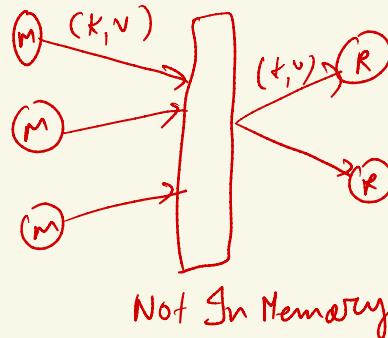

2018



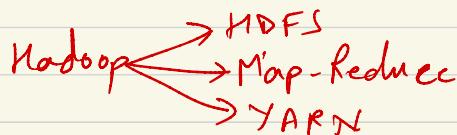
Hadoop vs Spark

Drawbacks of Hadoop

- It was made for Batch processing
- Slow processing because not In memory computation



- Not fault tolerant In terms of computation
- No mechanism for Caching, persistence etc.



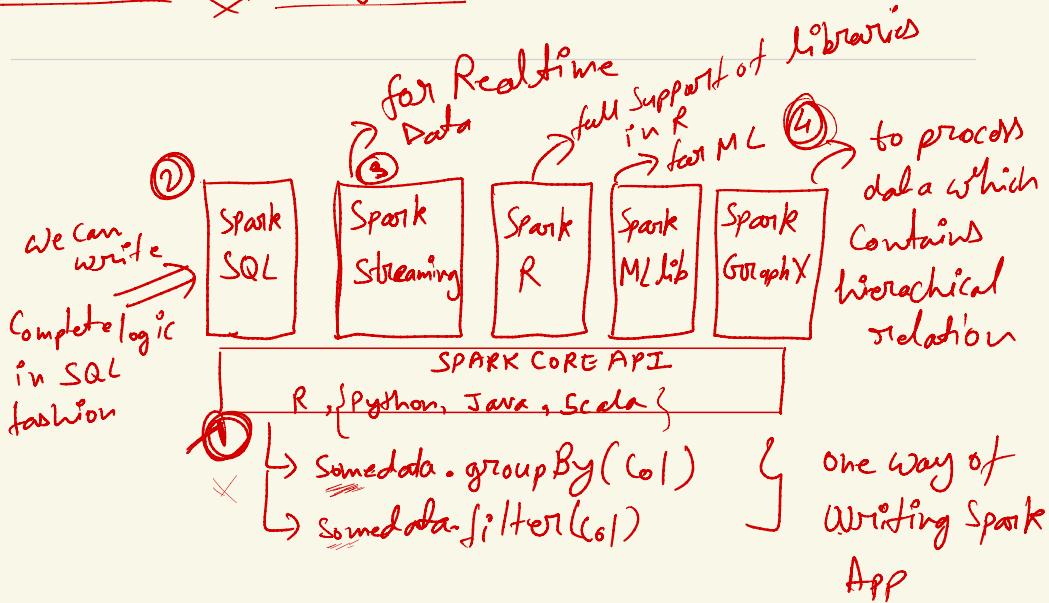
Apache Spark

- It is an distributed computation framework.
- It doesn't have any separate storage or file system.

Features of Spark

- ① Speed: Spark is 100 times faster than Hadoop
- ② Powerful Caching: Capabilities to persist data in memory
- ③ Deployment: More other resource managers can also work with Spark.
 - YARN (mostly used)
 - Mesos
 - Kubernetes
 - Spark own cluster manager
- ④ Good fit for Batch & Real time processing
- ⑤ Polyglot → Spark provides high level APIs or libraries for different languages like
 - Python (PySpark)
 - Java
 - Scala
 - R

