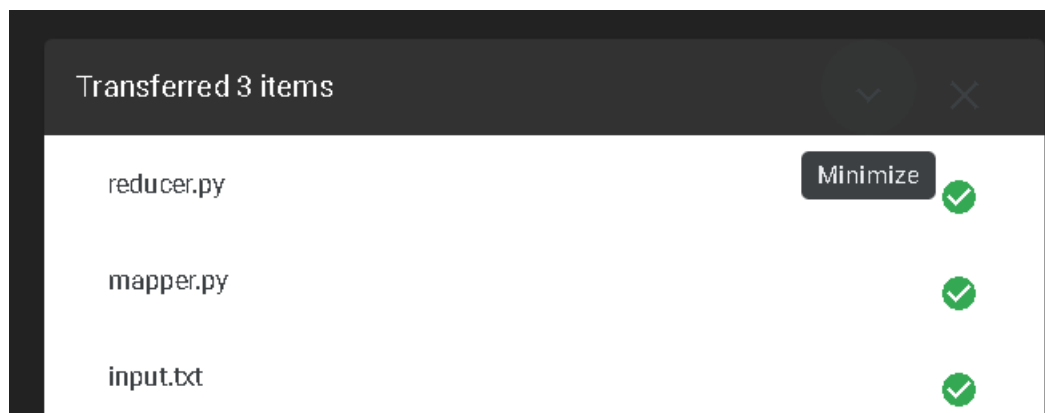
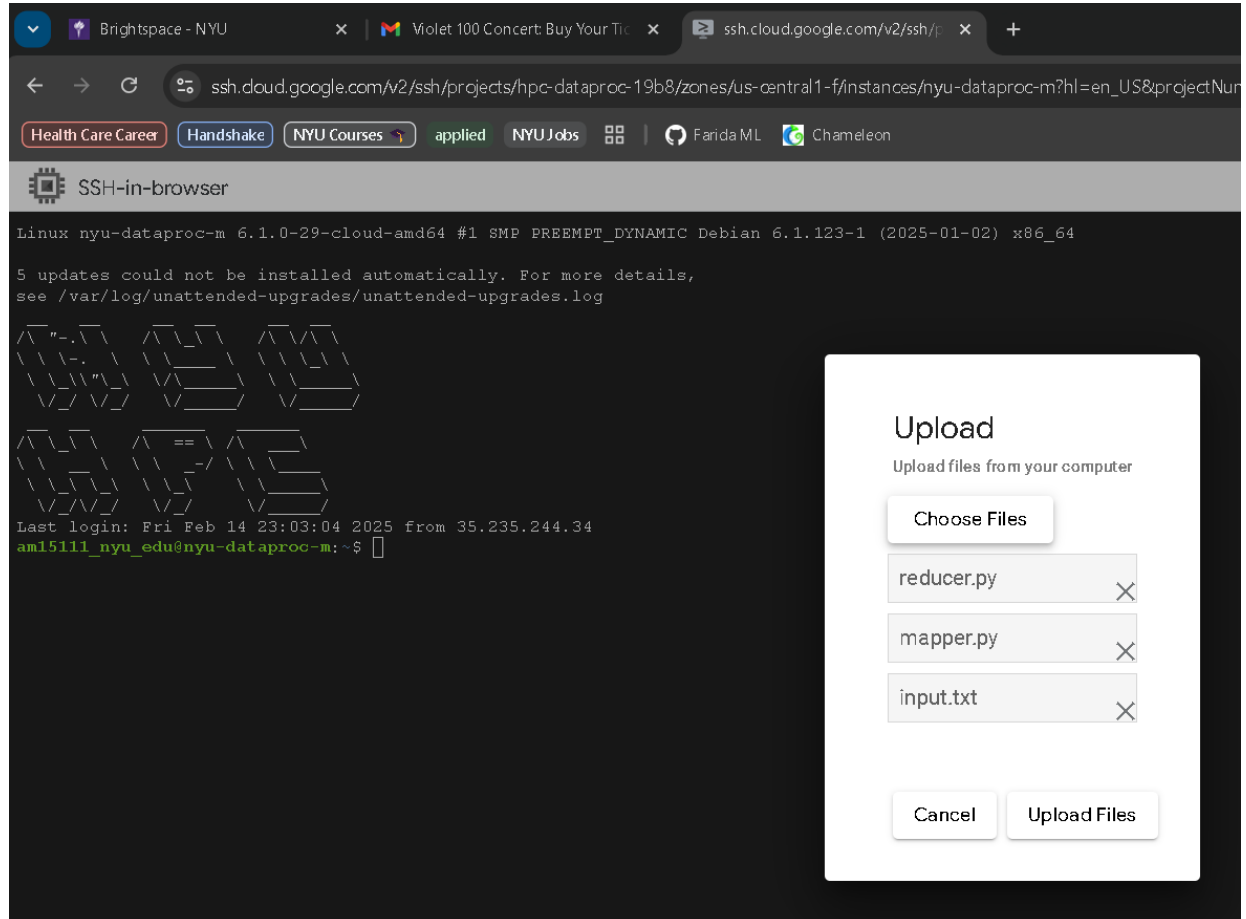


GHW#1- Big Data Assignment

Akshat Mishra | netid: am15111 | Spring 2025

Step 1: Upload the files to the NYU Dataproc local disk.



