

Hadoop

What is Hadoop?

Hadoop is an open-source framework developed by the Apache Software Foundation that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from a single server to thousands of machines, each offering local computation and storage. Hadoop is highly fault-tolerant and is designed to detect and handle failures at the application layer, delivering a highly available service on top of a cluster of computers.

What Can We Do with Hadoop?

Store Large Data Sets:

HDFS (Hadoop Distributed File System) allows you to store massive amounts of data across multiple machines. It is designed to handle large files (gigabytes to petabytes) and supports high-throughput access to data.

Process Large Data Sets:

MapReduce is Hadoop's core processing model that allows you to process large data sets in parallel across a distributed cluster.

It breaks down tasks into smaller sub-tasks and processes them concurrently.

Data Analysis:

Hadoop supports various data analysis tools like Apache Hive for SQL-like querying and Apache Pig for scripting, making it easier to analyze large volumes of data.

Data Warehousing:

Hadoop can be used as a data warehouse with tools like Hive and HBase, which support structured and semi-structured data storage and querying.

Data Integration:

Tools like Apache Sqoop and Flume can be used to import data from relational databases or streaming sources into Hadoop for further processing and analysis.

Machine Learning and Data Mining:

Hadoop can integrate with tools like Apache Mahout or Spark MLlib to perform machine learning tasks on large data sets.

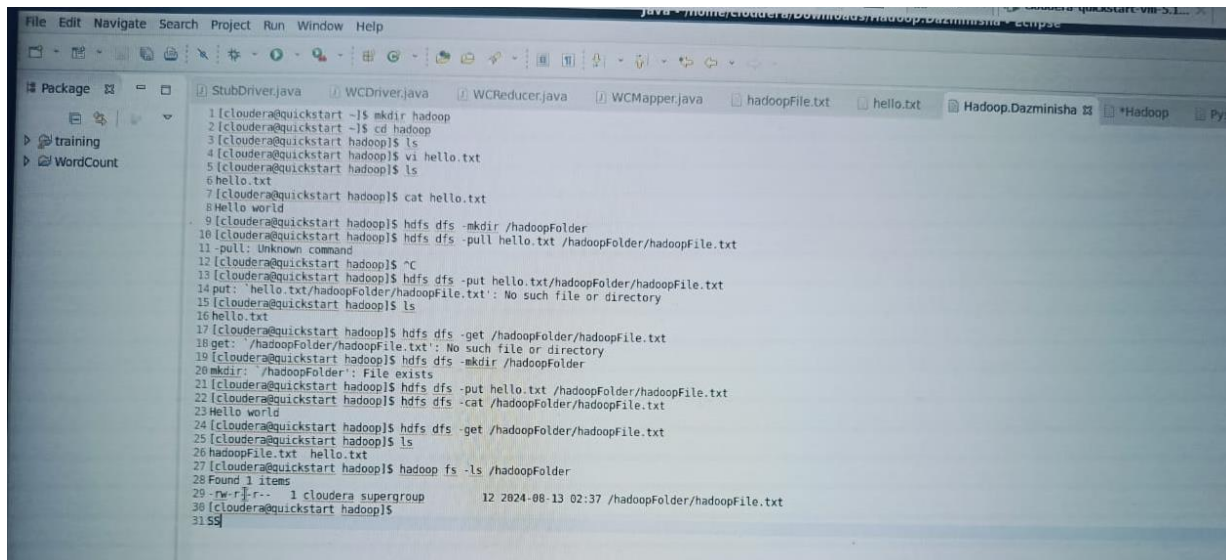
Summary

Hadoop is a powerful framework for distributed storage and processing of big data.

With Hadoop, you can store large volumes of data across multiple machines, process it using parallel computation, analyze it using tools like Hive and Pig,

and integrate data from various sources. Hadoop commands allow you to interact with HDFS,

