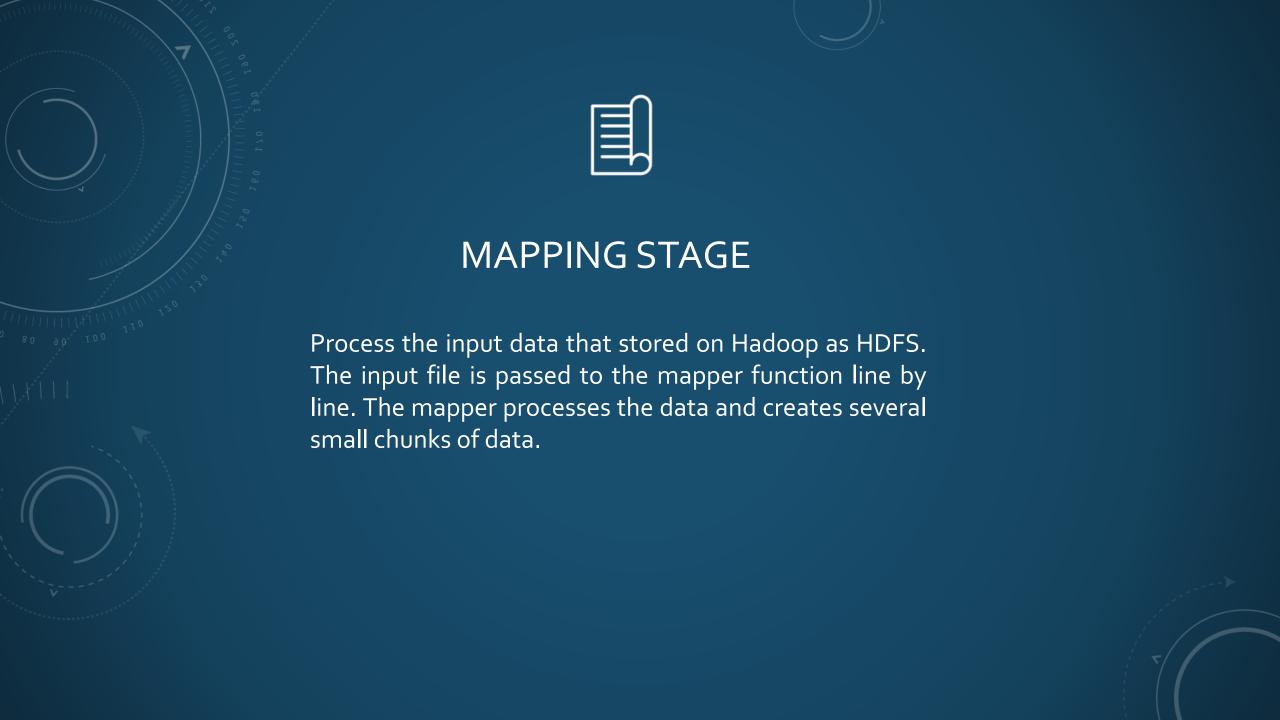
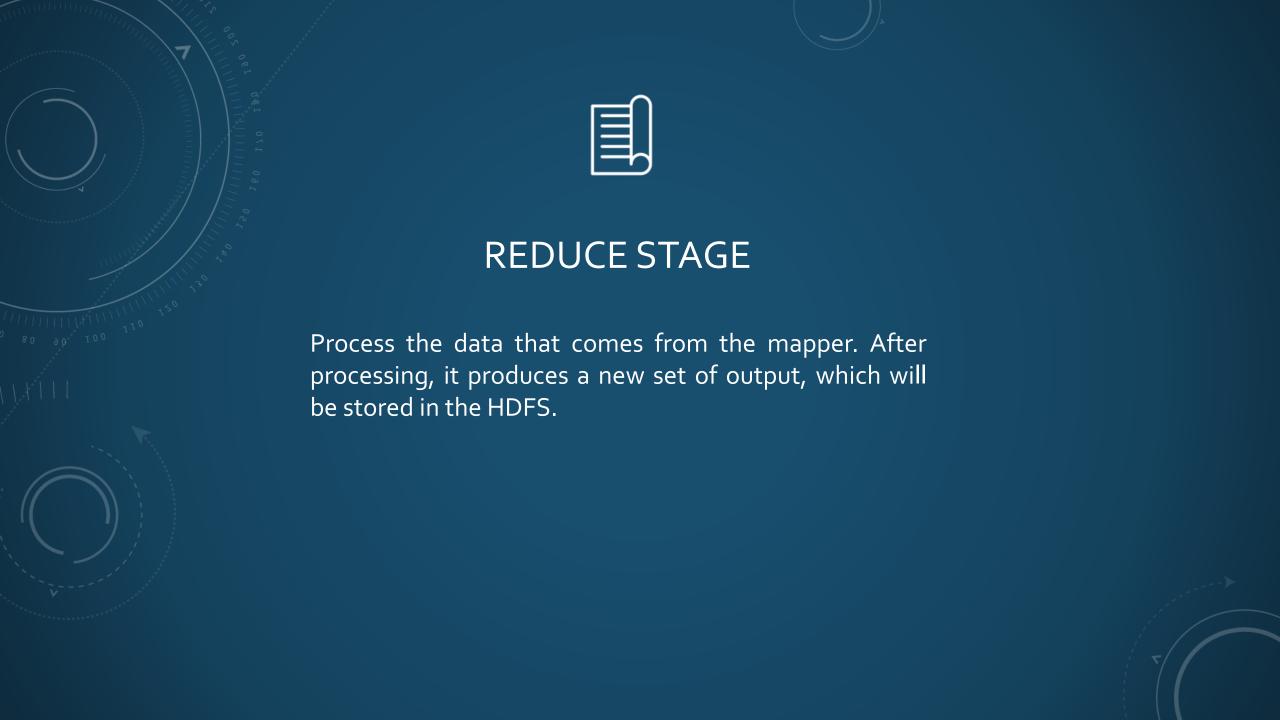




