# Pi Calculation Using MapReduce & PySpark

Natnael Haile
ID: 20007
June 19, 2024

# Table of Contents
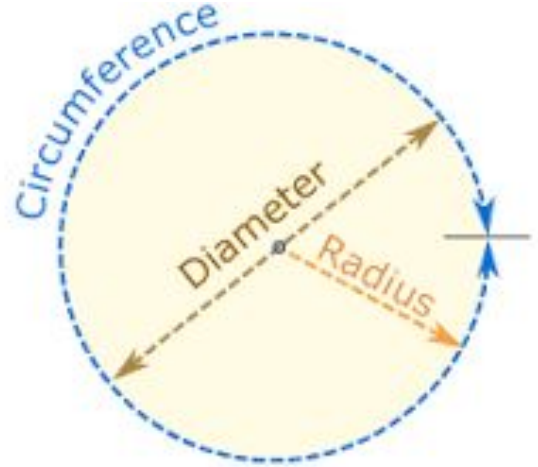
# Introduction

- **Our primary objective in this project is to calculate the value of Pi using two distinct distributed computing frameworks:**
  - **Hadoop MapReduce and**
  - **PySpark.**
- **While Hadoop is typically not the first choice for computationally intensive tasks, this exercise serves as a practical application to understand and utilize Hadoop's capabilities.**
- **Additionally, we will leverage PySpark to illustrate its efficiency and versatility in handling similar tasks.**
- **By undertaking this project, we will explore and compare the frameworks' abilities to handle large-scale data processing and computational tasks.**

$$\frac{Circumference}{Diameter} = \pi = 3.14159\ldots$$

# Introduction cont.

**Why Calculate Pi?**

- Calculating Pi is a well-known problem in computer science and mathematics, offering a straightforward way to demonstrate the power of distributed computing.
- The simplicity of the Monte Carlo method for estimating Pi makes it an ideal candidate for this exercise.
- It involves random sampling, which can be easily parallelized, providing a clear illustration of the frameworks' data processing and computational capabilities.



$$\frac{Circumference}{Diameter} = \pi = 3.14159\ldots$$

# Introduction cont.

**Hadoop MapReduce Program**

- **Purpose: Demonstrate Hadoop's framework in a computational context.**
- **Methodology: Use the Monte Carlo method to estimate the value of Pi through a series of parallel computations distributed across a Hadoop cluster.**
- **Outcome: Gain hands-on experience with Hadoop's MapReduce model and its potential for handling large-scale data processing tasks.**

**PySpark Program**

- **Purpose: Highlight PySpark's efficiency and ease of use in distributed computing.**
- **Methodology: Implement the Monte Carlo method to estimate Pi, leveraging PySpark's in-memory processing capabilities and simplified API.**
- **Outcome: Understand PySpark's advantages in terms of speed, scalability, and programming simplicity compared to Hadoop.**

# Introduction cont.

**Significance**

- **By comparing the two approaches, we will gain a deeper understanding of:**
    - **The strengths and limitations of Hadoop MapReduce and PySpark in computational tasks.**
    - **Practical insights into setting up and running distributed computations.**
    - **The nuances of performance, scalability, and ease of implementation in both frameworks.**
- **This project not only provides a hands-on experience with Hadoop and PySpark but also lays the groundwork for understanding their applications in more complex data processing and analytical tasks in the real world.**



$$\frac{\text{Circumference}}{\text{Diameter}} = \pi = 3.14159\ldots$$

# Design



**Pi Concept:**

- To estimate the value of Pi, we will use a Monte Carlo method. In this approach, we will throw 'N' darts at a square dartboard. Each dart lands at a random position (x, y) within the square.
- To determine whether a dart has landed inside the circle inscribed within the square, we use the equation of a circle. Specifically, a dart is inside the circle if the condition `x^2 + y^2 < r^2` is satisfied, where 'r' is the radius of the circle.
- By calculating the ratio of darts that land inside the circle to the total number of darts thrown, we can approximate the value of Pi.

# Design

**Calculating Pi: Process Design**

- Generate random coordinates (x, y) within a square.
- Map each (x, y) pair to a result: determine if the point lies inside the circle (assign 1) or outside (assign 0).
- Count the number of points that fall inside the circle.
- Sum up the values to determine the proportion of points inside the circle, which will be used to approximate the value of Pi.

