# CLONING THE CLUSTER:



# STARTED & WAITED:

```
Last login: Thu Mar 28 12:07:01 on ttys000
                    MacBook-Air ~ % ssh -i ~/Downloads/MyNewMacKey.pem hadoop@ec2-3-85-241-38.compute-1.amazonaws.com
The authenticity of host 'ec2-3-85-241-38.compute-1.amazonaws.com (3.85.241.38)' can't be established.
ED25519 key fingerprint is SHA256:zh7U2tcYBIlF1eyO3rlzk41J8oCSdSk5z+diISrhroo.
This key is not known by any other names
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-3-85-241-38.compute-1.amazonaws.com' (ED25519) to the list of known hosts.

A newer release of "Amazon Linux" is available.
  Version 2023.4.20240319:
Run "/usr/bin/dnf check-release-update" for full release and version update info
     ,     #_
   ~\_  ####_        Amazon Linux 2023
  ~~  \_#####\
  ~~     \###|
  ~~       \#/ ___   https://aws.amazon.com/linux/amazon-linux-2023
   ~~       V~' '->
    ~~~         /
      ~~._.   _/
         _/ _/
       _/m/'
Last login: Thu Mar 28 17:36:41 2024

EEEEEEEEEEEEEEEEEEEEE MMMMMMMM         MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::::M        M:::::::M R::::::::::::::R
EE:::::EEEEEEEEE:::E M::::::::M        M::::::::M R:::::RRRRRR:::::R
  E::::E       EEEEE M:::::::::M      M:::::::::M RR::::R      R::::R
  E::::E            M::::::M::M     M:::M::::::M   R:::R       R::::R
  E:::::EEEEEEEEEE   M:::::M M:::M M:::M M:::::M   R:::RRRRRR::::R
  E::::::::::::::E   M:::::M  M:::M:::M  M:::::M   R:::::::::::RR
  E:::::EEEEEEEEEE   M:::::M   M:::::M   M:::::M   R:::RRRRRR:::R
  E::::E             M:::::M    M:::M    M:::::M   R:::R       R::::R
  E::::E       EEEEE M:::::M     MMM     M:::::M   R:::R       R::::R
EE:::::EEEEEEEEE:::E M:::::M             M:::::M   R:::R       R::::R
E::::::::::::::::::E M:::::M             M:::::M RR::::R       R::::R
EEEEEEEEEEEEEEEEEEEEE MMMMMMM            MMMMMMM RRRRRRR        RRRRRR

[hadoop@ip-172-31-68-213 ~]$ ls
[hadoop@ip-172-31-68-213 ~]$ ls
NYSE_DATA.txt
[hadoop@ip-172-31-68-213 ~]$ nano NYSE_*
[hadoop@ip-172-31-68-213 ~]$ hadoop fs -mkdir MR
```
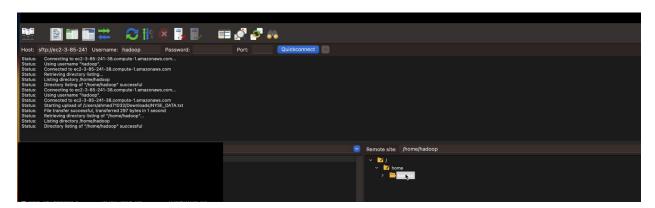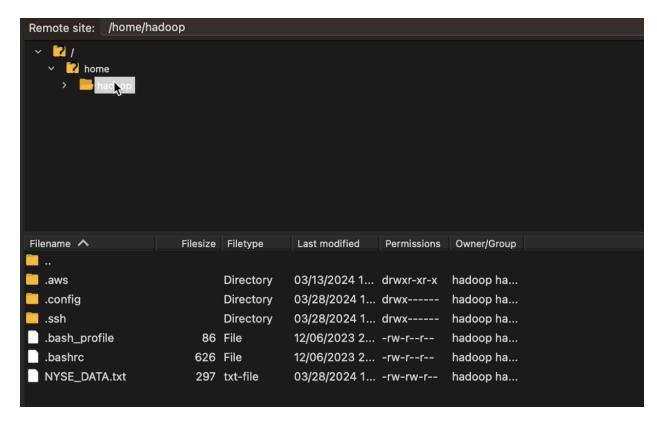
FILE ZILLA STEPS:

| Filename ∧ | Filesize | Filetype | Last modified | Permissions | Owner/Group |
|---|---|---|---|---|---|
| .. | | | | | |
| .aws | | Directory | 03/13/2024 1... | drwxr-xr-x | hadoop ha... |
| .config | | Directory | 03/28/2024 1... | drwx------ | hadoop ha... |
| .ssh | | Directory | 03/28/2024 1... | drwx------ | hadoop ha... |
| .bash_profile | 86 | File | 12/06/2023 2... | -rw-r--r-- | hadoop ha... |
| .bashrc | 626 | File | 12/06/2023 2... | -rw-r--r-- | hadoop ha... |
| NYSE_DATA.txt | 297 | txt-file | 03/28/2024 1... | -rw-rw-r-- | hadoop ha... |