submit and monitor jobs, query data, and manage resources within a Hadoop cluster.



```
File Edit Navigate Search Project Run Window Help
StubDriver.java WCDriver.java WCReducer.java WCMapper.java hadoopFile.txt hello.txt Hadoop.Dazminisha Hadoop Py
Package
training
WordCount
1 | cloudera@quickstart ~$ mkdir hadoop
2 | cloudera@quickstart ~$ cd hadoop
3 | cloudera@quickstart ~$ ls
4 | cloudera@quickstart ~$ vi hello.txt
5 | cloudera@quickstart ~$ ls
6 | hello.txt
7 | cloudera@quickstart ~$ cat hello.txt
8 | Hello world
9 | cloudera@quickstart ~$ hdfs dfs -mkdir /hadoopFolder
10 | cloudera@quickstart ~$ hdfs dfs -pull hello.txt /hadoopFolder/hadoopFile.txt
11 | -pull: Unknown command
12 | cloudera@quickstart ~$ ~C
13 | cloudera@quickstart ~$ hdfs dfs -put hello.txt/hadoopFolder/hadoopFile.txt
14 | put: 'hello.txt/hadoopFolder/hadoopFile.txt': No such file or directory
15 | cloudera@quickstart ~$ ls
16 | hello.txt
17 | cloudera@quickstart ~$ hdfs dfs -get /hadoopFolder/hadoopFile.txt
18 | get: /hadoopFolder/hadoopFile.txt: No such file or directory
19 | cloudera@quickstart ~$ hdfs dfs -mkdir /hadoopFolder
20 | mkdir: /hadoopFolder: File exists
21 | cloudera@quickstart ~$ hdfs dfs -put hello.txt /hadoopFolder/hadoopFile.txt
22 | cloudera@quickstart ~$ hdfs dfs -cat /hadoopFolder/hadoopFile.txt
23 | Hello world
24 | cloudera@quickstart ~$ hdfs dfs -get /hadoopFolder/hadoopFile.txt
25 | cloudera@quickstart ~$ ls
26 | hadoopFile.txt hello.txt
27 | cloudera@quickstart ~$ hadoop fs -ls /hadoopFolder
28 | Found 1 items
29 | -rw-r--r-- 1 cloudera supergroup 12 2024-08-13 02:37 /hadoopFolder/hadoopFile.txt
30 | cloudera@quickstart ~$
31 | $
```

```
StubDriver.java  WCDriver.java  WCReducer.java  WCMapper.java  hadoopFile.txt  hello.txt  Hadoop.Dazminisha  *Hadoop 2

Package  training  WordCount

1 [cloudera@quickstart ~]$ ls
2 avro enterprise-deployment.json Pictures
3 cloudera-manager express-deployment.json Public
4 cm_api.py external_jars rakeshdata
5 codegen_categories.java external-unified rakeshdata
6 datal hadoop sparkjars_exec
7 Dazminisha input.txt Templates
8 Desktop kerberos Videos
9 devices.json lib WordCount.jar
10 Documents Music workspace
11 Downloads parcels zeyo_tab.java
12 eclipse parquet.write
13 emp.java part.dir
14 [cloudera@quickstart ~]$ cat > /home/cloudera/WCFile.txt
15 TVS team
16 TVS team
17 Hadoop
18 Hadoop
19 Cloudera
20 Hadoop
21 Cloudera
22 Hadoop
23 Cloudera*Z
24 [1]+ Stopped cat > /home/cloudera/WCFile.txt
25 [cloudera@quickstart ~]$ hdfs dfs -mkdir /inputWC
26 mkdir: /: File exists
27 [cloudera@quickstart ~]$ hdfs dfs -ls /
28 Found 7 items
29 drwxrwxrwx - hdfs supergroup 0 2017-10-23 10:29 /benchmarks
30 drwxr-xr-x - cloudera supergroup 0 2024-08-13 02:37 /hadoopFolder
31 drwxr-xr-x - hbase supergroup 0 2024-08-13 01:11 /hbase
32 drwxr-xr-x - solr solr 0 2017-10-23 10:32 /solr
33 drwxrwxrwt - hdfs supergroup 0 2020-05-22 02:53 /tmp
34 drwxr-xr-x - hdfs supergroup 0 2020-05-22 02:53 /user
35 drwxr-xr-x - hdfs supergroup 0 2017-10-23 10:31 /var
36 [cloudera@quickstart ~]$ hdfs dfs -mkdir /inputWC
37 [cloudera@quickstart ~]$ hdfs dfs -ls /
38 Found 8 items
39 drwxrwxrwx - hdfs supergroup 0 2017-10-23 10:29 /benchmarks
40 drwxr-xr-x - cloudera supergroup 0 2024-08-13 02:37 /hadoopFolder
41 drwxr-xr-x - hbase supergroup 0 2024-08-13 01:11 /hbase
42 drwxr-xr-x - cloudera supergroup 0 2024-08-13 03:15 /inputWC
43 drwxr-xr-x - solr solr 0 2017-10-23 10:32 /solr
44 drwxrwxrwt - hdfs supergroup 0 2020-05-22 02:53 /tmp
45 drwxr-xr-x - hdfs supergroup 0 2020-05-22 02:53 /user
46 drwxr-xr-x - hdfs supergroup 0 2017-10-23 10:31 /var
47 [cloudera@quickstart ~]$ hdfs dfs -put /home/cloudera/WCFile.txt /inputWC/
48 [cloudera@quickstart ~]$ hdfs dfs -cat /inputWC/WCFile.txt
49 TVS team
50 TVS team
51 Hadoop
52 Hadoop
53 Cloudera
54 Hadoop
55 Cloudera
56 Hadoop
```

```
File Edit Navigate Search Project Run Window Help
7441 - /home/3cloudera/Downloads/Hadoop - 5.1.1...
Package StubDriver.java WCDriver.java WCRReducer.java WCMapper.java hadoopFile.txt hello.txt Hadoop.Dazminisha *Hadoop Pyspark
training
WordCount
55 Cloudera
56 Hadoop
57 |cloudera@quickstart ~| $ hadoop jar /home/cloudera/WordCount.jar WCDriver /inputWC/WCFile.txt /outWC
58 24/08/13 03:23:25 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
59 24/08/13 03:23:26 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
60 24/08/13 03:23:26 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this
61 24/08/13 03:23:26 INFO mapred.FileInputFormat: Total input paths to process : 1
62 24/08/13 03:23:26 WARN hdfs.DFSClient: Caught exception
63 java.lang.InterruptedException
64   at java.lang.Object.wait(Native Method)
65   at java.lang.Thread.join(Thread.java:1252)
66   at java.lang.Thread.join(Thread.java:1326)
67   at org.apache.hadoop.hdfs.DFSOutputStreamDataStreamer.closeResponder(DFSOutputStream.java:967)
68   at org.apache.hadoop.hdfs.DFSOutputStreamDataStreamer.endBlock(DFSOutputStream.java:795)
69   at org.apache.hadoop.hdfs.DFSOutputStreamDataStreamer.run(DFSOutputStream.java:894)
70 24/08/13 03:23:26 INFO mapreduce.JobSubmitter: number of splits:2
71 24/08/13 03:23:26 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1723536650149_0001
72 24/08/13 03:23:27 INFO impl.YarnClientImpl: Submitted application application_1723536650149_0001
73 24/08/13 03:23:27 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1723536650149_0001/
74 24/08/13 03:23:27 INFO mapreduce.Job: Running job: job_1723536650149_0001
75 24/08/13 03:23:36 INFO mapreduce.Job: Job job_1723536650149_0001 running in uber mode : false
76 24/08/13 03:23:36 INFO mapreduce.Job:  map 0% reduce 0%
77 24/08/13 03:23:45 INFO mapreduce.Job:  map 50% reduce 0%
78 24/08/13 03:23:46 INFO mapreduce.Job:  map 100% reduce 0%
79 24/08/13 03:23:50 INFO mapreduce.Job:  map 100% reduce 100%
80 24/08/13 03:23:51 INFO mapreduce.Job: Job job_1723536650149_0001 completed successfully
81 24/08/13 03:23:51 INFO mapreduce.Job: Counters: 49
82   File System Counters
83     FILE: Number of bytes read=130
84     FILE: Number of bytes written=430807
85     FILE: Number of read operations=0
86     FILE: Number of large read operations=0
87     FILE: Number of write operations=0
88     HDFS: Number of bytes read=308
89     HDFS: Number of bytes written=33
90     HDFS: Number of read operations=0
91     HDFS: Number of large read operations=0
92     HDFS: Number of write operations=2
93   Job Counters
94     Launched map tasks=2
95     Launched reduce tasks=1
96     Data-local map tasks=2
97     Total time spent by all maps in occupied slots (ms)=15310
98     Total time spent by all reduces in occupied slots (ms)=3055
99     Total time spent by all map tasks (ms)=15310
100    Total time spent by all reduce tasks (ms)=3055
101    Total vcore-milliseconds taken by all map tasks=15310
102    Total vcore-milliseconds taken by all reduce tasks=3055
103    Total megabyte-milliseconds taken by all map tasks=15677440
104    Total megabyte-milliseconds taken by all reduce tasks=3128320
105   Map-Reduce Framework
106     Map input records=8
107     Map output records=10
108     Map output bytes=104
109     Map output materialized bytes=136
110     Input split bytes=204
```

```

Edit Navigate Search Project Run Window Help
- - - - -
package StubDriver.java WCDriver.java WCReducer.java WCMapper.java hadoopFile.txt hello.txt Hadoop.Dazminisha *Hadoop Pyspark
/ training
/ WordCount
188 HDFS: Number of bytes read=300
189 HDFS: Number of bytes written=33
190 HDFS: Number of read operations=9
191 HDFS: Number of large read operations=0
192 HDFS: Number of write operations=2
193 Job Counters
194   Launched map tasks=2
195   Launched reduce tasks=1
196   Data-local map tasks=2
197   Total time spent by all maps in occupied slots (ms)=15310
198   Total time spent by all reduces in occupied slots (ms)=3055
199   Total time spent by all map tasks (ms)=15310
200   Total time spent by all reduce tasks (ms)=3055
201   Total vcore-milliseconds taken by all map tasks=15310
202   Total vcore-milliseconds taken by all reduce tasks=3055
203   Total megabyte-milliseconds taken by all map tasks=15677440
204   Total megabyte-milliseconds taken by all reduce tasks=3128320
205 Map-Reduce Framework
206   Map input records=8
207   Map output records=10
208   Map output bytes=104
209   Map output materialized bytes=136
210   Input split bytes=204
211   Combine input records=0
212   Combine output records=0
213   Reduce input groups=4
214   Reduce shuffle bytes=136
215   Reduce input records=10
216   Reduce output records=4
217   Spilled Records=20
218   Shuffled Maps =2
219   Failed Shuffles=0
220   Merged Map outputs=2
221   GC time elapsed (ms)=312
222   CPU time spent (ms)=2010
223   Physical memory (bytes) snapshot=647938048
224   Virtual memory (bytes) snapshot=8249905152
225   Total committed heap usage (bytes)=587014144
226 Shuffle Errors
227   BAD_ID=0
228   CONNECTION=0
229   IO_ERROR=0
230   WRONG_LENGTH=0
231   WRONG_MAP=0
232   WRONG_REDUCE=0
233 File Input Format Counters
234   Bytes Read=96
235 File Output Format Counters
236   Bytes Written=33
237 0
238 [cloudera@quickstart ~]$ hdfs dfs -ls /outWC
239 Found 2 items
240 -rw-r--r-- 1 cloudera supergroup 0 2024-08-13 03:23 /outWC/_SUCCESS
241 -rw-r--r-- 1 cloudera supergroup 33 2024-08-13 03:23 /outWC/part-00000
242 [cloudera@quickstart ~]$
243
Cloudera Live: Welco... java - from cloudera/ Writable Inse
```

Hadoop commands

Ls command to list the files/folders

```
ls
```

mkdir Hadoop: to create a folder

```
mkdir
```

cd : to change the directory

```
cd Hadoop/
```

Now create a file inside this hadoop folder as below:

Vi command to create a new file.

Press `i` to go on edit mode in this file and type any content as per your wish.

Once text is typed, press `escape` and then :wq

esc and ":wq"

To see the content of the file do as below: use cat command

```
cat hi.txt
```

The above are all simple commands.

Now let`s create a folder on our HDFS system (Hadoop file system) and learn the new commands:

hdfs dfs is Hadoop specific commands. (Hadoop Distributed File system)

```
hdfs dfs -mkdir /hadoopFolder
```

```
hdfs dfs -mkdir /newFolder
```

Lets copy the local file hello.txt onto the HDFS system as below:

```
hdfs dfs -put hello.txt /hadoopFolder/hadoopFile.txt
```

here we are copying the data of hello.txt onto a new file on HDFS system called hadoopFile.txt

put: command to copy a file from local machine to HDFS system.

```
hdfs dfs -put hello.txt/hadoopFolder/hadoopFile.txt
```

This way we were able to move a local file of our computer onto the HDFS system.

open hdfs file

```
hdfs dfs -cat /hadoopFolder/hadoopFile.txt
```

get command: to move the file from HDFS to local system.

```
hdfs dfs -get /hadoopFolder/hadoopFile.txt
```

The below commands are similar to that of put command which copies the local file into hdfs system.

```
hadoop fs -copyFromLocal hello.txt /hadoopFolder
```

```
hadoop fs -ls /hadoopFolder
```

HADOOP MAPREDUCE example using a WordCount program:

We need to now create a simple text file for Word counting purpose:

Create a WCFile.txt using cat command as below.

Once done with entering the text , press ctrl-z to exit.

```
cat > /home/cloudera/WCFile.txt
```

```
hdfs dfs -mkdir /inputWC
```



```
hdfs dfs -ls /
```

Now lets create a folder in the hdfs system to store this WCFFile.txt

Named the folder as inputWC

```
hdfs dfs -put /home/cloudera/WCFFile.txt /inputWC/
```

Lets put this WCFFile.txt onto the HDFS system:

Use the put command as below and also cross check the data using cat command

```
hdfs dfs -cat /inputWC/WCFFile.txt
```

Now lets execute the jar file for wordCounting

```
hadoop jar /home/cloudera/WordCount.jar WCDriver /inputWC/WCFFile.txt /outWC
```

The command to execute the jar file

```
hadoop jar /home/cloudera/WordCount.jar WCDriver /inputWC/WCFile.txt /outWC
```

Explaining the above command:

Hadoop jar -> this is a command to execute the jar file

/home/cloudera/WordCount.jar -> this is the location where our jar is located. WordCount.jar is the one which had exported from eclipse.

WCDriver -> This is a main java class which will be executed. [Refer the below code in eclipse]

/inputWC/WCFile.txt -> this is the file on hdfs which has the data to be counted.

/outWC -> this folder gets created once the below command is processed and stores the output.

To check the output execute as below

```
hdfs dfs -ls /outWC
```