# Implementation - Environment

# Implementation - Prepare Input Data

```
nhaile96456@cs570:~$ ls
hadoop-3.4.0-src   hadoop-3.4.0.tar.gz
nhaile96456@cs570:~$ mkdir PiCalculation
nhaile96456@cs570:~$ cd PiCalculation/
nhaile96456@cs570:~/PiCalculation$ vi GenerateRandomNumbers.java
nhaile96456@cs570:~/PiCalculation$ javac GenerateRandomNumbers.java
nhaile96456@cs570:~/PiCalculation$ java -cp . GenerateRandomNumbers
How many random numbers to generate:
1000000
What's the radius?
200
nhaile96456@cs570:~/PiCalculation$ ls
GenerateRandomNumbers.class   GenerateRandomNumbers.java   PiCalculationInput
nhaile96456@cs570:~/PiCalculation$
```

```
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user/nhaile96456
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user/nhaile96456/picalculate
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user/nhaile96456/picalculate/input
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -put ../PiCalculation/PiCalculationInput /user/nhaile96456/picalculate/input
nhaile96456@cs570:~/hadoop-3.4.0-src$ vi PiCalculation.java
nhaile96456@cs570:~/hadoop-3.4.0-src$
```

# Implementation - Prepare Input Data



```
nhaile96456@cs570:~$ ls
hadoop-3.4.0-src    hadoop-3.4.0.tar.gz
nhaile96456@cs570:~$ mkdir PiCalculation
nhaile96456@cs570:~$ cd PiCalculation/
nhaile96456@cs570:~/PiCalculation$ vi GenerateRandomNumbers.java
nhaile96456@cs570:~/PiCalculation$ javac GenerateRandomNumbers.java
nhaile96456@cs570:~/PiCalculation$ java -cp . GenerateRandomNumbers
How many random numbers to generate:
1000000
What's the radius?
200
nhaile96456@cs570:~/PiCalculation$ ls
GenerateRandomNumbers.class   GenerateRandomNumbers.java   PiCalculationInput
nhaile96456@cs570:~/PiCalculation$
```

```
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user/nhaile96456
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user/nhaile96456/picalculate
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -mkdir /user/nhaile96456/picalculate/input
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hdfs dfs -put ../PiCalculation/PiCalculationInput /user/nhaile96456/picalculate/input
nhaile96456@cs570:~/hadoop-3.4.0-src$ vi PiCalculation.java
nhaile96456@cs570:~/hadoop-3.4.0-src$
```

# Implementation - Prepare Input Data

# Implementation - MapReduce Program

```
nhaile96456@cs570:~/hadoop-3.4.0-src$ vi PiCalculation.java
nhaile96456@cs570:~/hadoop-3.4.0-src$ bin/hadoop com.sun.tools.javac.Main PiCalculation.java
nhaile96456@cs570:~/hadoop-3.4.0-src$ jar cf wc.jar PiCalculation*class
nhaile96456@cs570:~/hadoop-3.4.0-src$ ls
 LICENSE-binary                        PiCalculation.class            WordCount.class     index.html.1    index.html.7    licenses-binary    temp
 LICENSE.txt                           PiCalculation.java             WordCount.java      index.html.2    index.html.8    logs               wc.jar
 NOTICE-binary                         README.txt                     bin                 index.html.3    index.html.9    output
 NOTICE.txt                            WordCount                      etc                 index.html.4    input           output1
'PiCalculation$IntSumReducer.class'   'WordCount$IntSumReducer.class'    include          index.html.5    lib             sbin
'PiCalculation$TokenizerMapper.class' 'WordCount$TokenizerMapper.class'  index.html       index.html.6    libexec         share
nhaile96456@cs570:~/hadoop-3.4.0-src$
```

```java
public void reduce(Text key, Iterable<IntWritable> values,
        Context context) throws IOException, InterruptedException {
    int sum = 0;
    for (IntWritable val : values) {
        sum += val.get();
    }
    result.set(sum);
    context.write(key, result);
}
```

# Implementation - Execution

$ bin/hadoop jar wc.jar PiCalculation /user/nhaile96456/picalculate/input /user/nhaile96456/picalculate/output
$ bin/hdfs dfs -ls /user/nhaile96456/picalculate/output
$ bin/hdfs dfs -cat /user/nhaile96456/picalculate/output/part-r-00000

