 1) (job, 1)

(how, 1) (now, 1)
(now, 1) (what, 1)
(what, 1) (is, 1)
(is, 1) (in, 1)
(in, 1) (is, 1)
(is, 1) (what, 1)
(what, 1) (the, 1)
(the, 1) (time, 1)
(time, 1) (hadop, 1)

shuffling & sorting

Reducer

Record-writer

File Input Formats:

1. Text Input Format
2. Key Value Text Input Format
3. Sequence File Input Format
4. Sequence File As Text Input Format.

Data Types:

<u>Wrapper class</u>	<u>Primitive Types</u>	<u>Box Classes</u>
Integer	int	IntWritable
Long	long	LongWritable
Float	float	FloatWritable
Double	double	DoubleWritable
String	String	Text
Character	char	CharWritable

Process Flow in Word-Count Job

1. Input file is divided into blocks of 64MB
2. "Record Reader", takes each input split and produces, `{Key, Value}` pairs based on "input file" file format.
For example: TextInputFormat file generates, `< LongWritable, Text >` as output, where, "ByteOffset of each line" is `LongWritable` key & "content of each line" as `text`.
3. Record reader does it in a loop till all the lines in input split are completed.
4. Next, Mapper takes the output of Record-reader and generates the

key value pairs of below format

(byteoffset, text)
(0, hi how are you)
↓
Mapper
↓
(hi -1) (how -1)
(how -1) (is -1)
(are -1) (your -1)
(you -1) (jobs -1)

5. Next in the shuffling phase, it combines all the keys to form unique keys. For example
(hi, 1) (now, 1) (are, 1) (is, 1) (your, 1) (jobs, 1)
6. Next, in sorting phase, all the keys will be arranged in the sorting order
7. Next, in Reducer, phase all be reduced to unique keys & their respective counts.
(hi -1) (now, 2) (are, 1) (is, 1) (your, 1) (jobs, 1)
8. Then using the "RecordWriter", output will be written to file in the output directory

file formats:

1. In "Text InputFormat", files will be made into (key, value) pairs of <Byte-offset, Text>.

* In Hadoop, everything is handled using the <KEY, VALUE> pairs format.

2. In "Key Value TextInputFormat", file will be split based on tab. For example,

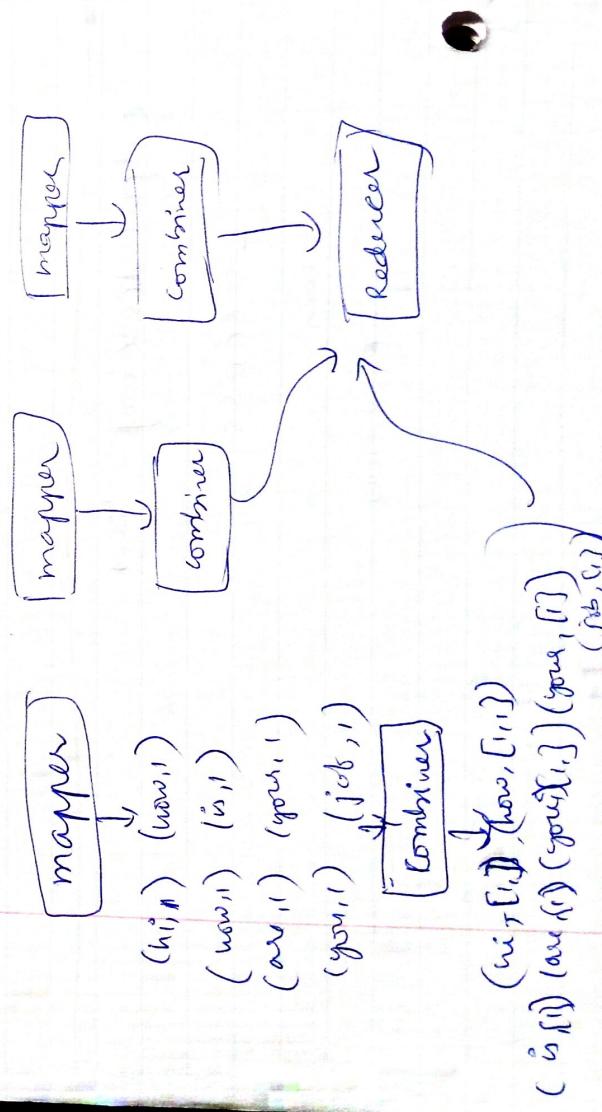
101 Ram 20 <101, (Ram, 20)>
102 Krish 30 => <102, (Krish 30)>
103 Rahim 40 <103, (Rahim 40)>

