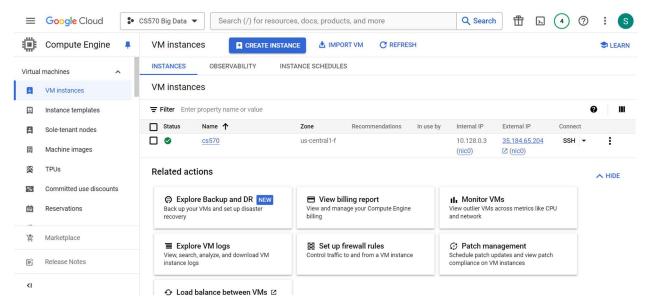
Week 3: Homework 1: Project: Creating MapReduce program to calculating Pi

GCP instance:



1. Creating folder Pi calculations and program to generate random numbers in (x,y) format.

```
ssh.cloud.google.com/v2/ssh/projects/cs570-big-data-424809/zones/us-central1-f/instances/cs570?authuser=0&hl=en_US&projectNumb...

ssh.cloud.google.com/v2/ssh/projects/cs570-big-data-424809/zones/us-central1-f/instances/cs570?authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-big-data-424809/zones/us-central1-f/instances/cs570?authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-big-data-424809/zones/us-central1-f/instances/cs570?authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-big-data-424809/zones/us-central1-f/instances/cs570?authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-authuser=0&hl=en_US&sh.cloud.google.com/v2/ssh/projects/cs570-authuser=0&hl=en_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shleen_US&shle
```

The program to generate random numbers:

Generate an input file to the Pi MapReduce program and compiling the program:

```
shagos90499@cs570:~/Pi_Calculations$ vi GenerateRandomNumbers.java
shagos90499@cs570:~/Pi_Calculations$ java -cp . GenerateRandomNumbers
How many random numbers to generate:
1000000
What's the radius?
200
shagos90499@cs570:~/Pi_Calculations$ ls
GenerateRandomNumbers.class GenerateRandomNumbers.java PiCalculationInput
shagos90499@cs570:~/Pi_Calculations$
```

Random numbers generated:

7,297) (265,396) (236,33) (159,310) (323,287) (56,67) (102,8) (303,198) (335,34) (16,27) (260,1) (348,287) (324,282) (89,392) (331,286) (134,345) (347,334) (106,328) (330,525) (170,333) (246,345) (226,283) (295,181) (133,11) (265,139) (106,392) (135,281) (306,275) (285,54) (272,16) (329,149) (204,310) (348,285) (275,244) (127,319) (267,75) (362,87) (46,99) (111,71) (311,78) (318,95) (162,287) (322,152) (39,351) (310,340) (320,66) (104,378) (144,399) (236,113) (366,141) (68,190) (171,35) (325,339) (225,26) (197,313) (38,376) (144,245) (29,180) (92,199) (55,979) (235,155) (146,442) (22,357) (250,132) (234,60) (185,377) (99,179) (341,368) (270,267) (125,269) (358,204) (365,183) (165,330) (122,18) (94,104) (87,351) (376,24) (344,352) (233,80) (185,377) (99,179) (341,368) (270,267) (12,269) (90,200) (266,373) (388,27) (313,320) (117,24) (179,169) (4,287) (144,11) (296,169) (352,302) (353,124) (316,755) (120,731) (13,16) (253,65) (273,113) (336,2) (211,121) (206,272) (144,235) (114,366) (112,337) (184,257) (187,135) (76,341) (342,341) (42,212) (284,75) (211,224) (94,65) (139,152) (352,373) (257,222) (356,76) (91,158) (355,822) (174,100) (156,150) (203,9) (505,11) (382,99) (382,00) (171,134) (66,157) (305,322) (268,181) (150,44) (212,191) (257,382) (113,130) (198,256) (10,399) (267,103) (321,49) (187,262) (285,155) (41,159) (121,224) (223,382) (203,137) (13,377) (188,166) (53,324) (113,130) (198,256) (10,399) (267,103) (321,49) (187,262) (285,155) (41,159) (122,238) (203,137) (13,377) (188,166) (33,341) (10,76) (49,196) (123,298) (803,137) (13,377) (188,166) (33,447) (117,60) (383,147) (117,700) (383,1) (366,49) (198,266) (23,296) (374,221) (22,388) (203,137) (13,377) (188,166) (33,447) (117,60) (383,14) (366,49) (198,266) (23,296) (374,221) (22,388) (203,137) (13,377) (188,166) (33,347) (197,248) (188,186) (33,347) (197,248) (188,186) (33,347) (197,248) (188,186) (33,347) (197,248) (188,186) (33,347) (197,248) (188,186) (33,348) (188,186) (33,348) (188,186) (33,348) (188,186) (33,348) (188,186) (33,348) (188,186) (33

Now lets go back to Hadoop and see if we can connect to the localhost:

```
shagos90499@cs570:~$ pwd
/home/shagos90499
shagos90499@cs570:~$ ls
Pi Calculations hadoop-3.4.0 hadoop-3.4.0.tar.gz
shagos90499@cs570:~$ cd hadoop-3.4.0/
shagos90499@cs570:~/hadoop-3.4.0$ ssh localhost
Welcome to Ubuntu 20.04.6 LTS (GNU/Linux 5.15.0-1060-gcp x86 64)
 * Documentation: https://help.ubuntu.com
                  https://landscape.canonical.com
 * Management:
 * Support:
                  https://ubuntu.com/pro
 System information as of Mon Jun 3 09:47:33 UTC 2024
  System load:
               0.0
                                                         109
                                  Processes:
               53.7% of 9.51GB Users logged in:
 Usage of /:
                                                         1
                                 IPv4 address for ens4: 10.128.0.3
 Memory usage: 8%
 Swap usage:
Expanded Security Maintenance for Applications is not enabled.
15 updates can be applied immediately.
7 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable
Enable ESM Apps to receive additional future security updates.
See https://ubuntu.com/esm or run: sudo pro status
New release '22.04.3 LTS' available.
Run 'do-release-upgrade' to upgrade to it.
Last login: Mon Jun 3 09:20:39 2024 from 127.0.0.1
shagos90499@cs570:~$
```

Next, we will create the HDFS directories required to execute MapReduce jobs:

```
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs dfs -mkdir /user
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs dfs -mkdir /user/shagos90499
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs dfs -mkdir /user/shagos90499/picalculate
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs dfs -mkdir /user/shagos90499/picalculate/input
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs dfs -put ../Pi_Calculations/PiCalculationInput /user/shagos90499/picalculate/input
shagos90499@cs570:~/hadoop-3.4.0$
```

Step 2: create a MapReduce program to calculate the numbers of inside darts and outside darts.

```
shagos90499@cs570:~/hadoop-3.4.0$ vi PiCalculation.java
shagos90499@cs570:~/hadoop-3.4.0$ bin/hadoop com.sun.tools.javac.Main PiCalculation.java
shagos90499@cs570:~/hadoop-3.4.0$ jar cf wc.jar PiCalculation*class
shagos90499@cs570:~/hadoop-3.4.0$ ls
LICENSE-binary 'PiCalculation;IntSumReducer.class' LICENSE.txt 'PiCalculation;TokenizerMapper.class'
                                                            README.txt
                                                                          index.html
                                                                                                           share
                                                                          input
NOTICE-binary PiCalculation.class
                                                                                        output
                  PiCalculation.java
NOTICE.txt
                                                            include
                                                                          libexec
                                                                                        sbin
shagos90499@cs570:~/hadoop-3.4.0$
```

Step 3: Use the file generated in Step 1.2 as the input to execute the MapReduce program created in Step 2:

```
2024-06-03 10:32:44,608 INFO impl.MetricsSording: Loaded properties from hadoop-metrics2.properties
2024-06-03 10:32:44,603 INFO impl.MetricsSording: Loaded properties from hadoop-metrics2.properties
2024-06-03 10:32:44,603 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2024-06-03 10:32:44,603 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2024-06-03 10:32:44,603 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2024-06-03 10:32:44,605 MARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execut e your application with ToolRunner to remedy this.
2024-06-03 10:32:45,255 INFO input.FileInputFormat: Total input files to process: 1
2024-06-03 10:32:45,293 INFO mapreduce.JobSubmitter: number of splits;
2024-06-03 10:32:45,498 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1485709741_0001
2024-06-03 10:32:45,804 INFO mapreduce.JobSubmitter: Executing with tokens: []
2024-06-03 10:32:45,807 INFO mapreduce.Job: Running job: job_local1485709741_0001
2024-06-03 10:32:45,807 INFO mapreduce.Job: Running job: job_local1485709741_0001
2024-06-03 10:32:45,807 INFO mapreduce.JobCunner: OutputCommitter for in config null
2024-06-03 10:32:45,804 INFO output.PathOutputCommitter: File Output Committer factory defined, defaulting to FileOutputCommitterFactory
2024-06-03 10:32:45,805 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-06-03 10:32:45,805 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2024-06-03 10:32:45,855 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
```

Step 4: Calculate Pi in the driver program based on the numbers of inside darts and outside darts

Output:

Test results:

```
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=9449789
File Output Format Counters
Bytes Written=29
inside 785377
outside 214623
Inside:785377, Outside:214623
PI:3.141508
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs
```

Stop machine:

```
shagos90499@cs570:~/hadoop-3.4.0$ sbin/stop-dfs.sh
Stopping namenodes on [localhost]
Stopping datanodes
Stopping secondary namenodes [cs570]
shagos90499@cs570:~/hadoop-3.4.0$
```

Detailed result:

```
ຊ ssh.cloud.google.com/v2/ssh/projects/cs570-big-data-424809/zones/us-central1-f/instances/cs57... 👚
 ssh.cloud.google.com/v2/ssh/projects/cs570-big-data-424809/zones/us-central1-f/instances/cs...
                                                                                                                                                               Ø
  SSH-in-browser
                                                                          ↑ UPLOAD FILE
                                                                                                         HDFS: Number of read operations=15
HDFS: Number of large read operations=0
HDFS: Number of write operations=4
HDFS: Number of bytes read erasure-coded=0
            Map-Reduce Framework
                        Map input records=1
                        Map output records=1000000
Map output bytes=11214623
                        Map output materialized bytes=33
                        Input split bytes=140
Combine input records=1000000
                        Combine output records=2
Reduce input groups=2
Reduce shuffle bytes=33
                        Reduce input records=2
Reduce output records=2
                        Spilled Records=4
                        Shuffled Maps =1
                        Failed Shuffles=0
                        Merged Map outputs=1
GC time elapsed (ms)=568
Total committed heap usage (bytes)=1191182336
            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=9449789
File Output Format Counters
                        Bytes Written=29
inside 785377
Inside:785377, Outside:214623
PI:3.141508
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs dfs -ls /user/lchen/picalculate/output5
ls: `/user/lchen/picalculate/output5': No such file or directory
shagos90499@cs570:~/hadoop-3.4.0$ bin/hdfs dfs -ls /user/shagos90499/picalculate/output5
```

Shut down VM:

