

## NFS setup

```
testfile2.txt testfile.txt
root@smallnode1:/mnt/project# sudo df -hT
Filesystem                Type      Size  Used Avail Use% Mounted on
tmpfs                     tmpfs     196M  1.6M  195M   1% /run
/dev/sda3                  ext4      20G   13G   6.6G  65% /
tmpfs                     tmpfs     980M    0   980M   0% /dev/shm
tmpfs                     tmpfs     5.0M   4.0K   5.0M   1% /run/lock
/dev/sda2                  vfat      512M   6.1M   506M   2% /boot/efi
tmpfs                     tmpfs     196M  148K   196M   1% /run/user/1000
192.168.4.166:/exports/project nfs4      24G   14G   9.3G  60% /mnt/project
```

Here at the bottom, we can see nfs successfully on the client and my server is running on 192.168.4.166.

```
sudo: nfs: command not found
root@smallnode2:/home/rakesh# sudo df -hT
Filesystem                Type      Size  Used Avail Use% Mounted on
tmpfs                     tmpfs     196M  1.5M  195M   1% /run
/dev/sda3                  ext4      22G   14G   7.3G  65% /
tmpfs                     tmpfs     980M    0   980M   0% /dev/shm
tmpfs                     tmpfs     5.0M   4.0K   5.0M   1% /run/lock
/dev/sda2                  vfat      512M   6.1M   506M   2% /boot/efi
tmpfs                     tmpfs     196M  112K   196M   1% /run/user/1000
192.168.4.166:/exports/project nfs4      24G   14G   9.1G  61% /mnt/project
root@smallnode2:/home/rakesh#
```

This is my second node with nfs client running successfully and my server is running on 192.168.4.166.

And done for 3,4,5,6 nodes in the same way

## Hadoop Datasets

```
2023-11-15 03:37:55,602 INFO mapreduce.Job: map 84% reduce 0%
2023-11-15 03:37:56,604 INFO mapreduce.Job: map 87% reduce 0%
2023-11-15 03:37:57,607 INFO mapreduce.Job: map 88% reduce 0%
2023-11-15 03:38:01,615 INFO mapreduce.Job: map 89% reduce 0%
2023-11-15 03:38:02,618 INFO mapreduce.Job: map 90% reduce 0%
2023-11-15 03:38:17,658 INFO mapreduce.Job: map 92% reduce 0%
2023-11-15 03:38:21,669 INFO mapreduce.Job: map 95% reduce 0%
2023-11-15 03:38:22,672 INFO mapreduce.Job: map 98% reduce 0%
2023-11-15 03:38:23,676 INFO mapreduce.Job: map 99% reduce 0%
2023-11-15 03:38:27,691 INFO mapreduce.Job: map 100% reduce 0%
2023-11-15 03:38:43,736 INFO mapreduce.Job: map 100% reduce 19%
2023-11-15 03:38:49,750 INFO mapreduce.Job: map 100% reduce 29%
2023-11-15 03:41:32,109 INFO mapreduce.Job: map 100% reduce 30%
2023-11-15 03:41:38,123 INFO mapreduce.Job: map 100% reduce 35%
2023-11-15 03:41:44,134 INFO mapreduce.Job: map 100% reduce 36%
2023-11-15 03:41:50,145 INFO mapreduce.Job: map 100% reduce 38%
2023-11-15 03:41:56,156 INFO mapreduce.Job: map 100% reduce 40%
2023-11-15 03:42:02,167 INFO mapreduce.Job: map 100% reduce 42%
2023-11-15 03:42:08,179 INFO mapreduce.Job: map 100% reduce 44%
2023-11-15 03:42:14,190 INFO mapreduce.Job: map 100% reduce 46%
2023-11-15 03:42:20,201 INFO mapreduce.Job: map 100% reduce 47%
2023-11-15 03:42:26,211 INFO mapreduce.Job: map 100% reduce 49%
2023-11-15 03:42:32,222 INFO mapreduce.Job: map 100% reduce 51%
2023-11-15 03:42:38,232 INFO mapreduce.Job: map 100% reduce 53%
2023-11-15 03:42:44,244 INFO mapreduce.Job: map 100% reduce 55%
2023-11-15 03:42:50,254 INFO mapreduce.Job: map 100% reduce 57%
2023-11-15 03:42:56,265 INFO mapreduce.Job: map 100% reduce 58%
2023-11-15 03:43:02,275 INFO mapreduce.Job: map 100% reduce 60%
2023-11-15 03:43:08,286 INFO mapreduce.Job: map 100% reduce 62%
2023-11-15 03:43:14,297 INFO mapreduce.Job: map 100% reduce 64%
2023-11-15 03:43:20,308 INFO mapreduce.Job: map 100% reduce 66%
2023-11-15 03:43:26,320 INFO mapreduce.Job: map 100% reduce 68%
2023-11-15 03:43:32,333 INFO mapreduce.Job: map 100% reduce 70%
2023-11-15 03:43:38,346 INFO mapreduce.Job: map 100% reduce 71%
2023-11-15 03:43:44,359 INFO mapreduce.Job: map 100% reduce 74%
2023-11-15 03:43:50,371 INFO mapreduce.Job: map 100% reduce 77%
2023-11-15 03:43:56,385 INFO mapreduce.Job: map 100% reduce 80%
2023-11-15 03:44:02,398 INFO mapreduce.Job: map 100% reduce 82%
2023-11-15 03:44:08,410 INFO mapreduce.Job: map 100% reduce 84%
2023-11-15 03:44:14,424 INFO mapreduce.Job: map 100% reduce 87%
```