Step 2: Use “ls” command to check the presence of files in the directory.

```
am15111_nyu_edu@nyu-dataproc-m:~$ ls
input.txt  mapper.py  reducer.py
```

Step 3: Move the input to HDFS using the following command:

```
hdfs dfs -put input.txt
```

Perform “ls” on the HDFS using the following command and we should be getting the input.txt that we just moved.

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls
Found 1 items
-rw-r--r--  1 am15111_nyu_edu am15111_nyu_edu    2615571 2025-02-14 23:22 input.txt
```

Step 4: Inside Dataproc, all env variables like \$JAVA_HOME and \$HADOOP_HOME are already set for us. Listing the contents of the \$HADOOP_HOME directory, we get

```
am15111_nyu_edu@nyu-dataproc-m:~$ ls $HADOOP_HOME
bin                hadoop-azure-datalake.jar      hadoop-extras-3.3.6.jar      hadoop-registry.jar
client             hadoop-azure.jar              hadoop-extras.jar            hadoop-resourceestimator-3.3.6.jar
etc                hadoop-client-3.3.6.jar        hadoop-fs2img-3.3.6.jar      hadoop-resourceestimator.jar
hadoop-aliyun-3.3.6.jar  hadoop-client.jar            hadoop-fs2img.jar            hadoop-rumen-3.3.6.jar
hadoop-aliyun.jar    hadoop-common-3.3.6-tests.jar  hadoop-google-secret-manager-credential-provider-3.3.6.jar  hadoop-rumen.jar
hadoop-annotations-3.3.6.jar  hadoop-common-3.3.6.jar      hadoop-google-secret-manager-credential-provider.jar  hadoop-shaded-guava.jar
hadoop-annotations.jar  hadoop-common.jar            hadoop-gridmix-3.3.6.jar      hadoop-shaded-protobuf.jar
hadoop-archive-logs-3.3.6.jar  hadoop-datajoin-3.3.6.jar    hadoop-gridmix.jar            hadoop-sls-3.3.6.jar
hadoop-archive-logs.jar  hadoop-datajoin.jar          hadoop-kafka-3.3.6.jar        hadoop-sls.jar
hadoop-archives-3.3.6.jar  hadoop-distcp-3.3.6.jar      hadoop-kms-3.3.6.jar          hadoop-streaming-3.3.6.jar
hadoop-archives.jar     hadoop-distcp.jar            hadoop-kms.jar                hadoop-streaming.jar
hadoop-auth-3.3.6.jar    hadoop-dynamometer-blockgen-3.3.6.jar  hadoop-minicluster-3.3.6.jar  lib
hadoop-auth.jar        hadoop-dynamometer-blockgen.jar  hadoop-minicluster.jar        libexec
hadoop-aws-3.3.6.jar    hadoop-dynamometer-infra-3.3.6.jar  hadoop-nfs-3.3.6.jar          sbin
hadoop-aws.jar         hadoop-dynamometer-infra.jar      hadoop-nfs.jar
hadoop-azure-3.3.6.jar  hadoop-dynamometer-workload-3.3.6.jar  hadoop-registry-3.3.6.jar
hadoop-azure-datalake-3.3.6.jar  hadoop-dynamometer-workload.jar
```

Step 5: We will use the hadoop-streaming-3.2.2.jar file present inside this directory to execute a python based MapReduce as show with the below command:

```
hadoop jar $HADOOP_HOME/hadoop-streaming-3.3.6.jar -input input.txt -output outputpython
-mapper "python mapper.py" -reducer "python reducer.py" -file mapper.py -file reducer.py
```