MapReduce is a processing technique and a program model for distributed computing based on java. It is consist of a framework using which we can write applications to process huge amounts of data, in parallel, on large clusters of commodity hardware in a reliable manner. Map reduce consist of mapping, shuffle and reduce stage.









```
ProcessUnits.java
data >
       package hadoop;
  1
       import java.util.*;
  3
  4
       import java.io.IOException;
       import java.io.IOException;
  6
       import org.apache.hadoop.fs.Path;
  8
       import org.apache.hadoop.conf.*;
       import org.apache.hadoop.io.*;
 10
       import org.apache.hadoop.mapred.*;
 11
 12
       import org.apache.hadoop.util.*;
```

data/ProcessUnit.java

- Q
- On datanode volumes, add ./data:/hadoop/data/
- This step is for mounting data folder on docker volumes so it can running through docker

```
datanode:
    image: bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8
    container_name: datanode
    restart: always
    volumes:
        - hadoop_datanode:/hadoop/dfs/data
        - ./data:/hadoop/data/
```



Run docker-compose up –d

```
D:\docker-hadoop>docker-compose up -d
Docker Compose is now in the Docker CLI, try `docker compose up`

Creating network "docker-hadoop_default" with the default driver
Creating historyserver ... done

Creating historyserver ... done

Creating nodemanager ... done

Creating nodemanager ... done

Creating namenode ... done
```

- For get into root enter docker exec —it datanode bash and cd to folder "data"
 D:\docker-hadoop>docker exec -it datanode bash root@782cdfa5dac5:/# cd hadoop/data
- Create folder units to compile java code

root@782cdfa5dac5:/hadoop/data# mkdir units

- Q
- javac -classpath hadoop-core-1.2.1.jar -d units ProcessUnits.java
- jar -cvf units.jar -C units/.

```
root@782cdfa5dac5:/hadoop/data# javac -classpath hadoop-core-1.2.1.jar -d units ProcessUnits.java
root@782cdfa5dac5:/hadoop/data# jar -cvf units.jar -C units/ .
added manifest
adding: hadoop/(in = 0) (out= 0)(stored 0%)
adding: hadoop/ProcessUnits$E_EMapper.class(in = 1980) (out= 814)(deflated 58%)
adding: hadoop/ProcessUnits$E_EReduce.class(in = 1661) (out= 678)(deflated 59%)
adding: hadoop/ProcessUnits.class(in = 1565) (out= 767)(deflated 50%)
```

