# *Pache Spark:
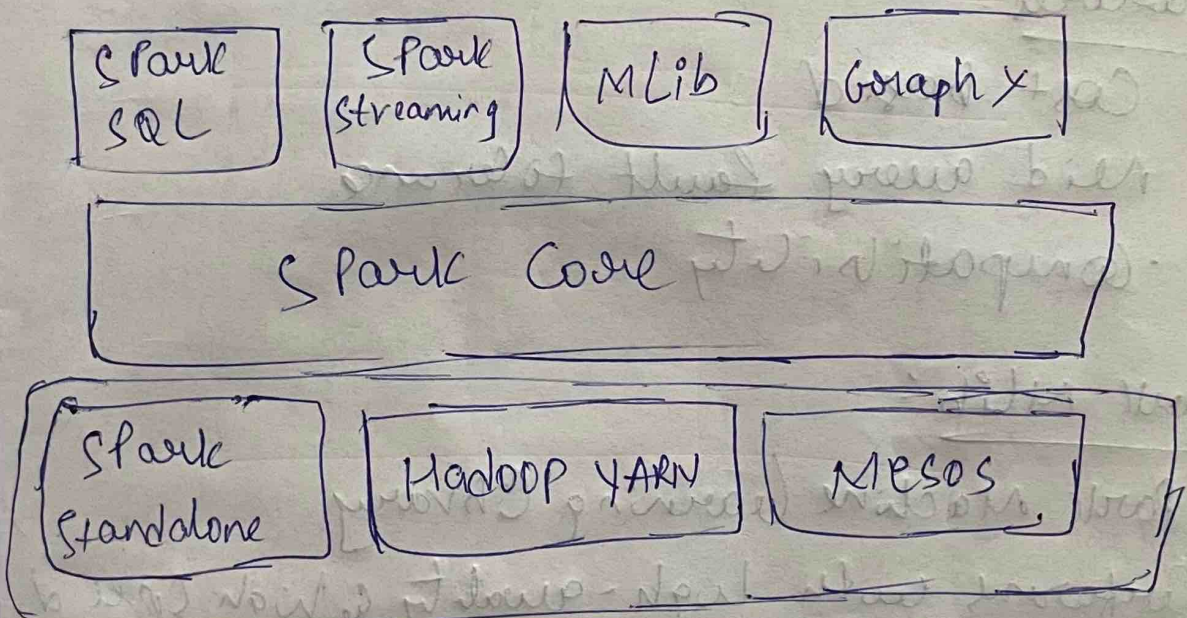
- cluster computing system
- supports execution graph
- In Python / Scala / used.
- It can own alone or by any existing cluster manager.

## Spark Components

SPark Core, SPark SQL, SPark Streaming, SPark MLib, SPark Graph x & Spark R.

| SPark SQL | SPark Streaming | MLib | Graph x |
|---|---|---|---|

| SPark Core |
|---|

| SPark Standalone | Hadoop YARN | Mesos |
|---|---|---|

→ All functionalities are built on "SPark core"

features:-

① essential Input/output & functionalities

② fault recovery

③ Task dispatching

④ Observing "spark cluster"

# Spark Streaming:-

→ It is an add on to core Spark API

→ Scalable, fault tolerant stream processing

# Spark SQL:-

→ distributed framework for "Structured data"

→ This it can perform extra optimization

→ It is spark module for data processing

## Features:

① Cost based

② Mid query fault tolerance

③ Compatibility

# Spark Mlib:-

→ Spark Machine learning library

→ Performs with high-quality & high speed

→ It can perform various Implementation of machine learning platforms.

## Core Concepts

1) **Job:-** It is piece of code which reads input from HDFS.

2) **Stages:-** It is divided into stages

3) **Tasks:** every task has some work.

4) **DAG:** "Directed Acrylic Graph".

5) **Executor:** Executing task

## Spark Components

1) Spark Driver

2) Spark Context

3) DAG Scheduler

4) Task Scheduler