Big Data

Doker (Jobsheet 7)

Dosen :

M. Hasyim Ratsanjani, S.Kom., M.Kom.



Dibuat Oleh :
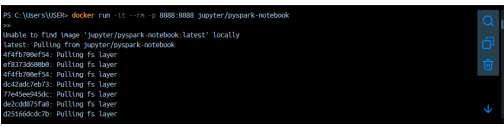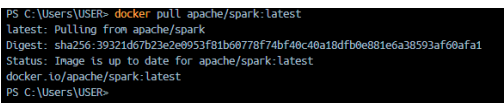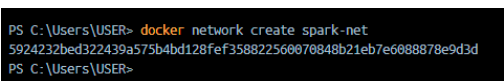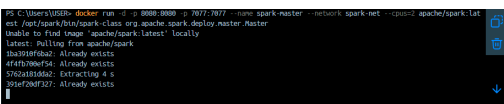Nurhaliza Anindya Putri        2241720016

Kelas : TI – 3D

POLITEKNIK NEGERI MALANG

2025

## 1. Implementasi Python Spark (di Dockers)

| No | Gambar | Deskripi |
|----|--------|----------|
| |  | |
| |  | Pull Image Spark Resmi |
| |  | Buat docker network |
| |  | Menjalankan Spark Master |
| |  | Menjalankan Spark Worker |
| |  | Hasil menggunakan 3 worker |
| |  | Mengakses Spark Web UI |
| |  | Menjalankan Spark Shell |
| |  | Menggunakan Jupyter Notebook dengan Spark |

| | | |
|---|---|---|
| |  | akses Jupyter Notebook di: http://localhost:8888 |
| | | |
| |  | Jalankan Spark Shell di Docker seperti contoh di atas |
| |  | Ketikkan kode berikut di Spark Shell: |
| |  | Hasil |
| |  | Menggunakan PySpark (Python) |
| |  | Ketikkan kode Python berikut: |
| |  | Menggunakan Jupyter Notebook |

| | | |
|---|---|---|
| | ```python
from pyspark.sql import SparkSession

if __name__ == "__main__":
    spark = SparkSession.builder.appName("WordCount").getOrCreate()

    data = ["Hello Spark", "Hello Docker", "Spark is awesome"]
    lines = spark.sparkContext.parallelize(data)

    counts = lines.flatMap(lambda x: x.split(" ")) \
                  .map(lambda x: (x, 1)) \
                  .reduceByKey(lambda a, b: a + b)

    output = counts.collect()
    for (word, count) in output:
        print("%s: %i" % (word, count))

    spark.stop()
``` | Menjalankan Program sebagai Script<br>Buat file wordcount.py dengan isi berikut: |
| | ```
PS C:\Users\USER> cd "D:\Kuliah\Semester 6\BIGDATA"
PS D:\Kuliah\Semester 6\BIGDATA> docker run --rm --name spark-job --network spark-net -v "${PWD}:/app" -w /app apache/spark:
latest /opt/spark/bin/spark-submit --master spark://spark-master:7077 wordcount.py
25/04/22 09:54:12 INFO SparkContext: Running Spark version 3.5.5
25/04/22 09:54:12 INFO SparkContext: OS info Linux, 5.15.167.4-microsoft-standard-WSL2, amd64
25/04/22 09:54:12 INFO SparkContext: Java version 11.0.26
25/04/22 09:54:12 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classe
s where applicable
25/04/22 09:54:13 INFO ResourceUtils: ==============================================================
25/04/22 09:54:13 INFO ResourceUtils: No custom resources configured for spark.driver.
25/04/22 09:54:13 INFO ResourceUtils: ==============================================================
``` | Jalankan script, jangan lupa juga mendifinisikan network spark-net |
| | ```
25/04/22 09:54:41 INFO DAGScheduler: Job 0 is finished. Cancelling potential speculative or zombie tasks for this job
25/04/22 09:54:41 INFO TaskSchedulerImpl: Killing all running tasks in stage 1: Stage finished
25/04/22 09:54:41 INFO DAGScheduler: Job 0 finished: collect at /app/wordcount.py:13, took 23.321105 s
Hello: 2
Spark: 2
is: 1
awesome: 1
Docker: 1
25/04/22 09:54:41 INFO SparkContext: SparkContext is stopping with exitCode 0.
25/04/22 09:54:41 INFO SparkUI: Stopped Spark web UI at http://9d42fc54638d:4040
25/04/22 09:54:41 INFO StandaloneSchedulerBackend: Shutting down all executors
``` | Program-program di atas akan menghasilkan output seperti: |