

Cluster:- group of virtual computer

## Apache Spark Streaming

Processing

Data Storage

DataStream in Spark

classmate

Date  
Page

## Apache Spark MLlib (ML)

- Spark Context
- DAG Scheduler
- Task Scheduler
- Scheduler Backend

→ -----

i) What is spark apache?

ii) Why apache spark? what problem does it solve?

Apache Spark is a unified computing engine and set of libraries for parallel data processing

iii) What is unified?

iv) What is computing engine?

v) What is libraries?

vi) What is parallel data processing?

Unified

→ Spark is designed to support wide range of task over the same computing engine

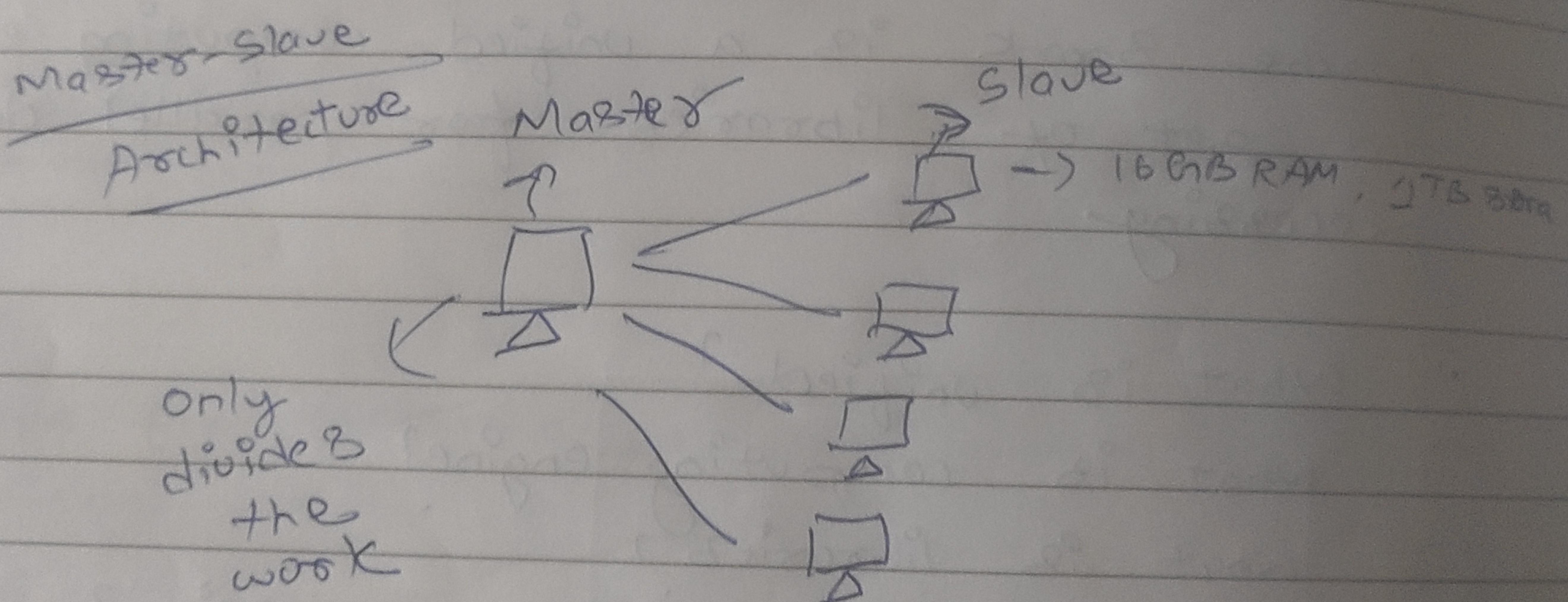
for e.g. Data scientist, Data engineer, Data analyst all can use the same platform for

there transformation or modelling

## Computing Engine ~~Storage~~

- Spark is limited to a computing engine  
It doesn't store the data.
- Spark can connect with different data sources S3, HDFS, JDBC/ODBC, Azure storage etc
- Spark work with almost all data storage system

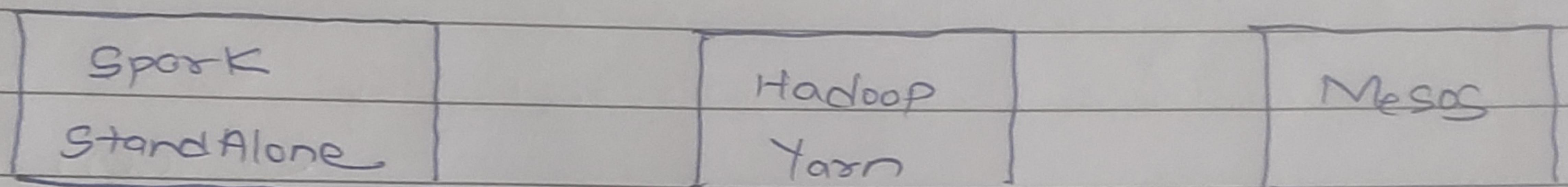
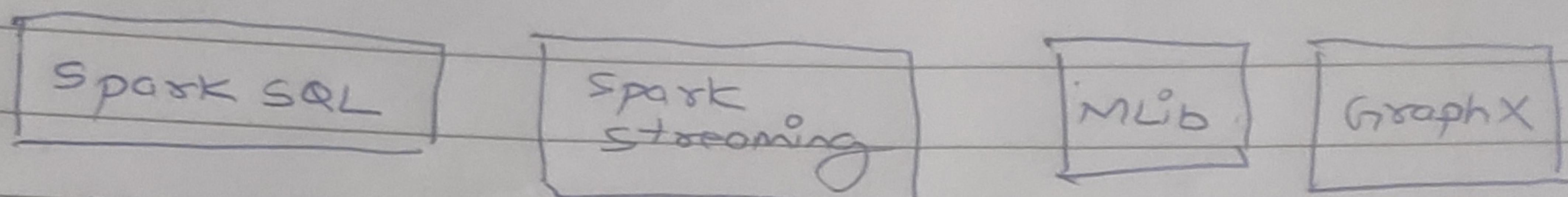
## Computer Cluster



## Spark Features

- Spark written in Scala, runs in JVM
- API : Scala, JAVA, Python, R
- Interactive Shell : Scala and Python
- Data source : SQL, NoSQL, S3, HDFS, Local File System, etc.
- Good fit for iterative task like ML

## Spark Component



## How Spark works