CREATING DIRECTORY AND MAKING EDITS TO THE MAPPER & REDUCER PYTHON FILES AND EXECUTING IT:

```
2024-03-28 18:18:30,425 INFO mapreduce.Job:  map 100% reduce 100%
2024-03-28 18:18:30,434 INFO mapreduce.Job: Job job_1711647601769_0001 completed successfully
2024-03-28 18:18:30,552 INFO mapreduce.Job: Counters: 55
        File System Counters
                FILE: Number of bytes read=204
                FILE: Number of bytes written=3245514
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=2332
                HDFS: Number of bytes written=35
                HDFS: Number of read operations=39
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=6
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Killed map tasks=1
                Launched map tasks=8
                Launched reduce tasks=3
                Rack-local map tasks=8
                Total time spent by all maps in occupied slots (ms)=4642416
                Total time spent by all reduces in occupied slots (ms)=1611264
                Total time spent by all map tasks (ms)=96717
                Total time spent by all reduce tasks (ms)=16784
                Total vcore-milliseconds taken by all map tasks=96717
                Total vcore-milliseconds taken by all reduce tasks=16784
                Total megabyte-milliseconds taken by all map tasks=148557312
                Total megabyte-milliseconds taken by all reduce tasks=51560448
        Map-Reduce Framework
                Map input records=25
                Map output records=25
                Map output bytes=174
                Map output materialized bytes=608
                Input split bytes=976
                Combine input records=0
                Combine output records=0
                Reduce input groups=5
                Reduce shuffle bytes=608
                Reduce input records=25
                Reduce output records=5
                Spilled R
                Shuffled
                Failed Shuffles=0
                Merged Map outputs=24
                GC time elapsed (ms)=886
                CPU time spent (ms)=14820
                Physical memory (bytes) snapshot=4742402048
                Virtual memory (bytes) snapshot=39247581184
                Total committed heap usage (bytes)=4715446272
                Peak Map Physical memory (bytes)=511586304
                Peak Map Virtual memory (bytes)=3199328256
                Peak Reduce Physical memory (bytes)=296660992
                Peak Reduce Virtual memory (bytes)=4553580544
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
```

FINAL RESULT:

```
                          Bytes Written=35
2024-03-28 18:18:30,553 INFO streaming.StreamJob: Output directory: MR/output
[hadoop@ip-172-31-68-213 ~]$ hadoop fs -ls MR
Found 2 items
-rw-r--r--   1 hadoop hdfsadmingroup        299 2024-03-28 18:10 MR/NYSE_DATA.txt
drwxr-xr-x   - hadoop hdfsadmingroup          0 2024-03-28 18:18 MR/output
[hadoop@ip-172-31-68-213 ~]$ hadoop fs -ls MR/output
Found 4 items
-rw-r--r--   1 hadoop hdfsadmingroup          0 2024-03-28 18:18 MR/output/_SUCCESS
-rw-r--r--   1 hadoop hdfsadmingroup         14 2024-03-28 18:18 MR/output/part-00000
-rw-r--r--   1 hadoop hdfsadmingroup         21 2024-03-28 18:18 MR/output/part-00001
-rw-r--r--   1 hadoop hdfsadmingroup          0 2024-03-28 18:18 MR/output/part-00002
[hadoop@ip-172-31-68-213 ~]$ hadoop fs -cat MR/output/p*
GE      15
MCD     175
BAC     26
CAH     56
PFE     45
[hadoop@ip-172-31-68-213 ~]$
Broadcast message from root@localhost (Thu 2024-03-28 18:24:40 UTC):

The system will power off now!

Connection to ec2-3-85-241-38.compute-1.amazonaws.com closed by remote host.
Connection to ec2-3-85-241-38.compute-1.amazonaws.com closed.
```

## TERMINATING THE CLUSTER & ENDING THE LAB: