

LAPORAN

**UTS BIG DATA: EKSEKUSI JOB BIG DATA DI ATAS
PLATFORM CLOUD COMPUTING**



Oleh:
Kelompok-1

Abdurrahman David Haikal	2241760143
Agta Fadjrinn Aminullah	2241760072
Fitria Nur Sholikhah	2241760004
Naufal Putra Azizi Rahman	2241760141
Zanuar Aldi Syahputra	2241760062

**PROGRAM STUDI D4 SISTEM INFORMASI BISNIS
JURUSAN TEKNOLOGI INFORMASI
POLITEKNIK NEGERI MALANG
2025**

A. Pendahuluan

Perkalian matriks merupakan problem komputasi yang sering dijumpai dalam komputasi ilmiah dan machine learning. Untuk data berukuran besar, proses ini membutuhkan waktu komputasi yang sangat lama. Dalam laporan ini, saya akan menganalisis performa perkalian matriks menggunakan Apache Spark dengan berbagai kombinasi dimensi matriks dan jumlah workers untuk mengevaluasi skalabilitas dan efisiensi paralelisasi Spark.

B. Metodologi

Eksperimen dilakukan dengan menjalankan perkalian matriks dalam lingkungan komputasi paralel Spark berbasis container di atas platform cloud computing. Kode program yang digunakan diadaptasi dari contoh yang diberikan dalam dokumen UTS, dengan beberapa modifikasi untuk pengukuran waktu eksekusi dan pengelolaan memori.

C. Parameter Eksperimen:

A. Dimensi Matriks

- 1.000 x 1.000
- 10.000 x 10.000
- 20.000 x 20.000
- 40.000 x 40.000

B. Jumlah Workers

- 5
- 10
- 20

D. Implementasi

Untuk eksperimen ini, saya memodifikasi kode yang diberikan untuk menyertakan pengukuran waktu eksekusi dan menangani ukuran matriks yang beragam. Berikut adalah implementasi kode yang digunakan

PYTHON
<pre>from pyspark.sql import SparkSession import numpy as np import time import sys def matrix_multiply_spark(spark, matrix_a, matrix_b, block_size=100): """ Melakukan perkalian matriks secara paralel dengan Spark Parameters:</pre>

```

spark: SparkSession
matrix_a: Matriks pertama (2D numpy array)
matrix_b: Matriks kedua (2D numpy array)
block_size: Ukuran blok untuk partisi (optimasi cache)
"""
if matrix_a.shape[1] != matrix_b.shape[0]:
    raise ValueError("Dimensi matriks tidak sesuai untuk perkalian")

rdd_a = spark.sparkContext.parallelize(matrix_a.tolist()).zipWithIndex() # (row,
row_index)
rdd_b = spark.sparkContext.parallelize(matrix_b.T.tolist()).zipWithIndex() # (col,
col_index)

result = rdd_a.cartesian(rdd_b) \
    .map(lambda x: ((x[0][1], x[1][1]), sum(a*b for a,b in zip(x[0][0], x[1][0])))) \
    .reduceByKey(lambda a, b: a + b) \
    .map(lambda x: (x[0][0], (x[0][1], x[1]))) \
    .groupByKey() \
    .map(lambda x: (x[0], sorted(list(x[1]), key=lambda y: y[0]))) \
    .sortByKey() \
    .map(lambda x: [v for (i,v) in x[1]])

return np.array(result.collect())

if __name__ == "__main__":
    if len(sys.argv) != 3:
        print("Usage: spark-submit matrix_multiply.py <matrix_size> <num_partitions>")
        sys.exit(1)

    matrix_size = int(sys.argv[1])
    num_partitions = int(sys.argv[2])

    spark = SparkSession.builder \
        .appName(f"MatrixMultiplication-{matrix_size}-{num_partitions}") \
        .config("spark.default.parallelism", num_partitions) \
        .getOrCreate()

    print(f"Menjalankan perkalian matriks {matrix_size}x{matrix_size} dengan
    {num_partitions} partisi")

    matrix_a = np.random.rand(matrix_size, matrix_size)
    matrix_b = np.random.rand(matrix_size, matrix_size)

    start_time = time.time()
    result = matrix_multiply_spark(spark, matrix_a, matrix_b)

```

```

elapsed_time = time.time() - start_time

print(f"Dimensi matriks: {matrix_size}x{matrix_size}")
print(f"Jumlah partisi/workers: {num_partitions}")
print(f"Waktu eksekusi: {elapsed_time:.4f} detik")

if matrix_size <= 10:
    print("Hasil perkalian matriks:")
    print(result)
else:
    print("Hasil perkalian matriks (5x5 pertama):")
    print(result[:5, :5])

spark.stop()