```
am15111_nyu_edu@nyu-dataproc-m:~$ hadoop jar $HADOOP_HOME/hadoop-streaming-3.3.6.jar -input input.txt -output outputpython -mapper "python mapper.py" -reducer "python reducer.py" -file mapper.py -file reducer.py
2025-02-14 23:24:30,143 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [mapper.py, reducer.py] (/usr/lib/hadoop/hadoop-streaming-3.3.6.jar) /tmp/streamjob373956415624652333.jar tmpDir=null
2025-02-14 23:24:31,325 INFO client.DefaultHARMPFailoverProxyProvider: Connecting to ResourceManager at nyu-dataproc-m.local./192.168.1.93:8032
2025-02-14 23:24:31,504 INFO client.AHSProxy: Connecting to Application History server at nyu-dataproc-m.local./192.168.1.93:10200
2025-02-14 23:24:32,299 INFO client.DefaultHARMPFailoverProxyProvider: Connecting to ResourceManager at nyu-dataproc-m.local./192.168.1.93:8032
2025-02-14 23:24:32,300 INFO client.AHSProxy: Connecting to Application History server at nyu-dataproc-m.local./192.168.1.93:10200
2025-02-14 23:24:32,547 INFO mapreduce.JobResourceUploader: Disabling Brause Coding for path: /tmp/hadoop-yarn/staging/am15111_nyu_edu/.staging/job_1737068921170_1231
2025-02-14 23:24:32,987 INFO mapred.FileInputFormat: Total input files to process : 1
2025-02-14 23:24:33,051 INFO mapreduce.JobSubmitter: number of splits:6
2025-02-14 23:24:33,652 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1737068921170_1231
2025-02-14 23:24:33,653 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-02-14 23:24:33,808 INFO conf.Configuration: resource-types.xml not found
2025-02-14 23:24:33,808 INFO resource.ResourceUtil: Unable to find 'resource-types.xml'.
2025-02-14 23:24:33,907 INFO impl.YarnClientImpl: Submitted application application_1737068921170_1231
2025-02-14 23:24:33,940 INFO mapreduce.Job: The url to track the job: http://nyu-dataproc-m.local.:8088/proxy/application_1737068921170_1231/
2025-02-14 23:24:33,942 INFO mapreduce.Job: Running job: job_1737068921170_1231
2025-02-14 23:24:45,036 INFO mapreduce.Job: Job job_1737068921170_1231 running in uber mode : false
2025-02-14 23:24:45,037 INFO mapreduce.Job: map 0% reduce 0%
2025-02-14 23:24:56,131 INFO mapreduce.Job: map 17% reduce 0%
2025-02-14 23:24:57,138 INFO mapreduce.Job: map 100% reduce 0%
```

Step 6: Since the mentioned folder does not exist, This should trigger a job like this:

After successful completion of the job, you should see a ‘Job completed successfully’ message:

```
am15111_nyu_edu@nyu-dataproc-m:~$ hadoop jar (HADOOP_HOME)/hadoop-streaming-3.3.6.jar -input input.txt -output outputpython -mapper "python mapper.py" -reducer "python reducer.py" -file mapper.py -file reducer.py
2025-02-14 22:24:26.142 INFO streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packagedJobJar: (mapper.py, reducer.py) [(usr/lib/hadoop/hadoop-streaming-3.3.6.jar) (/tmp/steamm)sh77865641562452321.jar tmpbin=null
2025-02-14 22:24:26.222 INFO client.DefaultHadoopFileSystemDelegator: Connecting to ResourceManager at nyu-dataproc-m.local (/192.168.1.92:8032)
2025-02-14 22:24:26.564 INFO client.AMRProxy: Connecting to Application History server at nyu-dataproc-m.local (/192.168.1.92:10200)
2025-02-14 22:24:27.288 INFO client.DefaultHadoopFileSystemDelegator: Connecting to ResourceManager at nyu-dataproc-m.local (/192.168.1.92:8032)
2025-02-14 22:24:27.288 INFO client.AMRProxy: Connecting to Application History server at nyu-dataproc-m.local (/192.168.1.92:10200)
2025-02-14 22:24:27.547 INFO mapreduce.JobHadoopClosable: Disabling Resource Coding for path: /tmp/hadoop-nyu-dataproc-m15111_nyu_edu/_staging/job_1737060921170_1221
2025-02-14 22:24:27.897 INFO mapred.FileInputFormat: Total input files to process : 1
2025-02-14 22:24:27.951 INFO mapreduce.JobSubmitter: number of splits: 8
2025-02-14 22:24:27.952 INFO mapreduce.JobSubmitter: Submitting tasks for job: job_1737060921170_1221
2025-02-14 22:24:27.952 INFO mapreduce.JobSubmitter: Resolving with tasks: {}
2025-02-14 22:24:27.985 INFO conf.Configuration: resourceTypes.xml not found.
2025-02-14 22:24:27.985 INFO resource.ResourceUtil: Unable to find 'resourceTypes.xml'.
2025-02-14 22:24:27.987 INFO impl.TaskClientImpl: Submitted application application_1737060921170_1221
2025-02-14 22:24:27.986 INFO mapreduce.Job: The url to track the job: http://nyu-dataproc-m.local:8082/psmp/application_1737060921170_1221/
2025-02-14 22:24:27.982 INFO mapreduce.Job: Running job: job_1737060921170_1221
2025-02-14 22:24:44.616 INFO mapreduce.Job: job_1737060921170_1221 running in uber mode : false
2025-02-14 22:24:45.837 INFO mapreduce.Job: map 8% reduce 0%
2025-02-14 22:24:54.131 INFO mapreduce.Job: map 17% reduce 0%
2025-02-14 22:24:57.131 INFO mapreduce.Job: map 18% reduce 0%
2025-02-14 22:25:04.174 INFO mapreduce.Job: map 18% reduce 50%
2025-02-14 22:25:05.132 INFO mapreduce.Job: map 18% reduce 100%
2025-02-14 22:25:07.208 INFO mapreduce.Job: Job job_1737060921170_1221 completed successfully
2025-02-14 22:25:07.288 INFO mapreduce.Job: Counters: 55
File System Counters
  FILE: Number of bytes read=4335449
  FILE: Number of bytes written=11814320
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=3534872
  HDFS: Number of bytes written=4335465
  HDFS: Number of read operations=26
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=0
  HDFS: Number of bytes read across-racks=0
Job Counters
  Killed map tasks=1
  Launched map tasks=6
  Launched reduce tasks=2
  Map-local map tasks=6
  Total time spent by all maps in occupied slots (ms)=226204
  Total time spent by all reducers in occupied slots (ms)=48240
  Total time spent by all map tasks (ms)=58971
  Total time spent by all reduce tasks (ms)=18048
  Total volume-milliseconds taken by all map tasks=58971
  Total volume-milliseconds taken by all reduce tasks=18048
  Total mapbyte-milliseconds taken by all map tasks=241954216
  Total mapbyte-milliseconds taken by all reduce tasks=41205768
Map-Reduce Framework
  Map input records=42123
  Map output records=440718
  Map output bytes=3451899
  Map output materialized bytes=4335508
  Input split bytes=424
  Combine input records=0
  Combine output records=0
  Reduce input groups=35236
  Reduce shuffle bytes=4335508
  Reduce input records=460718
  Reduce output records=351516
  Spilled Records=321418
  Shuffled Maps=12
  Failed Shuffles=0
  Merged Map outputs=12
  CPU time elapsed (ms)=555
  CPU time spent (ms)=21560
  Physical memory (bytes) usage=211355616
  Virtual memory (bytes) usage=4265214848
  Total committed heap usage (bytes)=7802046726
  Peak Map Physical memory (bytes)=125915148
  Peak Map Virtual memory (bytes)=1017214976
  Peak Reduce Physical memory (bytes)=455076992
  Peak Reduce Virtual memory (bytes)=301673132
Shuffle Reads
  RAO_ID=0
  COMBINATION=0
  IO_EXCEPTION=0
  UNKNOWN_EXCEPTION=0
  UNKNOWN_JOB=0
  UNKNOWN_RESOURCE=0
File Input Format Counters
  FILE: Number of bytes read=4335449
  FILE: Number of bytes written=11814320
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=3534872
  HDFS: Number of bytes written=4335465
  HDFS: Number of read operations=26
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=0
  HDFS: Number of bytes read across-racks=0
Job Counters
  Killed map tasks=1
  Launched map tasks=6
  Launched reduce tasks=2
  Map-local map tasks=6
  Total time spent by all maps in occupied slots (ms)=226204
  Total time spent by all reducers in occupied slots (ms)=48240
  Total time spent by all map tasks (ms)=58971
  Total time spent by all reduce tasks (ms)=18048
  Total volume-milliseconds taken by all map tasks=58971
  Total volume-milliseconds taken by all reduce tasks=18048
  Total mapbyte-milliseconds taken by all map tasks=241954216
  Total mapbyte-milliseconds taken by all reduce tasks=41205768
Map-Reduce Framework
  Map input records=42123
  Map output records=440718
  Map output bytes=3451899
  Map output materialized bytes=4335508
  Input split bytes=424
  Combine input records=0
  Combine output records=0
  Reduce input groups=35236
  Reduce shuffle bytes=4335508
  Reduce input records=460718
  Reduce output records=351516
  Spilled Records=321418
  Shuffled Maps=12
  Failed Shuffles=0
  Merged Map outputs=12
  CPU time elapsed (ms)=555
  CPU time spent (ms)=21560
  Physical memory (bytes) usage=211355616
  Virtual memory (bytes) usage=4265214848
  Total committed heap usage (bytes)=7802046726
  Peak Map Physical memory (bytes)=125915148
  Peak Map Virtual memory (bytes)=1017214976
  Peak Reduce Physical memory (bytes)=455076992
  Peak Reduce Virtual memory (bytes)=301673132
Shuffle Reads
  RAO_ID=0
  COMBINATION=0
  IO_EXCEPTION=0
  UNKNOWN_EXCEPTION=0
  UNKNOWN_JOB=0
  UNKNOWN_RESOURCE=0
File Input Format Counters
```

Step 7: Listing the contents of the output directory and moving output files from HDFS to Dataproc using:

hdfs dfs -get outputpython

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls outputpython
Found 3 items
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 0 2025-02-14 22:33 outputpython/_SUCCESS
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 327750 2025-02-14 22:33 outputpython/part-00000
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 326118 2025-02-14 22:33 outputpython/part-00001
```