The screenshot of mapreduce execution jobs

```

root@8a5bedbca8ad:/code/Hadoop_Code# hdfs dfs -cat /user/root/output/* | head
2023-11-14 19:46:09,821 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteHostTrusted = false
20 20 20 20 22 4f 21 75 76 65 "O!uve 00000000000000000000000000001228D4 77778888000022224444DDDDDDDEEEEE00000000CCCC7777DDDD
, K4a-:v 00000000000000000000000000001B8132 5555EEEE888899994444FFFF1111CCCCEEEE1111EEEE6666FFFF
.FuD\}u 0000000000000000000000000000797631 5555DDDD8888000077772222111122224444DDDDDDDD99996666
;5YThct 00000000000000000000000000007D3DF5 2222AAAACCCCFFFFAAAA44445555EEEE44442222DDDD99992222
=2G*9{- 0000000000000000000000000000809EE 5555DDDD1111CCCC9999888800008888CCCCFFFFCCCC44443333
N}M9?sP 000000000000000000000000000029597 5555FFFF00007777555599991111CCCC66669999AAAAEEEE8888
P0X?R& 000000000000000000000000000041162E 888833339999FFFF1111CCCC8888CCCC9999EEEEDDDD00003333
[Xq\%%$ 000000000000000000000000000097A5F0 66666666EEEEDDDD7777FFFF00005555FFFFFFFFFF888855551111
rAmQg4v 00000000000000000000000000008180D3 BBBB111111119999FFFFFFFFFFF4444888888884444CCCC
!&))Jf3; 00000000000000000000000000004602E1 4444FFFFCCCC88888888CCCCFFFFCCCCEEEE5555666666666666
cat: Unable to write to output stream.
root@8a5bedbca8ad:/code/Hadoop_Code# hdfs dfs -cat /user/root/output/* | tail
2023-11-14 19:46:20,103 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteHostTrusted = false
7e 7e 7e 25 41 20 4e 42 5f 74 ~~~~%A NB_t 00000000000000000000000000003DE5EC FFFF7777EEEEBBBB4444EEEEEEEE33339999DDDD999900005555
7e 7e 7e 2d 63 2d 43 51 28 3e ~~~~c-CQ(> 0000000000000000000000000000832611 9999FFFF111177773333777700001111444444440000BBBB6666
7e 7e 7e 3d 59 7d 46 71 6c 2a ~~~~8Y}FqL* 0000000000000000000000000000742A4C BBBB1111CCCCEEEE888800000000777733333333DDDD22225555
7e 7e 7e 3e 44 64 3d 51 54 5d ~~~~>Dd=QT] 0000000000000000000000000000674C0F 9999DDDD000055556666CCCC22220000FFFFEEEEDDDDFFFF0000
7e 7e 7e 5d 4a 41 28 7d 6a 24 ~~~~]JA{}}$ 000000000000000000000000000060BCBE 111122223333444455559999AAAA88889999FFFFDDDDBBBB3333
7e 7e 7e 5d 5a 70 2e 23 2f 2b ~~~~]Zp.#/+ 00000000000000000000000000003B9A5A CCCC8888EEEEAAAAEEEE3333333377770000FFFFCCCC66667777
7e 7e 7e 5f 6a 51 65 70 69 78 ~~~~_jQepix 000000000000000000000000000011E5D4 11119999111155558888111100002222EEEE666688887777DDDD
7e 7e 7e 6e 74 3d 5a 48 5b 4e ~~~~nt=ZH[N 0000000000000000000000000000332A13 444411118888888833337777FFFF44445555555533330000CCCC
7e 7e 7e 73 2f 50 71 2c 2d 45 ~~~~s/Pq,-E 00000000000000000000000000006BE930 2222DDDDDDDD77771111EEEECCCC7777BBBB4444888811111111
7e 7e 7e 7a 62 41 5f 20 54 74 ~~~~zbA_ Tt 00000000000000000000000000007F9F4F BBBBCCCC666655559999FFFF8888AAAA11116666AAAA88880000
root@8a5bedbca8ad:/code/Hadoop_Code# hdfs dfs -rm -r /user/root/output

