

## pyspark

### \* components :-

- ① Spark SQL
- ② Spark Streaming
- ③ Spark MLIB
- ④ Spark GraphX
- ⑤ Spark R
- ⑥ Spark core

### \* what is Apache Spark

- clustered computing system
- provide high level API
- Support general execution graph
- High level tools for structured data processing
- can run separate or with existing cluster manager

[Note] → cluster computing system



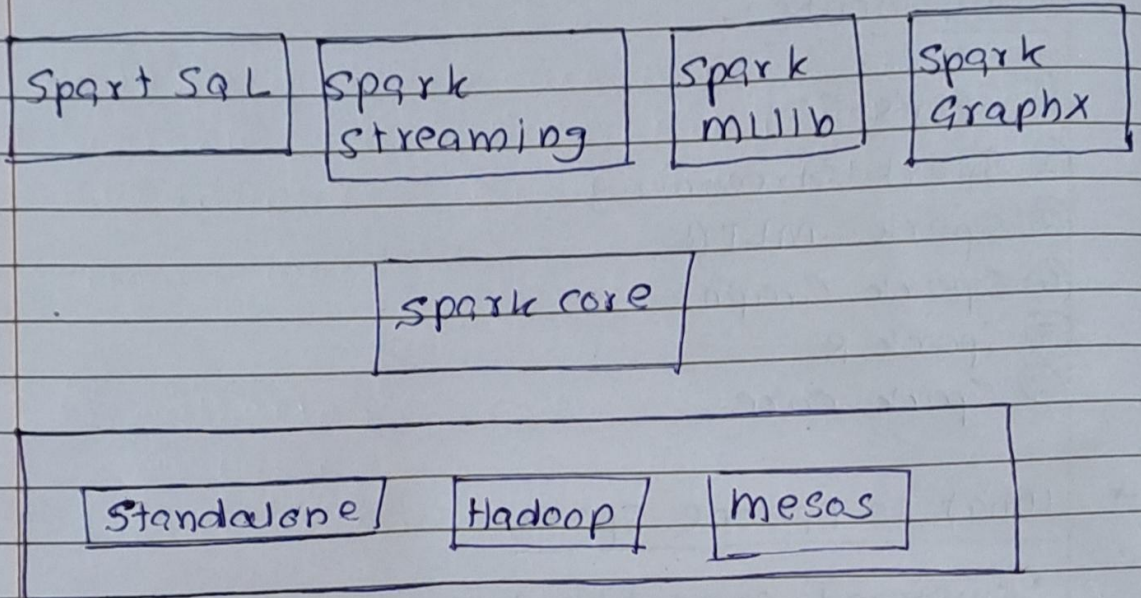
Group of virtual machines

### \* features :-

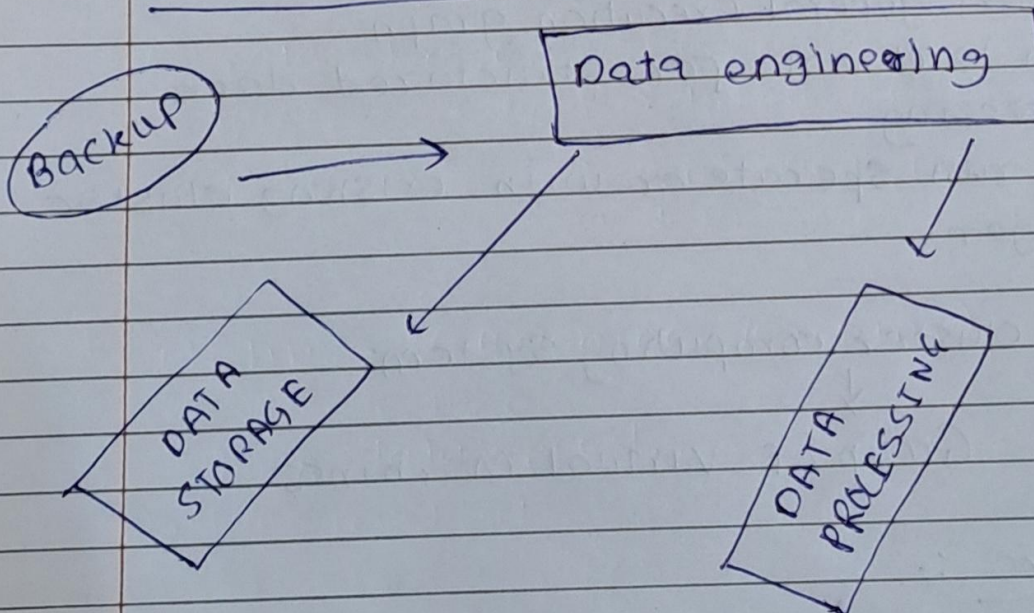
- Developed using Scala language
- Runs on JVM
- API → Scala, Java, Python
- Shell → Scala, Python
- Datasources → SQL, NO-SQL, Local files



\* Component diagram:-

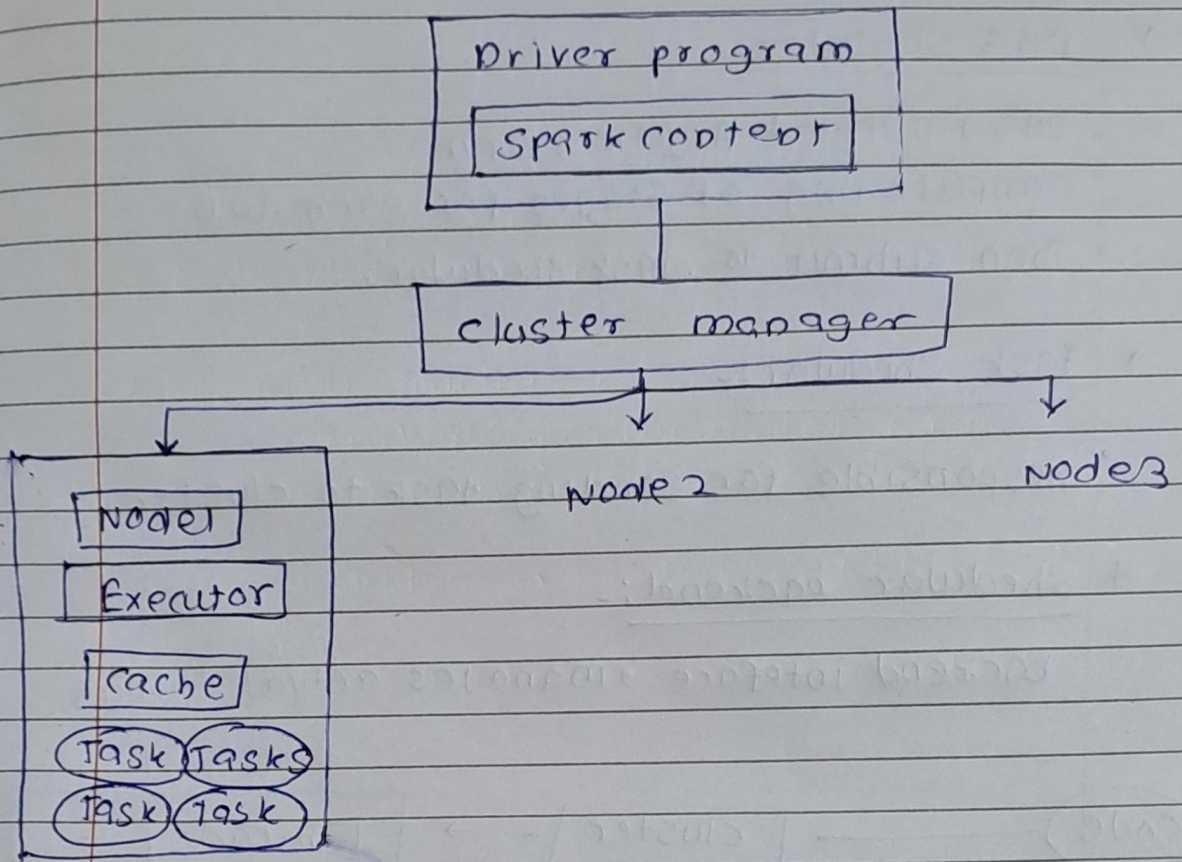


\* data E. concept :-



Parallel and distributed.





- Job
- stages → map, reduce
- DAG → Directed Acyclic graph
- Executor

### \* spark context :-

- connection to spark cluster
- can be use to create RPD.
- accumulator and broadcast variables on that cluster

## \* Streaming :-

- Add on to spark core API which allows Scalable, high throughput fault tolerant stream processing of live data stream.

### How it works ?

#### (a) Gathering :-

① Basic

② Advance

#### (b) processing :-

Gathered data processed.

#### (c) data storage :-

processed data → database

stream → continuous flow of data

## \* MLlib :- (machine learning)

- High data speed.
- High level algo's





## \* Spark core

- All functionalities are build on top of Apache Spark core.
- Embedded with RDD (Resilient distributed dataset)
- RDD is abstraction of Spark.
- two operations performed on RDD.  
(a) Transformation :-

Old RDD  $\longrightarrow$  Create New RDD

- (b) Action :-  
use to work with Actual dataset.

## \* Spark SQL :

- Act as SQL query engine.
- works with structured and semistructured data.
- Also use to access streaming and historical data.
- use for structured data processing.

## \* features :-

- ① cost Based optimizer.
- ② mid query fault tolerance.
- ③ compitible with hive.
- ④ Datafram + SQL  $\longrightarrow$  way to access data.



### \* DAG Scheduler:-

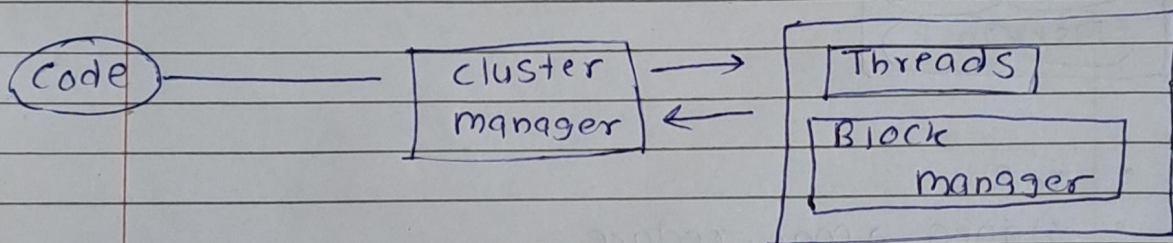
- DAG → Directed Acyclic Graph
- compute DAG of stages for each JOB
- Then submit to task scheduler.

### \* Task scheduler:-

- Responsible for sending task to cluster.

### \* Scheduler backend:-

Backend interface, manages actual things



### \* How it work?

- Small code base.
- divide in various layers
- 1st layer :- interpreter  
↳ Scala interpreter
- As we enter our code in spark console, spark create operator graph.
- Graph submitted to DAG Scheduler
- DAG's Optimised the graphs.
- Then task manager send to task to cluster.