We will see all the files in dataproc local now using “ls outputpython”

```
am15111_nyu_edu@nyu-dataproc-m:~$ ls outputpython
_SUCCESS part-00000 part-00001
```

Step 8: For consolidating all the outputs within a single sorted result file (without using a single reducer), we could use a cat command with a wildcard:

```
hdfs dfs -cat outputpython/part* | sort > concatenated_results_new.tx
```

```
cp concatenated_file.txt .
```

The screenshot shows a web browser window with multiple tabs. The active tab is an SSH session titled 'ssh.cloud.google.com/v2/ssh/'. The browser's address bar shows the URL 'ssh.cloud.google.com/v2/ssh/projects/hpc-dataproc-19b8/zones/us-central1-f/instances/nyu-dataproc-m?hl=en_US&projectNumb'. Below the address bar, there are several buttons: 'Health Care Career', 'Handshake', 'NYU Courses', 'applied', 'NYU Jobs', 'Farida ML', and 'Chameleon'. The main content area is titled 'SSH-in-browser' and displays the output of an SSH session. The output includes memory usage statistics, shuffle errors, and file input/output format counters. A terminal window shows the command 'hdfs dfs -ls outputpython' and its output, which lists three files. A download dialog box is open on the right side of the screen, titled 'Download'. It shows the path '/home/am15111_nyu_ed' and has 'Cancel' and 'Download' buttons.

```
CPU time spent (ms)=21940
Physical memory (bytes) snapshot=5171355648
Virtual memory (bytes) snapshot=40045314048
Total committed heap usage (bytes)=7902068736
Peak Map Physical memory (bytes)=725741568
Peak Map Virtual memory (bytes)=5017214976
Peak Reduce Physical memory (bytes)=485076992
Peak Reduce Virtual memory (bytes)=5018673152

Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0

File Input Format Counters
  Bytes Read=2636051
File Output Format Counters
  Bytes Written=653868
2025-02-14 23:25:07,290 INFO streaming.StreamJob: Output directory: outputpython
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls outputpython
Found 3 items
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 0 2025-02-14 23:25:07 2025-02-14 23:25:07
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 327750 2025-02-14 23:25:07 2025-02-14 23:25:07
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 326118 2025-02-14 23:25:07 2025-02-14 23:25:07
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -get outputpython
```

Download

Path /home/am15111_nyu_ed

Cancel Download

Output file of the concatenated file

```
am15111_nyu_edu@nyu-dataproc-m:~$ cat concatenated_file.txt
!"ren<:h      1
!"('){}1I('.\.I. 1
!1('Xl 1
!71111111)I'11(' .l.)'      1
!<'Nl]7&'t',      1
!IIfI,'Ct&' , (l 1
!IamI—'lwilI 1
!Iu'_/1'tfuI.u,m.      1
!WL'7t/'V-(wt) 1
!hc'r0j'ore      1
!ilce_/isft,      1
!l'Ii.r 1
!lIlIlL'C('S.i'aIy      1
!lm'funereal      1
!t'G.)'t)fl      1
!'H$e      1
!'aIa7.:'.n      1
"      159
"      1
"1      1
":vl_v 1
"(. 'Izildmrt,. 1
"(x'r.su:zl      1
".      1
".w      1
"/ilrigfzr,      1
"0      1
"Of"mplir,'dIl1zr      1
"1      1
":1      1
":4      1
"A      1
"Abdul, 1
"Ali,      1
```

The Downloaded file is concatenated_file.txt.

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls outputph
Found 3 items
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 0 2025-02-15 01:15 outputph/_SUCCESS
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 88 2025-02-15 01:15 outputph/part-00000
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 88 2025-02-15 01:15 outputph/part-00001
```

Need to concatenate these two files and take top 10 from both of these files, so running this command gives sorted based on maximum occurring frequency and saving it in file top_10_frequent_words.txt :

```
hdfs dfs -cat outputph/part-* | sort -k2,2nr | head -n 10 > top_10_frequent_words.txt
```

With the use of : **cat top_10_frequent_words.txt** we can view the file contents.

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat outputph/part-* | sort -k2,2nr | head -n 10 > top_10_frequent_words.txt
am15111_nyu_edu@nyu-dataproc-m:~$ ls
concatenated_file.txt  input.txt  mapper.py  mappertt.py  outputpython  reducer.py  reducertt.py  top_10_frequent_words.txt
am15111_nyu_edu@nyu-dataproc-m:~$ cat top_10_frequent_words.txt
the      21499
to       13071
and      12126
of       11301
a        8964
in       6630
you      4992
that    4545
is       4380
he       4221
```

Step 5: Download the top_10_frequent _words.txt and Now directory structure looks like this:

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls
Found 3 items
-rw-r--r--  1 am15111_nyu_edu am15111_nyu_edu    2615571 2025-02-14 23:22 input.txt
drwxr-xr-x  - am15111_nyu_edu am15111_nyu_edu         0 2025-02-15 01:15 outputph
drwxr-xr-x  - am15111_nyu_edu am15111_nyu_edu         0 2025-02-14 23:25 outputpython
```

The Downloaded file is present in the Folder. (top_10_frequent_words.txt)

C. Create a new text file which is somewhere between 12-13mb, but in order to achieve this, copy and paste a small portion of the previously created text file (like around the first 250-300 lines, not more) multiple times till you reach the desired size. Perform the same programs (a) and (b) on this file.

Step1: Uploading the text file named “duplicated_input.txt”. Directory structure after uploading this file:

```
am15111_nyu_edu@nyu-dataproo-m:~$ ls
concatenated_file.txt  duplicated_input.txt  input.txt  mapper.py  mappert.py  outputpython  reducer.py  reducertt.py  top_10_frequent_words.txt
am15111_nyu_edu@nyu-dataproo-m:~$ hdfs dfs -put duplicated_input.txt
```

Step 2: uploading this file to hdfs, command:

```
hdfs dfs -put duplicated_input.txt
```

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -put duplicated_input.txt
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls
Found 4 items
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 13521879 2025-02-15 01:36 duplicated_input.txt
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 2615571 2025-02-14 23:22 input.txt
drwxr-xr-x - am15111_nyu_edu am15111_nyu_edu 0 2025-02-15 01:15 outputph
drwxr-xr-x - am15111_nyu_edu am15111_nyu_edu 0 2025-02-14 23:25 outputpython
```

Step 3: running the mapper.py and reducer.py to get the word count for a. part:

```
mapred streaming -input duplicated_input.txt -output duplicate_op1 -mapper "python3
mapper.py" -reducer "python3 reducer.py" -file mapper.py -file reducer.py
```

```

aws1111_nyu_edu@nyu-dataproc-m:~$ mapred streaming -input duplicated_input.txt -output duplicate_op1 -mapper "python3 mapper.py" -reducer "python3 reducer.py" -file mapper.py -file reducer.py
WARN: HADOOP_JOB_HISTORYSERVER_OPTS has been replaced by MAPRED_HISTORYSERVER_OPTS. Using value of HADOOP_JOB_HISTORYSERVER_OPTS.
2025-02-15 01:39:00,935 WARN streaming.StreamJob: file option is deprecated, please use generic option files instead.
packageJobJar: [mapper.py, reducer.py] (/usr/lib/hadoop/hadoop-streaming-3.9.6.jar) (/tmp/stagingjob84755962134225562170.jar tmpDir=null)
2025-02-15 01:39:01,468 INFO client.AHSProxy: Connecting to Application History server at nyu-dataproc-m.local:192.168.1.93:8032
2025-02-15 01:39:02,398 INFO client.AHSProxy: Connecting to Application History server at nyu-dataproc-m.local:192.168.1.93:10200
2025-02-15 01:39:03,388 INFO client.DefaultHARMAFailureProxyProvider: Connecting to ResourceManager at nyu-dataproc-m.local:192.168.1.93:8032
2025-02-15 01:39:03,389 INFO client.AHSProxy: Connecting to Application History server at nyu-dataproc-m.local:192.168.1.93:10200
2025-02-15 01:39:03,643 INFO mapreduce.JobSubmitter: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/aw1111_nyu_edu/staging/job_1737068921170_1358
2025-02-15 01:39:04,041 INFO mapred.FileInputFormat: Total input files to process : 1
2025-02-15 01:39:04,168 INFO mapreduce.JobSubmitter: number of splits=6
2025-02-15 01:39:04,315 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1737068921170_1358
2025-02-15 01:39:04,315 INFO mapreduce.JobSubmitter: Executing with tokens: {}
2025-02-15 01:39:04,489 INFO conf.Configuration: resource-types.xml not found
2025-02-15 01:39:04,489 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2025-02-15 01:39:04,611 INFO impl.YarnClientImpl: Submitted application application_1737068921170_1358
2025-02-15 01:39:04,648 INFO mapreduce.Job: The url to track the job is http://nyu-dataproc-m.local:8080/proxy/application_1737068921170_1358/
2025-02-15 01:39:04,648 INFO mapreduce.Job: Running job: job_1737068921170_1358

```

[illegible]

Step 4: Make sure that file is being saved in the hdfs directory:

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls
Found 5 items
drwxr-xr-x - am15111_nyu_edu am15111_nyu_edu 0 2025-02-15 01:39 duplicate_op1
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 13521879 2025-02-15 01:36 duplicated_input.txt
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 2615571 2025-02-14 23:22 input.txt
drwxr-xr-x - am15111_nyu_edu am15111_nyu_edu 0 2025-02-15 01:15 outputph
drwxr-xr-x - am15111_nyu_edu am15111_nyu_edu 0 2025-02-14 23:25 outputpython

am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls duplicate_op1
Found 3 items
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 0 2025-02-15 01:39 duplicate_op1/_SUCCESS
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 329910 2025-02-15 01:39 duplicate_op1/part-00000
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 328335 2025-02-15 01:39 duplicate_op1/part-00001
```

Step 5: Now concatenating it and saving the text outside of hdfs and downloading it.

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat duplicate_op1/part* | sort > concatenated_file_duplicate.txt
am15111_nyu_edu@nyu-dataproc-m:~$ ls
concatenated_file.txt concatenated_file_duplicate.txt duplicated_input.txt input.txt mapper.py mappertt.py outputpython reducer.py reducertt.py top_10_frequent_words.txt
```

```
am15111_nyu_edu@nyu-dataproc-m:~$ cp concatenated_file_duplicate.txt .
cp: 'concatenated_file_duplicate.txt' and './concatenated_file_duplicate.txt' are the same file
am15111_nyu_edu@nyu-dataproc-m:~$ pwd
/home/am15111_nyu_edu
```

```
am15111_nyu_edu@nyu-dataproc-m:~$ cat concatenated_file_duplicate.txt
!"ren<:h 1
!'('()1I('.\.I. 1
!l('Xl 1
!7l1l1l1l)I'1l('l.)' 1
!<'Nl]7&'t' 1
!IIfI,'Ct$',(l 1
!IamI—'lwi1I 1
!Iu'/_l'tfuI.u,m. 1
!WL'7t/'V-(wt) 1
!hc'r0j'ore 1
!ilce_/isft, 1
!l'Ii.r 1
!lIlIlL'C('S.i'aIy 1
!lm'funereal 1
!t'G.)'t)fl 1
!'H$€ 1
!'aIa7.:'.n 1
" 159
" ' 1
"1 1
"! :vl_v 1
"(. 'Izildmrt,. 1
"(x' r.su:zl 1
". 1
".w 1
"/ilrigfzr, 1
"0 1
"Of"mplir,'dIl1zr 1
"1 1
":1 1
":4 1
"A 1
"Abdul, 1
"Ali, 1
"And 2
"April 2
```

The Downloaded file is concatenated_file_duplicate.txt.

Now doing b part running the mapper.py and reducer.py

Step 6: Command:

```
mapred streaming -input duplicate_op1/part-* -output duplicate_tt -mapper "python3 mapper.py" -reducer "python3 reducer.py" -file mapper.py -file reducer.py
```

```
am15111_nyu_edu@nyu-dataproc-m:~$ mapred streaming -input duplicate_op1/part-* -output duplicate_tt -mapper "python3 mapper.py" -reducer "python3 reducer.py" -file mapper.py -file reducer.py
WARNING: HADOOP_JOB_HISTORYSERVER_OPTS has been replaced by MAPRED_HISTORYSERVER_OPTS. Using value of HADOOP_JOB_HISTORYSERVER_OPTS.
2025-02-15 01:55:55.302 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [mapper.py, reducer.py] /usr/lib/hadoop/hadoop-streaming-3.2.6.jar /tmp/streamjob64559360706783975.jar tmpDir=null
2025-02-15 01:55:56.462 INFO client.DefaultNoHARMFalloverProxyProvider: Connecting to ResourceManager at nyu-dataproc-m.local./192.168.1.93:8032
2025-02-15 01:55:56.623 INFO client.AHSProxy: Connecting to Application History server at nyu-dataproc-m.local./192.168.1.93:10200
2025-02-15 01:55:57.377 INFO client.DefaultNoHARMFalloverProxyProvider: Connecting to ResourceManager at nyu-dataproc-m.local./192.168.1.93:8032
2025-02-15 01:55:57.378 INFO client.AHSProxy: Connecting to Application History server at nyu-dataproc-m.local./192.168.1.93:10200
2025-02-15 01:55:57.639 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/am15111_nyu_edu/.staging/job_1737068921170_1381
2025-02-15 01:55:58.828 INFO mapred.FileInputFormat: Total input files to process : 2
2025-02-15 01:55:58.881 INFO mapreduce.JobSubmitter: number of splits:6
2025-02-15 01:55:59.046 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1737068921170_1381
2025-02-15 01:55:59.046 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-02-15 01:55:59.195 INFO conf.Configuration: resource-types.xml not found
2025-02-15 01:55:59.195 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2025-02-15 01:55:59.258 INFO impl.YarnClientImpl: Submitted application application_1737068921170_1381
2025-02-15 01:55:59.303 INFO mapreduce.Job: The url to track the job: http://nyu-dataproc-m.local.:8088/proxy/application_1737068921170_1381/
2025-02-15 01:55:59.305 INFO mapreduce.Job: Running job: job_1737068921170_1381
2025-02-15 01:56:15.411 INFO mapreduce.Job: Job job_1737068921170_1381 running in uber mode : false
2025-02-15 01:56:13.412 INFO mapreduce.Job: map 0% reduce 0%
2025-02-15 01:56:21.493 INFO mapreduce.Job: map 17% reduce 0%
2025-02-15 01:56:25.516 INFO mapreduce.Job: map 100% reduce 0%
2025-02-15 01:56:32.566 INFO mapreduce.Job: map 100% reduce 50%
2025-02-15 01:56:35.584 INFO mapreduce.Job: map 100% reduce 100%
2025-02-15 01:56:37.603 INFO mapreduce.Job: Job job_1737068921170_1381 completed successfully
2025-02-15 01:56:37.694 INFO mapreduce.Job: Counters: 56
  File System Counters
    FILE: Number of bytes read=775969
    FILE: Number of bytes written=3895600
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=675243
    HDFS: Number of bytes written=191
    HDFS: Number of read operations=28
    HDFS: Number of large read operations=0
```

Step7: Since there are two files as seen using command: **hdfs dfs -ls duplicate_tt**

Need to concatenate these two files and take top 10 from both of these files, so running this command gives sorted based on maximum occurring frequency and saving it in file top_10_frequent_words.txt :

```
hdfs dfs -cat duplicate_tt/part-* | sort -k2,2nr | head -n 10 >
```

```
duplicate_top_10_frequent_words.txt
```

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -ls duplicate_tt
Found 3 items
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 0 2025-02-15 01:56 duplicate_tt/_SUCCESS
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 92 2025-02-15 01:56 duplicate_tt/part-00000
-rw-r--r-- 1 am15111_nyu_edu am15111_nyu_edu 99 2025-02-15 01:56 duplicate_tt/part-00001
```

```
am15111_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat duplicate_tt/part-* | sort -k2,2nr | head -n 10 > duplicate_top_10_frequent_words.txt
am15111_nyu_edu@nyu-dataproc-m:~$ ls
concatenated_file.txt      duplicate_top_10_frequent_words.txt  input.txt  mapper.py  reducer.py  top_10_frequent_words.txt
concatenated_file_duplicate.txt  duplicated_input.txt                  mapper.py  outputpython  reducer.py
```

Step7: Seeing the content of top 10 occurring words:

```
am15111_nyu_edu@nyu-dataproc-m:~$ cat duplicate_top_10_frequent_words.txt
the      103132
to       70998
a        63932
of       59982
and      47302
you      41622
is       33059
I        32594
in       26247
your     26059
```

The Downloaded file is duplicated_top_10_frequent_words.txt.

Q1. What difference do you notice, if any, in the output directory after running the normal word count program on the modified text file as compared to the when you ran it with the originally created text file?

The output directory shows a difference in word frequency when running the normal word count program on the modified text file compared to the original 2 MB file. Since 200-300 lines were duplicated multiple times to generate the 13 MB file, certain words appear more frequently in the modified version. As a result, both directories contain only a few other differences beyond this frequency variation.

Q2. Why do you think there is a difference in the output while running word count despite the input files being of similar size?

The difference in the output occurs because the modified text file was created by duplicating 200-300 lines multiple times, rather than generating entirely new content. This repetition causes certain words to appear with disproportionately higher frequencies compared to the original file, even though both files are of similar size. The normal word count program processes each word as it appears in the file, leading to a skewed distribution in the modified file, where some words dominate the frequency count.

Q3. What difference do you notice, if any, among the top 10 most frequently occurring words between the two files? Why do you think this is the case?

The top 10 most frequently occurring words differ between the 3 MB file and the 13 MB file due to the duplication of 200-300 lines multiple times in the larger file. The key differences observed are:

1. Increased Frequency for Common Words – Words like "the", "to", "a", and "of" appear much more frequently in the 13 MB file, reflecting the repeated content.
2. Ranking Changes – Some words, such as "he", "that", and "in", which were in the top 10 for the smaller file, are replaced by "I" and "your" in the larger file. This suggests that the duplicated lines contain these words more often.
3. Skewed Distribution – Since the larger file was created by duplication rather than unique content, certain words dominate the frequency count, altering the natural distribution seen in the 3 MB file.

Learnings from the Experiment:

1. **Impact of Data Duplication** – Repeating small portions of text inflates word frequency, altering natural distributions despite similar file sizes.
2. **Importance of Unique Data** – Larger file size doesn't mean more diverse data; redundancy skews analysis results.
3. **MapReduce Behavior** – It processes all words as they appear, making input structure crucial for accurate output.
4. **Changes in Word Rankings** – Frequently repeated words dominate the top 10 in the larger file.
5. **Efficient Hadoop Usage** – Reinforced HDFS and MapReduce operations for handling large datasets.
6. **Sorting & Aggregation** – Helps consolidate results for better analysis.