# Test

**Result:**

S = 784885

N = 1000000

Pi = 4 * S / N = 4 * 784885 / 1000000 = 3.13954

```
            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=9448752
            File Output Format Counters
                    Bytes Written=29
inside  784885
outside 215115
Inside:784885, Outside:215115
PI:3.13954
```

# Test cont.

```
nhaile96456@cs570:~/hadoop-3.4.0-src$ sbin/stop-dfs.sh
Stopping namenodes on [localhost]
Stopping datanodes
Stopping secondary namenodes [cs570]
nhaile96456@cs570:~/hadoop-3.4.0-src$
```

# Implementation - Environment

# Implementation - Environment

# Implementation - Prepare the Input File

In the cloud shell create the calculate_py.py script and copy it to the cloud bucket storage

```python
import argparse
import logging
from operator import import add
from random import random

from pyspark.sql import SparkSession

logger = logging.getLogger(__name__)
logging.basicConfig(level=logging.INFO, format='%(levelname)s: %(message)s')


def calculate_pi(partitions, output_uri):
    """
    Calculates pi by testing a large number of random numbers against a unit circle
    inscribed inside a square. The trials are partitioned so they can be run in
    parallel on cluster instances.

    :param partitions: The number of partitions to use for the calculation.
    :param output_uri: The URI where the output is written, typically a GCS bucket,
                       such as 'gs://bigdata_pyspark/pi-output'.
    """

    def calculate_hit(_):
        x = random() * 2 - 1
        y = random() * 2 - 1
        return 1 if x ** 2 + y ** 2 < 1 else 0

    tries = 100000 * partitions

    logger.info(
        "Calculating pi with a total of %s tries in %s partitions.", tries, partitions)
"calculate_pi.py" 58L, 2044B
```

# Implementation - Prepare the Input File

**Upload the Script to Cloud Storage bucket for easy access, i.e., your newly created bucket.**



```
nhaile96456@cloudshell:~ (cs570-big-data-analytics)$ gsutil cp calculate_pi.py gs://bigdata_pyspark/
Copying file://calculate_pi.py [Content-Type=text/x-python]...
/ [1 files][  2.0 KiB/  2.0 KiB]
Operation completed over 1 objects/2.0 KiB.
nhaile96456@cloudshell:~ (cs570-big-data-analytics)$
```

# Implementation - Running the Script

**Login to Google Cloud through the CLI.**

# Implementation - Running the Script

**Submitting the PySpark job to Dataproc.**

# Implementation - Verifying the Output

**Check the output.**

# Implementation - Verifying the Output

**Check the output.**



```
nhaile96456@cloudshell:~ (cs570-big-data-analytics)$ gsutil cat gs://bigdata_pyspark/pi-output/*.json
{"tries":400000,"hits":314260,"pi":3.1426}
nhaile96456@cloudshell:~ (cs570-big-data-analytics)$
```

# Enhancement Ideas

**Fixed Radius, Variable Random Number Size:**

- Maintain a constant radius size while varying the quantity of random points generated.
- Analyze how changes in the number of random points affect the accuracy of the Pi calculation.

**Fixed Random Number Size, Variable Radius:**

- Keep the number of random points constant while altering the radius size.
- Investigate how different radius sizes impact the precision of the Pi estimation.

**Accuracy Optimization:**

- Conduct comparative experiments to identify the configuration that yields the most accurate approximation of Pi.
- Use these findings to optimize the process for better precision.
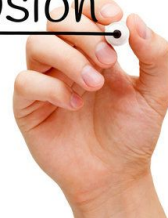
**Significance**

- Gain deeper insights into factors affecting the accuracy of distributed computations, improve other tasks relying on random sampling techniques, and assess the performance of Hadoop MapReduce and PySpark.

# Conclusion

- This project provided valuable hands-on experience with both Hadoop MapReduce and PySpark, highlighting their strengths and limitations in distributed computing.
- By comparing these frameworks, we gained a deeper understanding of their practical applications in large-scale data processing and computational tasks.
- This exploration not only enhanced our knowledge of distributed computing but also laid the groundwork for more advanced and precise computational projects in the future.

*Conclusion*

# References

Python - Calculating π number with Apache Spark

Time for action – using Hadoop to calculate Pi

Overview of Pi calculation using MapReduce

Pi Computation With MapReduce

Calculating Pi with Apache Spark

Run a program to estimate pi | Data Science with Apache Spark

# GitHub Link

- https://github.com/cur10usityDrives/Big-Data/tree/main/PySpark/Pi