output directory files

- success
- logs
- part-00000

Combiner:

1. # of input splits = # of mappers.
 2. # of reducer = # of output files
 2. By default only "one reducer"
- * To reduce the traffic on the network
- due to congestion (because of all mappers sending data (key value) pairs to reducers.)
- We use "combiners".
- * Combiner is a mini-reducer on each mapper node.



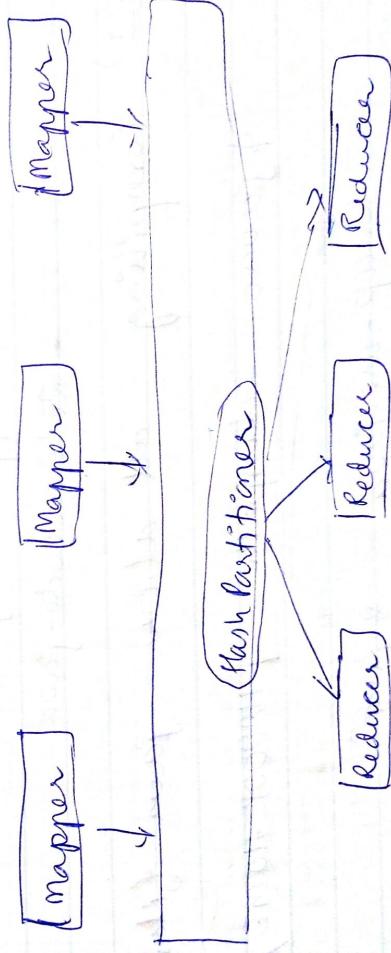
Partitioner

- If is used to send specific, $\langle \text{key}, \text{value} \rangle$ pairs to a specific Reducer.

2) It improves the performance & readability
(for example, all keys of length 3 to one partition,
keys of length 31 to the other etc...)

③ By default, 'HashPartitioner' is used, where
it randomly allocates the $\langle \text{key}, \text{value} \rangle$ pairs

to random reducers.
(based on Hash code of object)



HelloWorld Job (Sample Program)

- Step 1) Create a file ,file.txt
vi , (or) cat or touch
- Step 2) reading file.txt from local file system
to HDFS

\$ hadoop fs -put file.txt file
- Step 3) writing programs
 - DriverCode.java
 - MapperCode.java
 - ReducerCode.java
- Step 4) compiling all above java files

\$ javac -classpath \$ HADOOP_HOME /

↳ hadoop - code.java * .java.
- Step 5) Creating jar files
(Hadoop can only run the jar files)

```
$ jar cvf test.jar *.class
```

Step 6: Running above test.jar on file (which
is there in hdfs).

~~testjar~~
\$ hadoop jar test.jar DriverCode ~~file~~
file <test-output>

Hello Word Job

Word Count Job

Driver Code : WordCount.java

public class WordCount extends Configured implements
Tool {

public int run (String args[]) throws Exception

{
if (args.length < 2)

System.out.println ("Plz give ip & op
directories
properly")
return -1;

}

JobConf conf = new JobConf (WordCount.class);

conf.setJobName ("Word Count");

FileInputFormat.setInputPaths (conf, new
Path(args[0]));

FileOutputFormat.setOutputPath (conf, new
Path(args[1]));

```
conf.setMapperClass(WordMapper.class);  
conf.setReducerClass(WordReducer.class);  
conf.setMapOutputKeyClass(Text.class);  
conf.setMapOutputValueClass(IntWritable.class);  
conf.setOutputKeyClass(Text.class);  
conf.setOutputValueClass(IntWritable.class);  
JobClient.runJob(conf);  
return 0;  
}  
public static void main(String args) throws Exception {  
    int exitCode = ToolRunner.run(new  
        WordCount(), args);  
    System.exit(exitCode);  
}
```

Hadoop for wcfar



for more information on Main class, work
below "ToolRunner" class implementation

Class ToolRunner

of public static int run(Tool tool, String args)

```
{   tool.run(args); }
```

Mapper Code :- WordMapper.java

```
public class WordMapper extends Mapper<
```

```
MapReduceBase implements Mapper<  
LongWritable, Text, Text, IntWritable>  
    mapper tip key type  
    mapper tip value  
    tip value  
    tip type  
    tip key
```

```
    public void map (LongWritable key,  
                    Text value, OutputCollector <Text, IntWritable>
```

```
                    output, Reporter r) throws IOException {
```

```
    String s = value.toString();
```

```
    for (String word : s.split (" "))
```

```
    {  
        if (word.length () > 0)  
            output.collect (new Text (word),  
                           new IntWritable (1));  
    }
```

```
} // map and  
{} // clear and
```

Reducer Code :- Word Reducer-Java

public class WordReducer extends MapReduceBase
implements Reducer<Text, IntWritable, Text,
IntWritable>

{

 public void reduce(Text key, Iterator<
 IntWritable> values, OutputCollector<Text, IntWritable>
 output, Reporter r) throws IOException {

 while (values.hasNext())
 {
 IntWritable i = values.next();
 count += i.get(); // conversion from
 // obj type to
 // primitive
 // type
 }
 output.collect(key, new IntWritable(count));
 }

}

Partitioner code:

```
public class MyPartitioner implements
```

```
Partitioner < Text, IntWritable >
```

```
public void configure(JobConf conf)
```

```
public int getPartition(Text key,  
IntWritable value, int numRacks)
```

```
String s = key.toString();  
if (s.length() == 1)  
    return 0;  
if (s.length() == 2)  
    return 1;  
if (s.length() == 3)  
    return 2;  
else return 3;
```

configure Partitioner class in Driver
code, to run it.

```
conf.setPartitionerClass(MyPartitioner.class);
conf.setNumReduceTasks(4);
```

* Maximum of reducers

= 99999
Was
Output file can be
part-00000
↓
part-99999.

Hadoop distributed file System (HDFS)

commands:-

- 1) hadoop fs
- 2) hadoop fs -ls
 - list the metadata (directory & files)
- 3) if \downarrow multi-node cluster
hadoop fs -ls \Rightarrow hdfs://<hostname>/user/
 \swarrow user name
- 4) Create directory \Rightarrow myphen - added , compared to
unix systems
\$ hadoop fs -mkdir <dir-name>
- 5) Directory we cannot create file in
hdfs, we have to create it in local
file system & copy to hdfs.
- 6) Copy file to hdfs
local fs
\$ hadoop fs -put <src> <dest>
 \uparrow
(or)
~~-copyFromLocal~~
① moveFromLoc
② copyFromLoc

- ④ Get the file from hdfs file system to local file system
- local file system
- hdfs \uparrow \downarrow local file system
- \$ hadoop fs -get <src> <dest>
- copy To Local
 - move To Local
- ⑤ Working data directly into hdfs is not possible *
- ⑥ Viewing file content in hdfs
- \$ hadoop fs -cat <file-name>
- ⑦ remove file from hdfs
- \$ hadoop fs -rm <file-name>
- ⑧ remove file forcefully
- \$ hadoop fs -rmr <file-name> \downarrow directory

12) Copy ^{file} from one location to another

\$ hadoop fs -cp <src> <dest>

13) moving files
\$ hadoop fs -mv <src> <dest>

Disk Usage:

\$ hadoop fs -du <file-name>

14) Set replication factor

\$ hadoop fs -setrep <int> <file-name>

15)