# MapReduce in Calculating Pi

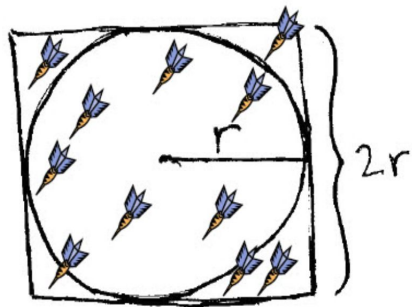Linh Bien
Professor: Henry Chang

# Table of content

- Introduction
- Design
- Implementation
- Test
- Enhancement ideas
- Conclusion
- Reference

# Introduction

In this project, the Hadoop environment is used to calculate Pi

# Design

- Throw $N$ darts on the board. Each dart lands at a random position $(x,y)$ on the board.
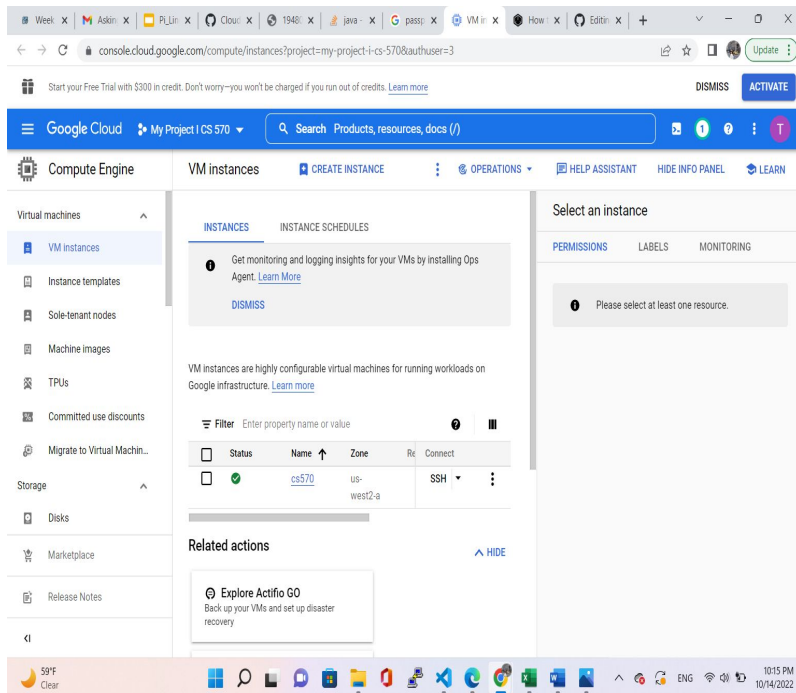
  - Note if each dart landed inside the circle or not
    - Check if $x^2+y^2<r$
  - Take the total number of darts that landed in the circle as $S$
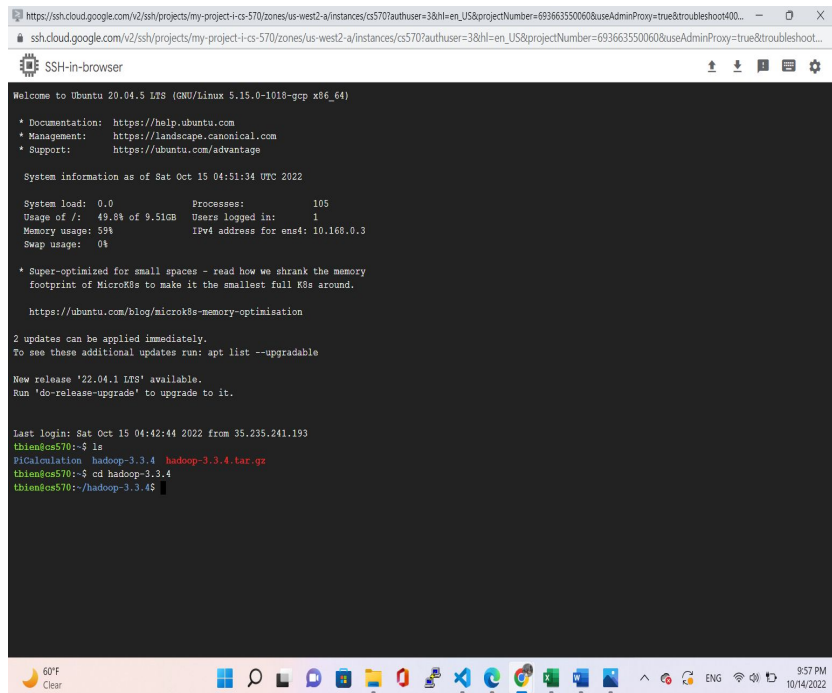