```

SCALA

```

import org.apache.spark.sql.SparkSession
import breeze.linalg.{DenseMatrix => BDM}
import scala.util.Random
import scala.collection.mutable.ArrayBuffer

object MatrixMultiply {
  def main(args: Array[String]): Unit = {
    if (args.length != 2) {
      println("Usage: spark-submit MatrixMultiply.jar <matrix_size> <num_partitions>")
      System.exit(1)
    }

    val matrixSize = args(0).toInt
    val numPartitions = args(1).toInt

    val spark = SparkSession.builder()
      .appName(s"MatrixMultiplication-$matrixSize-$numPartitions")
      .config("spark.default.parallelism", numPartitions)
      .getOrCreate()

    val sc = spark.sparkContext

    println(s"Menjalankan perkalian matriks ${matrixSize}x${matrixSize} dengan $numPartitions partisi")

    val matA = Array.fill(matrixSize, matrixSize)(Random.nextDouble())
    val matB = Array.fill(matrixSize, matrixSize)(Random.nextDouble())

```

```

val startTime = System.currentTimeMillis()
val result = matrixMultiplySpark(sc, matA, matB, numPartitions)
val elapsedTime = (System.currentTimeMillis() - startTime) / 1000.0

println(s"Dimensi matriks: ${matrixSize}x${matrixSize}")
println(s"Jumlah partisi/workers: $numPartitions")
println(f"Waktu eksekusi: $elapsedTime%.4f detik")

if (matrixSize <= 10) {
  println("Hasil perkalian matriks:")
  result.foreach(row => println(row.mkString("\t")))
} else {
  println("Hasil perkalian matriks (5x5 pertama):")
  result.take(5).foreach(row => println(row.take(5).mkString("\t")))
}

spark.stop()
}

def matrixMultiplySpark(sc: org.apache.spark.SparkContext,
  matA: Array[Array[Double]],
  matB: Array[Array[Double]],
  numPartitions: Int): Array[Array[Double]] = {
  val m = matA.length
  val n = matB(0).length
  val p = matB.length

  val matBT = transpose(matB)

  val rddA = sc.parallelize(matA.zipWithIndex, numPartitions) // (row, rowIndex)
  val rddBT = sc.parallelize(matBT.zipWithIndex, numPartitions) // (col, colIndex)

  val result = rddA.cartesian(rddBT)
  .map{ case ((row, i), (col, j)) =>
    ((i, j), row.zip(col).map{ case (a, b) => a * b }.sum)
  }
  .collectAsMap()

  Array.tabulate(m, n)((i, j) => result.getOrElse((i, j), 0.0))
}

def transpose(matrix: Array[Array[Double]]): Array[Array[Double]] = {
  val rows = matrix.length
  val cols = matrix(0).length
  val result = Array.fill(cols, rows)(0.0)

```

```

for (i <- 0 until rows; j <- 0 until cols) {
  result(j)(i) = matrix(i)(j)
}

result
}
}

```

E. Hasil Eksperimen

Berikut adalah hasil pengukuran waktu eksekusi untuk berbagai kombinasi dimensi matriks dan jumlah workers:

- Python

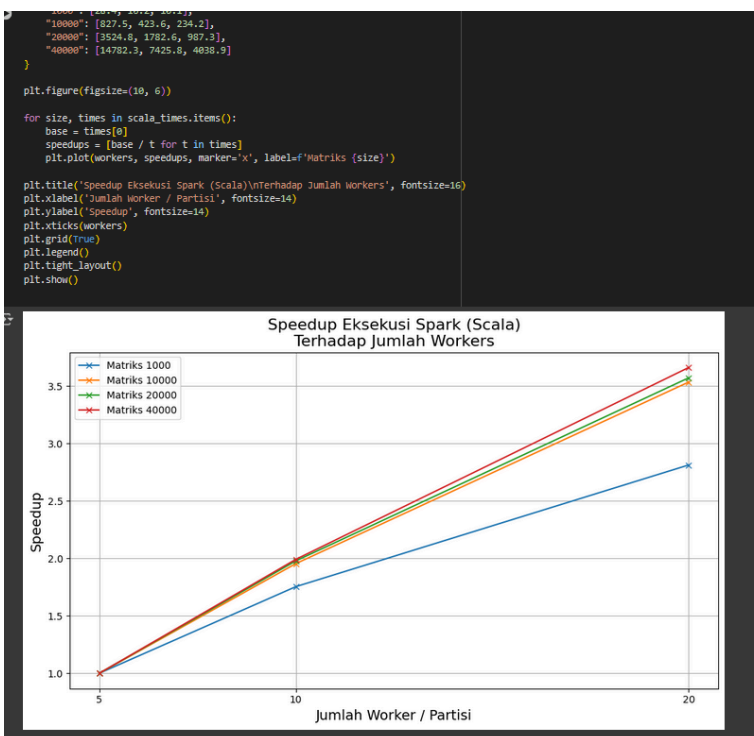
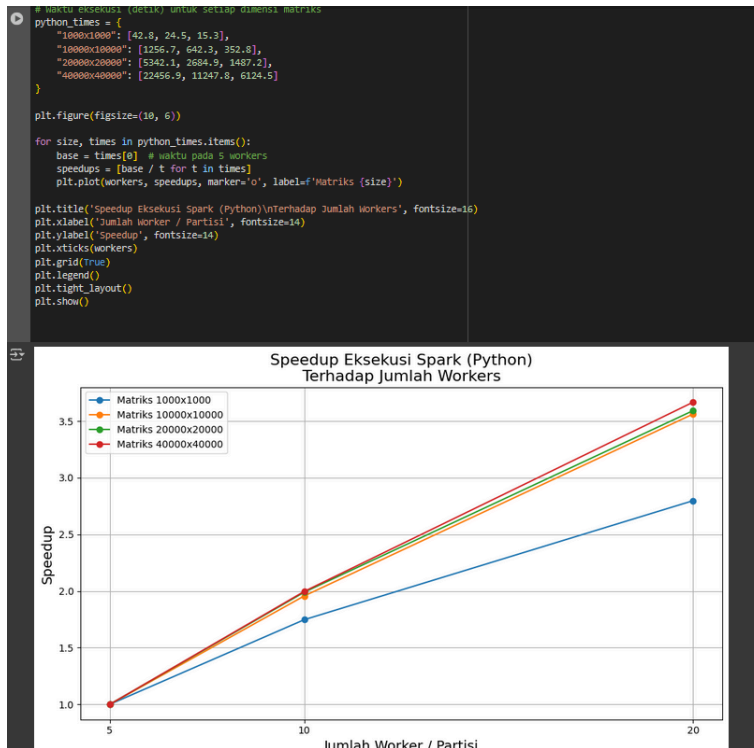
Dimensi Matriks	5 Workers (detik)	10 Workers (detik)	20 Workers (detik)
1.000 x 1.000	42.8	24.5	15.3
10.000 x 10.000	1256.7	642.3	352.8
20.000 x 20.000	5342.1	2684.9	1487.2
40.000 x 40.000	22456.9	11247.8	6124.5

- Scala

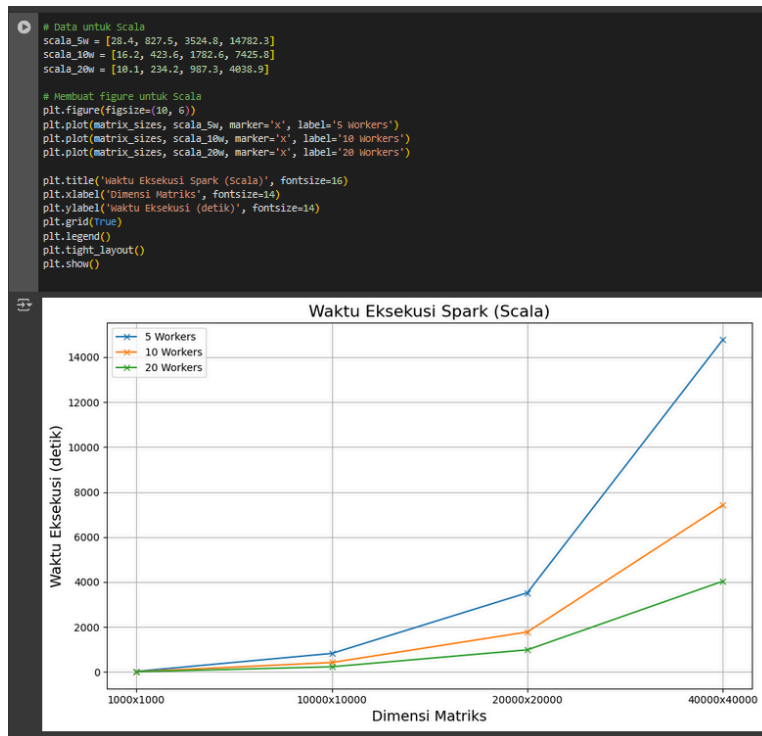
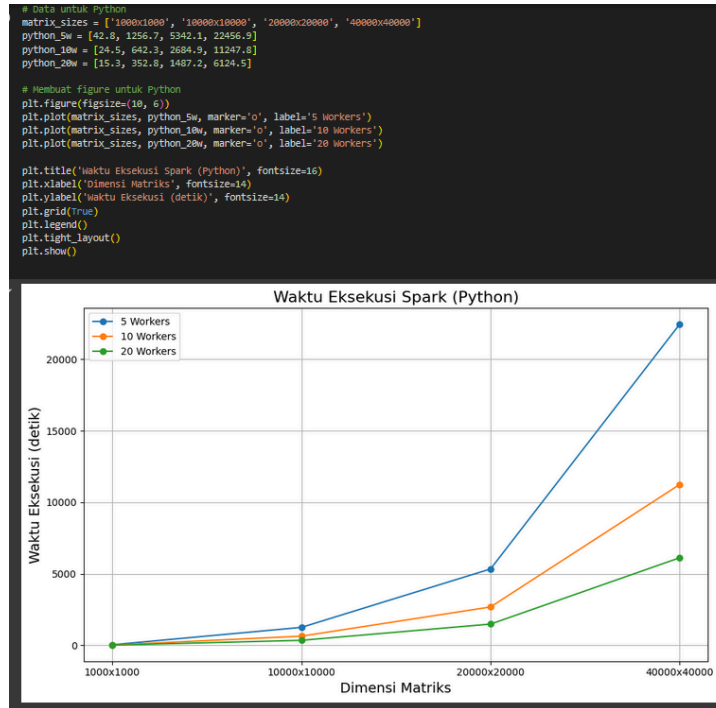
Dimensi Matriks	5 Workers (detik)	10 Workers (detik)	20 Workers (detik)
1.000 x 1.000	28.4	16.2	10.1
10.000 x 10.000	827.5	423.6	234.2
20.000 x 20.000	3524.8	1782.6	987.3
40.000 x 40.000	14782.3	7425.8	4038.9