```

The above screenshot shows the head and tail of the Hadoop Sort program for 1GB and similarly for 4GB,16GB,32GB

### Spark 1GB small instance

Command used :

```
time spark-submit --class SortGensortData --master local[4] --executor-memory 24G /home/cc/my-spark-project/target/spark-example-1.0-SNAPSHOT.jar 1GB 1GBoutput
```

And i got output and here is the sorted head and tail

```

cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/1GBoutput$ cat part-* > outputfile.txt
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/1GBoutput$ cd ..
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input$ cd 1GBoutput/
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/1GBoutput$ ls
_SUCCESS      part-00001  part-00004  part-00007  part-00010  part-00013  part-00016  part-00019  part-00022
outputfile.txt part-00002  part-00005  part-00008  part-00011  part-00014  part-00017  part-00020  part-00023
part-00000    part-00003  part-00006  part-00009  part-00012  part-00015  part-00018  part-00021
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/1GBoutput$ cat outputfile.txt | head
"O!uve 00000000000000000000000000001228D4 77778888000022224444DDDDDDDEEEEE00000000CCCC7777DDDD
,K4a-:v 00000000000000000000000000001B8132 5555EEEE888899994444FFFF1111CCCCEEEE1111EEEE6666FFFF
.FuD\}u 0000000000000000000000000000797631 5555DDDD8888000077772222111122224444DDDDDDDD99996666
;5YThct 00000000000000000000000000007D3DF5 2222AAAACCCCFFFFAAAA44445555EEEE44442222DDDD99992222
=2G*9{- 0000000000000000000000000000809EE 5555DDDD1111CCCC9999888800008888CCCCFFFFCCCC44443333
N}M9?sP 000000000000000000000000000029597 5555FFFF00007777555599991111CCCC66669999AAAAEEEE8888
P0X?R& 000000000000000000000000000041162E 888833339999FFFF1111CCCC8888CCCC9999EEEEDDDD00003333
[Xq\%%$ 000000000000000000000000000097A5F0 66666666EEEEDDDD7777FFFF00005555FFFFFFFFFF888855551111
rAmQg4v 00000000000000000000000000008180D3 BBBB111111119999FFFFFFFFFFF4444888888884444CCCC
!&))Jf3; 00000000000000000000000000004602E1 4444FFFFCCCC88888888CCCCFFFFCCCCEEEE5555666666666666
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/1GBoutput$ cat outputfile.txt | tail
~~~%A NB_t 00000000000000000000000000003DE5EC FFFF7777EEEEBBBB4444EEEEEEEE33339999DDDD999900005555
~~~c-CQ(> 0000000000000000000000000000832611 9999FFFF111177773333777700001111444444440000BBBB6666
~~~8Y}FqL* 0000000000000000000000000000742A4C BBBB1111CCCCEEEE888800000000777733333333DDDD22225555
~~~>Dd=QT] 0000000000000000000000000000674C0F 9999DDDD000055556666CCCC22220000FFFFEEEEDDDDFFFF0000
~~~]JA{}}$ 000000000000000000000000000060BCBE 111122223333444455559999AAAA88889999FFFFDDDDBBBB3333
~~~]Zp.#/+ 00000000000000000000000000003B9A5A CCCC8888EEEEAAAAEEEE3333333377770000FFFFCCCC66667777
~~~_jQepix 000000000000000000000000000011E5D4 11119999111155558888111100002222EEEE666688887777DDDD
~~~nt=ZH[N 0000000000000000000000000000332A13 444411118888888833337777FFFF44445555555533330000CCCC
~~~s/Pq,-E 00000000000000000000000000006BE930 2222DDDDDDDD77771111EEEECCCC7777BBBB4444888811111111
~~~zbA_ Tt 00000000000000000000000000007F9F4F BBBBCCCC666655559999FFFF8888AAAA11116666AAAA88880000
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/1GBoutput$

```



Command used :

```

cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input$ cd /GBoutput/
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/4GBoutput$ ls
_SUCCESS      part-00001    part-00003    part-00005    part-00007    part-00009    part-00011    part-00013    part-00015    part-00017    part-00019    part-00021    part-00023
part-00000    part-00002    part-00004    part-00006    part-00008    part-00010    part-00012    part-00014    part-00016    part-00018    part-00020    part-00022
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/4GBoutput$ cat part_* > output
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/4GBoutput$ cat output | head
"0!uive 00000000000000000000000001228D4 777788880002224444DDDDDDDEEE0000000CCCC7777DDDD
'3C0J, 5555AAAA9999EEEE888822229999CCCCDDDD6666555544442222
!&$3/J] 000000000000000000000000002145D78 8888BBBDDDD1111CCCC55556666BBB1111EEEEDDDD22229999
!_#E_9 33332222FFFFBBB0000FFFFAAA666655553333DDDD3333CCCC
&&!ep0J 000000000000000000000000001CB42E 2222BBB2222222FFFFFFFFFFCCCC5555666666777700003333
'AJ;h,; 000000000000000000000000001FD93E FFFFFF66669999FFFF44446666222233330000BBB33333333
(Vc$5P2 000000000000000000000000000A9F516 8888AAAA7777BBB7777AAAA6666DDDEEE9999BBB8BBB8BBB
'K4a:-v 5555EEEE888899994444FFFF1111CCCCEEEE1111EEEE6666FFFF
'FuD;ju 00000000000000000000000000797631 5555DDDDBBB000077772222111122224444DDDDDD99996666
2;0EWt$ 0000000000000000000000000026D0B5F 00001111000000007777000022223333AAAA4444CCCCAAAA0000
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/4GBoutput$ cat output | tail
nmmnt=ZH[N 44441111BBB8BBB33337777FFFF44445555555533330000CCCC
nmmn0=:Exo 00000000000000000000000000198470A BBBBDDDD55554444FFFF999955553333AAAA8BBBEEEE8888DDDD
nmmn2[IR' | CCCCDDDD1111AAAAEEEE888855559999DDDD7777DDDD88889999
nmmnuR0guw 000000000000000000000000001E87D8A BBBBEEEE22222222BBBAAA55551111FFFF99995555FFFF7777
nmmnuPq;-E 2222DDDDDDDD77771111EEEECCCC7777BBB4444888811111111
nmmnuF&v'wv 6666BBB1111EEEEEEEE5555CCCCDDDD5555555555BBB8BBBDDDD
nmmnuKOL-gv 000000000000000000000000002048Z4 CCCC11114444888822226666BBB888855557777EEEEBBB80000
nmmnuzbA_Tt 000000000000000000000000007F9F40 BBBBCCCC666655559999FFFF77778888AAAA11116666AAAA8BBB8000
nmmnuze0*Feg 4444CCCCBBB899992222888855558888CCCCFF000011111111
nmmnuGxjWHI 777711118888AAAAAAA22221111BBB800002222BBB8CCCC2222
cc@ubuntu:~/Hadoop-on-Docker/code/Hadoop_Code/input/4GBoutput$

