**Lab 5 Report**
**Student:** Ubaidullauly Azamat
**Group:** IT-2305

I created a cluster and installed a dataset via this command:
# on master node
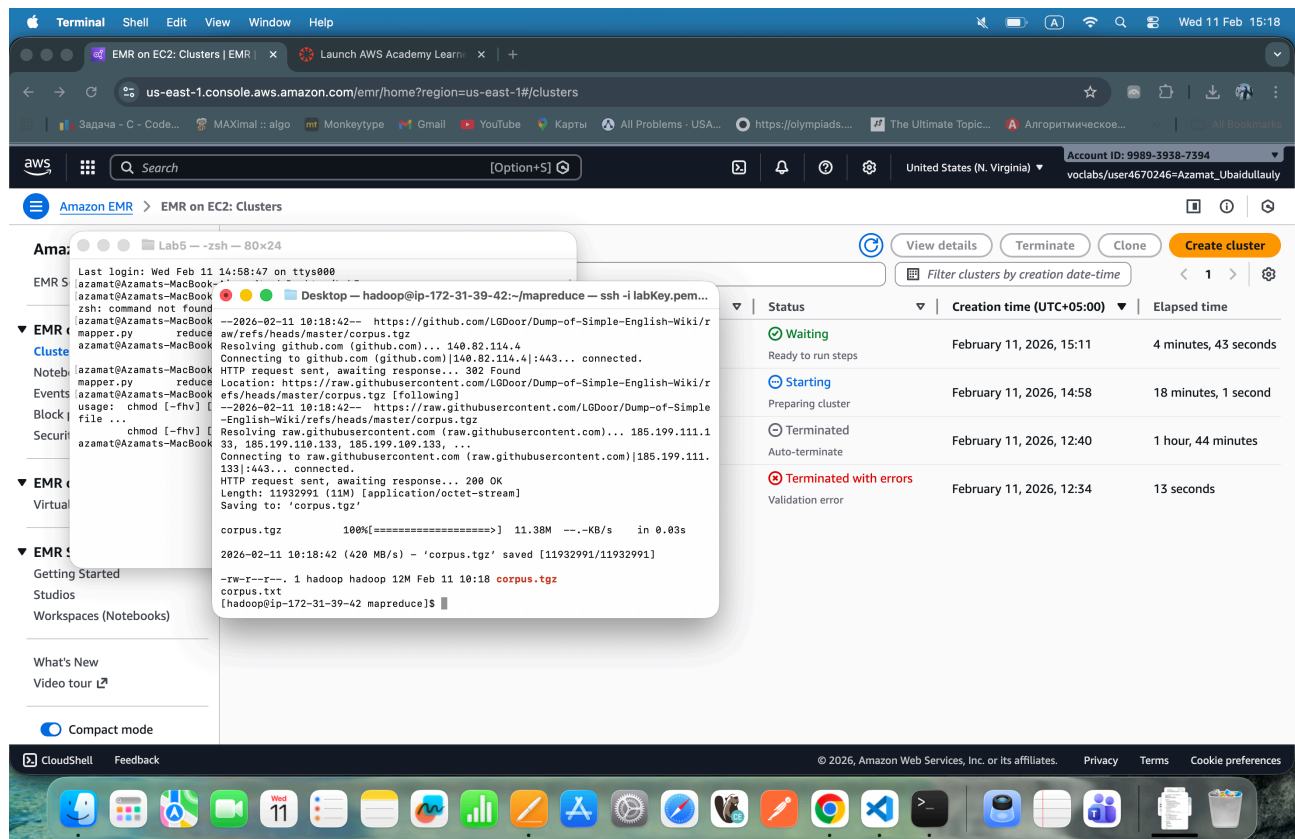cd ~
mkdir -p mapreduce && cd mapreduce

wget https://github.com/LGDoor/Dump-of-Simple-English-Wiki/raw/refs/heads/master/corpus.tgz
ls -lh corpus.tgz
tar -xvzf corpus.tgz



Then we put it into HDFS using the commands below:

```
hdfs dfs -mkdir -p /user/hadoop/input
hdfs dfs -put -f corpus.txt /user/hadoop/input/
hdfs dfs -ls /user/hadoop/input/
```

## Check cluster if it is running:

```
-files <file1,...>               specify a comma-separated list of files to be copied to the map reduce cluster
-libjars <jar1,...>              specify a comma-separated list of jar files to be included in the classpath
-archives <archive1,...>         specify a comma-separated list of archives to be unarchived on the compute machines

The general command line syntax is:
command [genericOptions] [commandOptions]


For more details about these options:
Use $HADOOP_HOME/bin/hadoop jar hadoop-streaming.jar -info

Try -help for more information
Streaming Command Failed!
[hadoop@ip-172-31-39-42 mapreduce]$ yarn node -list
2026-02-11 10:22:49,624 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at ip-172-31-39-42.ec2.internal/172.31.39.42:8032
2026-02-11 10:22:49,721 INFO client.AHSProxy: Connecting to Application History server at ip-172-31-39-42.ec2.internal/172.31.39.42:10200
Total Nodes:2
         Node-Id           Node-State Node-Http-Address       Number-of-Running-Containers
ip-172-31-47-213.ec2.internal:8041            RUNNING ip-172-31-47-213.ec2.internal:8042                     0
ip-172-31-37-21.ec2.internal:8041             RUNNING ip-172-31-37-21.ec2.internal:8042                     0
[hadoop@ip-172-31-39-42 mapreduce]$ hdfs dfsadmin -report
Configured Capacity: 62141726720 (57.87 GB)
Present Capacity: 62136997792 (57.87 GB)
DFS Remaining: 62080337115 (57.82 GB)
DFS Used: 56660677 (54.04 MB)
DFS Used%: 0.09%
Replicated Blocks:
        Under replicated blocks: 0
        Blocks with corrupt replicas: 0
        Missing blocks: 0
        Missing blocks (with replication factor 1): 0
        Low redundancy blocks with highest priority to recover: 0
        Pending deletion blocks: 0
Erasure Coded Block Groups:
        Low redundancy block groups: 0
        Block groups with corrupt internal blocks: 0
        Missing block groups: 0
        Low redundancy blocks with highest priority to recover: 0
        Pending deletion blocks: 0

-------------------------------------------------
Live datanodes (1):

Name: 172.31.47.213:9866 (ip-172-31-47-213.ec2.internal)
Hostname: ip-172-31-47-213.ec2.internal
Decommission Status : Normal
Configured Capacity: 62141726720 (57.87 GB)
DFS Used: 56660677 (54.04 MB)
Non DFS Used: 4728928 (4.51 MB)
DFS Remaining: 62080337115 (57.82 GB)
DFS Used%: 0.09%
DFS Remaining%: 99.90%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 0
Last contact: Wed Feb 11 10:22:57 UTC 2026
Last Block Report: Wed Feb 11 10:14:30 UTC 2026
Num of Blocks: 2


[hadoop@ip-172-31-39-42 mapreduce]$
```