Create directory for input and copy sample.txt to directory input

```
root@782cdfa5dac5:/hadoop/data# hadoop fs -mkdir /input_dir
root@782cdfa5dac5:/hadoop/data# hadoop fs -put sample.txt /input_dir/sample.txt
2021-06-13 14:15:46,838 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
```

- Q
- javac -classpath hadoop-core-1.2.1.jar -d units ProcessUnits.java
- jar -cvf units.jar -C units/.

```
root@782cdfa5dac5:/hadoop/data# javac -classpath hadoop-core-1.2.1.jar -d units ProcessUnits.java
root@782cdfa5dac5:/hadoop/data# jar -cvf units.jar -C units/ .
added manifest
adding: hadoop/(in = 0) (out= 0)(stored 0%)
adding: hadoop/ProcessUnits$E_EMapper.class(in = 1980) (out= 814)(deflated 58%)
adding: hadoop/ProcessUnits$E_EReduce.class(in = 1661) (out= 678)(deflated 59%)
adding: hadoop/ProcessUnits.class(in = 1565) (out= 767)(deflated 50%)
```

Create directory for input and copy sample.txt to directory input

```
root@782cdfa5dac5:/hadoop/data# hadoop fs -mkdir /input_dir
root@782cdfa5dac5:/hadoop/data# hadoop fs -put sample.txt /input_dir/sample.txt
2021-06-13 14:15:46,838 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
```



Sample.txt contain of data set

```
root@782cdfa5dac5:/hadoop/data# hadoop fs -cat /input dir/sample.txt
2021-06-13 15:03:52,432 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
                                                   26
1979
      23
                      2
                             43
                                            25
                                                           26
                                                                  26
                                                                          26
                                                                                 25
                                                                                        26
                                                                                                25
      26
                                    28
                                                                                                29
1980
              27
                                                                  31
                                    33
                                                                         34
                                                                                                34
1981
      31
                     32
                                                   35
                                                          36
                                                                  36
                                                                                 34
                                                                                        34
                                    39
                                           41
                                                   42
                                                          43
                                                                  40
1984
      39
                     39
                             39
                                                                         39
                                                                                 38
                                                                                        38
1985 38
                                                                                                45root@782cdfa5dac5:/hadoop/data#
```

Run hadoop jar units.jar hadoop.ProcessUnits input

root@782cdfa5dac5:/hadoop/data# hadoop jar units.jar hadoop.ProcessUnits /input_dir /output_dir 2021-06-13 14:16:33,427 INFO client.RMProxy: Connecting to ResourceManager at resourcemanager/172.25.0.3:8032

If job run successfully

```
2021-06-13 14:16:41,663 INFO mapreduce.Job: map 0% reduce 0%
2021-06-13 14:16:47,736 INFO mapreduce.Job: map 100% reduce 0%
2021-06-13 14:16:53,776 INFO mapreduce.Job: map 100% reduce 100%
2021-06-13 14:16:53,792 INFO mapreduce.Job: Job job_1623593386825_0002 completed successfully
```



• To view result enter hadoop fs –cat /output_dir/part-ooooo



FAILURE IN MAP REDUCE – TASK FAILURE

The child JVM reports the error back to its parent task tracker before it exits. The error ultimately makes it into the user logs. The task tracker marks the task attempt as failed, freeing up a slot to run another task.



FAILURE IN MAP REDUCE – TASK TRACKER FAILURE

Failure of a task tracker is another failure mode. If a task tracker fails by crashing or running very slowly, it will stop sending heartbeats to the job tracker or send them very infrequently



FAILURE IN MAP REDUCE – JOB TRACKER FAILURE

Failure of the job tracker is the most serious failure mode. Hadoop has no mechanism for dealing with job tracker failure it is a single point of failure so in this case all running jobs fail.