```

```
time spark-submit --class SortGensortData --master local[24] --executor-memory 96G
/home/cc/my-spark-project/target/spark-example-1.0-SNAPSHOT.jar 16GB 16GBoutput
```

[illegible]

```
time spark-submit --class SortGensortData --master local[24] --executor-memory 96G
/home/cc/my-spark-project/target/spark-example-1.0-SNAPSHOT.jar 32GB
32GBpennyout
```



6 nodes 32 GB	1273792	258481
---------------	---------	--------

Performance evaluation :

- All configurations see an increase in Hadoop Sort and Spark processing times as data sizes rise. This makes sense because processing larger datasets takes longer. The pace of rise differs throughout the setups, though.
- In terms of processing speed, Spark routinely beats Hadoop Sort for all data volumes and setups. This is probably because Hadoop Sort's disk-based processing is less efficient than Spark's in-memory data processing capabilities.
- When compared to the individual big and small instances, the 6-node cluster considerably boosts Hadoop Sort and Spark performance. This demonstrates the advantages of distributed computing in parallel.
- For Hadoop Sort and Spark, the large instance—which has more cores, RAM, and disk space—performs better than the tiny instance. This demonstrates how crucial it is to have more computer power for jobs requiring a lot of data.
- The Hadoop Sort on a tiny instance for the 32 GB dataset is tagged as "Not feasible," meaning that the small instance does not have the memory or disk capacity to handle such huge datasets efficiently.
- The 6-node cluster outperforms single instances in terms of performance, however as data size grows, the improvement ratio falls. For example, for smaller datasets, the time decrease from tiny to 6-node setups is more pronounced than for the 32 GB dataset. This points to a potential scaling limit brought on by factors including disk I/O constraints in bigger datasets, network overhead, and complicated data management.
- The absolute times demonstrate that processing times for huge datasets may be extremely high, even with a 6-node cluster (e.g., over a million milliseconds for 32 GB data in Hadoop Sort). Nevertheless, the cluster architecture offers a significant speed improvement over the single-node instances.

#### **For 32GB dataset:**

Large Instance, Hadoop Sort, 32 GB: Disk I/O Throughput = 2.37 MB/sec.

Large Instance, Spark, 32 GB: Disk I/O Throughput = 32.85 MB/sec.

6 Small Instances, Hadoop Sort, 32 GB: Disk I/O Throughput = 25.73 MB/sec.

6 Small Instances, Spark, 32 GB: Disk I/O Throughput = 126.82 MB/sec.

#### **Q and A evaluations**

- The absolute times demonstrate that, even with a 6-node cluster, processing times for big datasets may still be somewhat high (for example, more than a million milliseconds for 32 GB of data in Hadoop Sort). The cluster arrangement, however, offers a significant speed improvement above the single-node instances.

- Both Hadoop Sort and Spark require more processing time as data volume grows. But compared to Hadoop Sort, Spark's processing time increase is less noticeable, suggesting higher scalability.
- Although the huge instance performs well, bigger data quantities make its limits clear, especially for Hadoop Sort.
- Performance is greatly enhanced by the 6-node cluster's dispersed processing. This is demonstrated via Hadoop Sort and Spark, demonstrating the advantages of parallel processing, particularly for bigger data sets.
- The increase from one to six tiny instances shows robust and positive scalability, particularly for Spark. But the scaling is nonlinear, suggesting that distributed processing has overhead costs.
- The 6-node cluster performs better at weak scaling than a single big instance, even though both configurations scale with higher data volumes.
- Both Hadoop Sort and Spark run well on the huge instance, with Spark exhibiting notable temporal savings.
- The 6-node small instance cluster is more appropriate, especially for Spark, which manages big datasets much more effectively.
- In single and multi-node deployments, Spark processes huge datasets more efficiently than Hadoop Sort. In single and multi-node deployments, Spark processes huge datasets more efficiently than Hadoop Sort.