Apache Spark Ecosystem

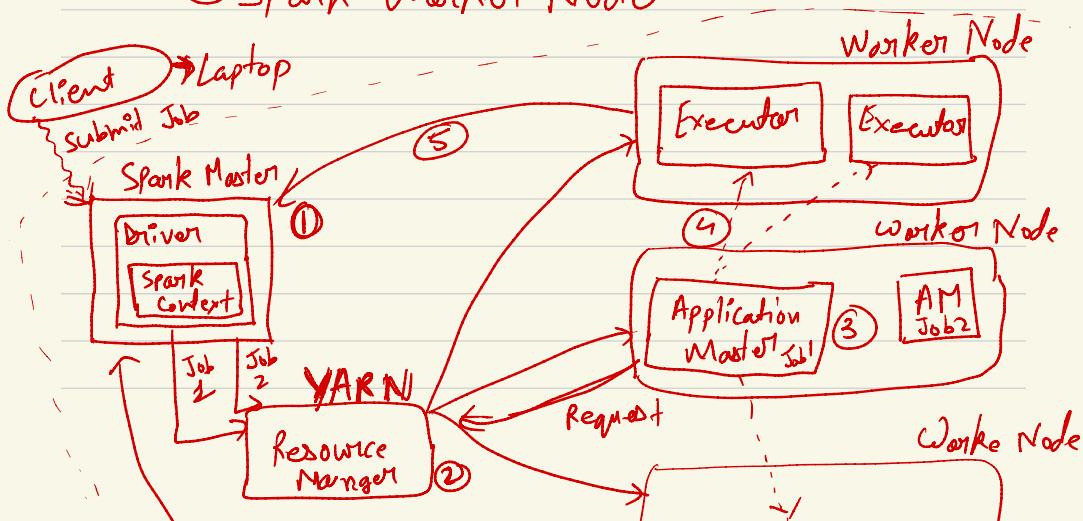


Spark Architecture

→ Spark follows Master-Slave Arch

① Spark Master

② Spark Worker Node



Executor

① Driver program

- It acts like a main method which will create the spark context
- The first method which will be created

② Spark Context or Spark Session

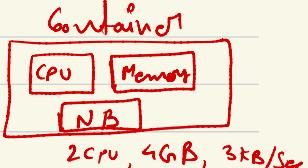
- It is the entry point to spark core of spark application
- It will connect our spark application with execution environment on cluster
- It creates job execution graph, manage partitions, scheduling etc.

Before spark 2.0.0 Spark Context was used separately

After Spark 2.0.0 + Spark Session does all the thing

③ Resource Manager

- It will create one container and start Application Master service in it
- Application Master will request to Resource Manager for complete resource allocation or to start the required executors.
 - Virtual entity which will have CPU, memory, NB

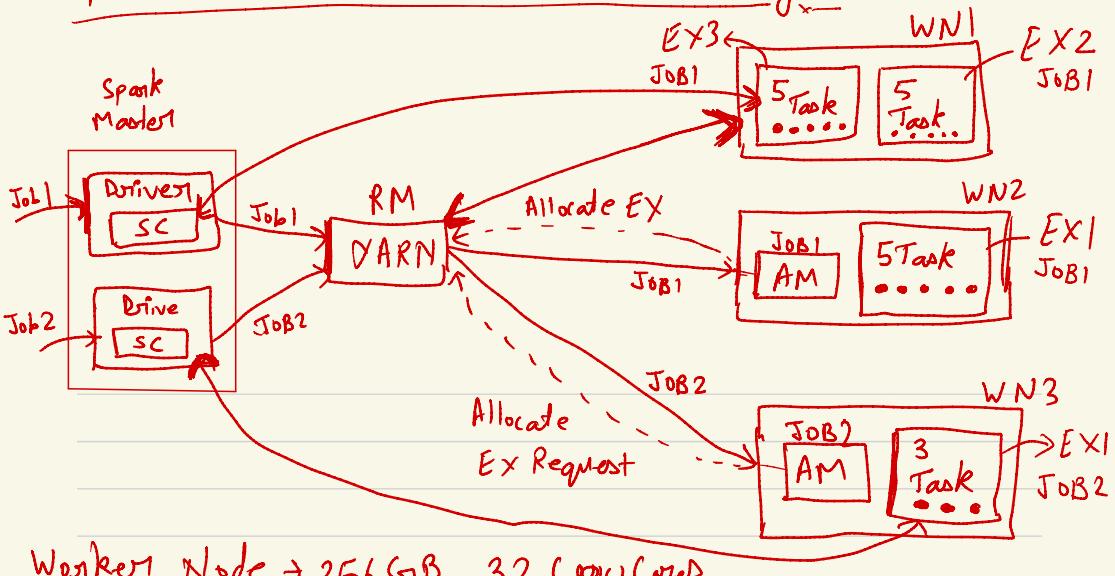


2 vcores, 4GB, 3KB/sec

↳ Executors does the actual computation on partitioned data

- ↳ RM will Create required number of executors for processing.
- ↳ Executors will start interacting with Driver program to send the updates for Job progress and other metadata info.