$$4\left(\frac{S}{N}\right) = \pi$$

**Formula:**

$4 * S / N = 4 * (pi * r * r) / (4 * r * r) = pi$

# Implementation

# Implementation

```
$ mkdir PiCalculation1
```

```
tbien@cs570:~/hadoop-3.3.4$ cd PiCalculation
tbien@cs570:~/hadoop-3.3.4/PiCalculation$ vi GenerateRandomNumbers.java
tbien@cs570:~/hadoop-3.3.4/PiCalculation$ javac GenerateRandomNumbers.java
tbien@cs570:~/hadoop-3.3.4/PiCalculation$ java -cp . GenerateRandomNumbers
How many random numbers to generate:
1000000
What's the radius?
200
```

# Implementation

Make HDFS directory

```
$ bin/hdfs dfs -mkdir /user

$ bin/hdfs dfs -mkdir /user/tbien

$ bin/hdfs dfs -mkdir /user/tbien/picalculate

$ bin/hdfs dfs -mkdir /user/tbien/picalculate/input2

$ bin/hdfs dfs -put ../PiCalculation1/PiCalculationInput /user/tbien/picalculate/input2
```

# Implementation

# Implementation - MapReduce program

Create Mapreduce to calculate number of inside and outside darts

```
tbien@cs570:~/hadoop-3.3.4$ vi PiCalculation.java
tbien@cs570:~/hadoop-3.3.4$ bin/hadoop com.sun.tools.javac.Main PiCalculation.java
tbien@cs570:~/hadoop-3.3.4$ jar cf wc.jar PiCalculation*class
tbien@cs570:~/hadoop-3.3.4$ ls
LICENSE-binary      PiCalculation                      PiCalculation.java      bin       input    licenses-binary   sbin
LICENSE.txt         'PiCalculation$IntSumReducer.class'      PiCalculation1     etc       input1   logs              share
NOTICE-binary       'PiCalculation$TokenizerMapper.class'    PiCalculation1.java  include   lib      output            wc.jar
NOTICE.txt          PiCalculation.class                README.txt              index.html libexec output1
tbien@cs570:~/hadoop-3.3.4$
```

# Implementation



```
tbien@cs570:~/hadoop-3.3.4$ bin/hadoop jar wc.jar PiCalculation /user/tbien/picalculate/input /user/tbien/picalculate/output3
2022-10-15 23:31:37,192 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2022-10-15 23:31:37,356 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2022-10-15 23:31:37,357 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2022-10-15 23:31:37,620 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute yo
ur application with ToolRunner to remedy this.
2022-10-15 23:31:37,780 INFO input.FileInputFormat: Total input files to process : 0
2022-10-15 23:31:37,790 INFO mapreduce.JobSubmitter: number of splits:0
2022-10-15 23:31:38,051 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1021437452_0001
2022-10-15 23:31:38,052 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-10-15 23:31:38,290 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2022-10-15 23:31:38,291 INFO mapreduce.Job: Running job: job_local1021437452_0001
2022-10-15 23:31:38,308 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2022-10-15 23:31:38,321 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2022-10-15 23:31:38,321 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup
 failures: false
2022-10-15 23:31:38,323 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2022-10-15 23:31:38,392 INFO mapred.LocalJobRunner: Waiting for map tasks
2022-10-15 23:31:38,392 INFO mapred.LocalJobRunner: map task executor complete.
2022-10-15 23:31:38,400 INFO mapred.LocalJobRunner: Waiting for reduce tasks
2022-10-15 23:31:38,401 INFO mapred.LocalJobRunner: Starting task: attempt_local1021437452_0001_r_000000_0
2022-10-15 23:31:38,450 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
```

# Test

```
        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=9450262
        File Output Format Counters
                Bytes Written=29
inside  784816
outside 215184
Inside:784816, Outside:215184
PI:3.139264
```

# Enhancement idea

Test more on larger numbers then we will have more accurate results

# Conclusion

It is important to write right codes and run the correct results

Test on larger numbers to have more accurate results

# References

Research Gate. https://www.researchgate.net/figure/MapReduce-calculation-process_fig2_359948761

Overview of Pi Calculation using MapReduce