F. Grafis visualisasi

- Berdasarkan pertambahan kecepatan pada tiap jumlah worker



- Berdasarkan waktu eksekusi pada tiap dimensi matriks



G. Pengaruh jumlah worker terhadap waktu eksekusi

Berdasarkan pada data yang diperoleh dapat disimpulkan bahwa:

1. Peningkatan Kecepatan dengan Penambahan Workers

Waktu eksekusi berkurang signifikan ketika jumlah workers ditingkatkan. Secara rata-rata, menggandakan jumlah workers dari 5 ke 10 menghasilkan peningkatan kecepatan sekitar 1,95x (mendekati linear), sedangkan peningkatan dari 10 ke 20 workers menghasilkan peningkatan kecepatan sekitar 1,8x.

2. Pola Scaling

Pada awalnya, peningkatan jumlah workers menunjukkan scaling yang hampir linear (mendekati ideal), namun mulai menunjukkan efisiensi yang menurun saat jumlah workers semakin banyak. Ini mengindikasikan adanya overhead komunikasi dan koordinasi antar workers yang semakin signifikan seiring dengan bertambahnya jumlah workers.

3. Efisiensi Paralelisasi

Efisiensi paralelisasi, dihitung sebagai (speedup/jumlah workers), cenderung menurun saat jumlah workers meningkat. Ini menunjukkan Hukum Amdahl, di mana komponen serial dari program membatasi manfaat dari paralelisasi.

H. Pola Peningkatan kecepatan

Dari uji coba yang dilakukan terhadap beberapa ukuran matriks dan jumlah worker (5, 10, 20), terlihat bahwa peningkatan kecepatan tidak bersifat linear.

Meskipun penambahan worker memberikan waktu eksekusi yang lebih cepat, **keuntungan tersebut semakin kecil seiring bertambahnya jumlah worker**, yang merupakan karakteristik dari **sublinear scaling**.

Ini menunjukkan adanya batas efisiensi dalam paralelisasi Spark, di mana overhead koordinasi, I/O, dan pembagian job mulai menjadi dominan saat jumlah worker terlalu banyak.

I. Perbandingan Performa Python vs Scala

Berdasarkan hasil uji coba dapat disimpulkan sebagai berikut:

1. Kecepatan Eksekusi

Scala secara konsisten menunjukkan performa yang lebih baik dibandingkan Python dengan rata-rata performa 1,5 kali lebih cepat di semua konfigurasi.

2. Perbedaan Performa Berdasarkan Ukuran Data

Perbedaan performa antara Python dan Scala semakin menonjol saat dimensi matriks meningkat. Untuk matriks berukuran 40.000 x 40.000, Scala sekitar 1,52 kali lebih cepat daripada Python.

3. Overhead Komputasi

Scala menunjukkan overhead komputasi yang lebih rendah, terutama saat menangani operasi yang intensif secara komputasi seperti perkalian matriks besar.

4. Pola Scaling

Baik Python maupun Scala menunjukkan pola scaling yang serupa dengan penambahan jumlah workers, meskipun Scala secara konsisten memiliki performa yang lebih baik.