Spark With YARN as Resource Manager



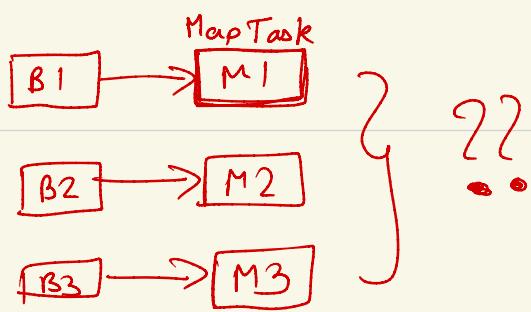
Worker Node \rightarrow 256 GB, 32 CPU Cores

What we need for Job1 \rightarrow 3 Executors, each with 5 CPU Cores & 32 GB Memory

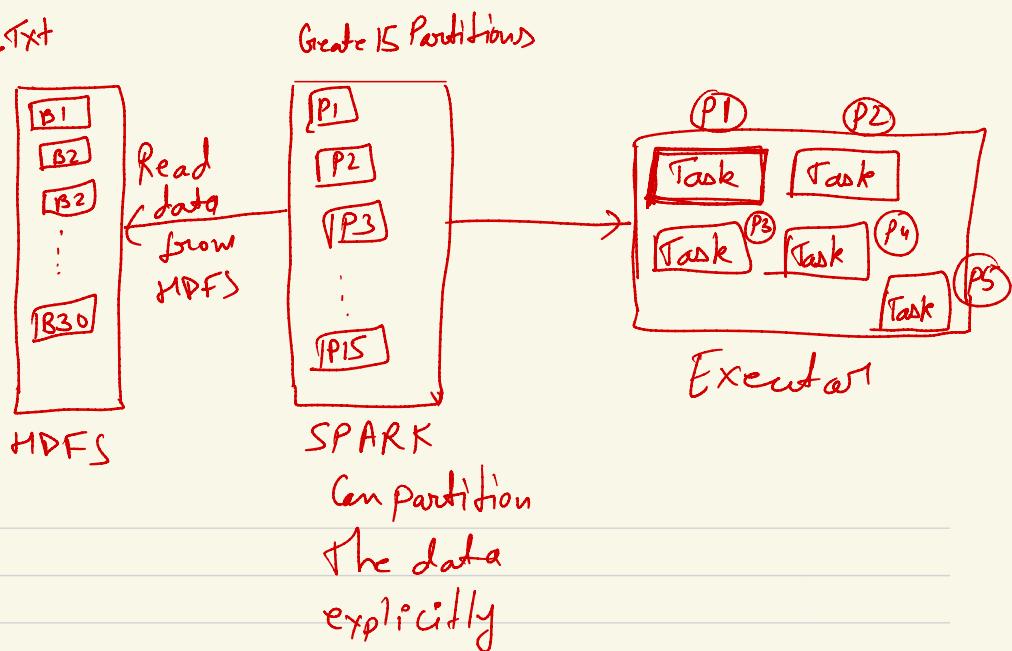
What we need for Job2 \rightarrow 1 Executor, each with 3 CPU Cores & 16 GB Memory

Q.) Parameters in spark-submit Command.

- i) num of executors
- ii) num of CPU cores in each executor
- iii) num of memory for each executor



For Hadoop



Task → ① Smallest unit of execution is known as Task.

② Each task will process one data partition.

③ One task will be executing on One CPU Core

④ One Executor can run multiple tasks and equal to number of CPU Cores

