

## Apache Spark =

### Introduction :-

First

Distributed Processing :- It involves the use of multiple computing resources such as computers or servers to solve a single problem. Basically, multiple machines can handle large datasets.

whereas

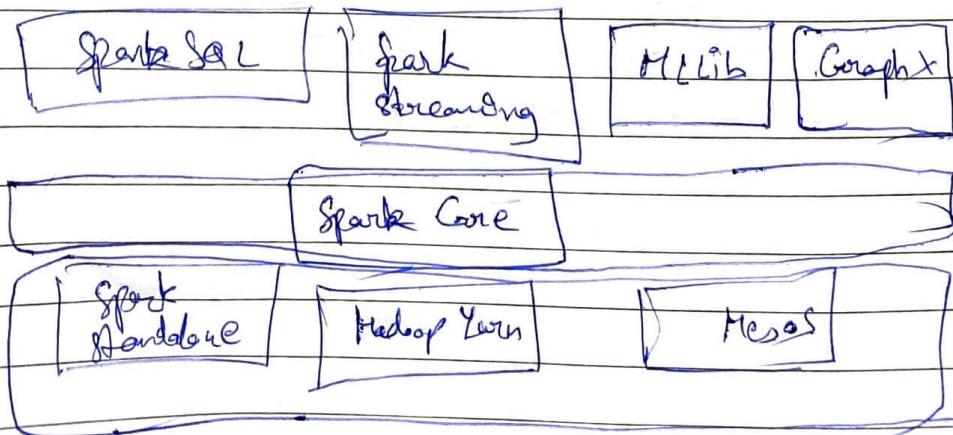
Parallel Processing :- In this one big task is divided into smaller subtasks that can be executed simultaneously.

Spark Components :-

Spark features :-

- ① It is written in Scala programming language and runs in Java
- ② It uses APIs like: Scala, Java, Python, R.
- ③ Interactive Shell or Lang. for Spark :- Scala and Python.
- ④ Data Sources :- SQL, NoSQL, HDFS etc.
- ⑤ It is good for ML algorithms.

Spark Components :-



Q) What is Spark :-

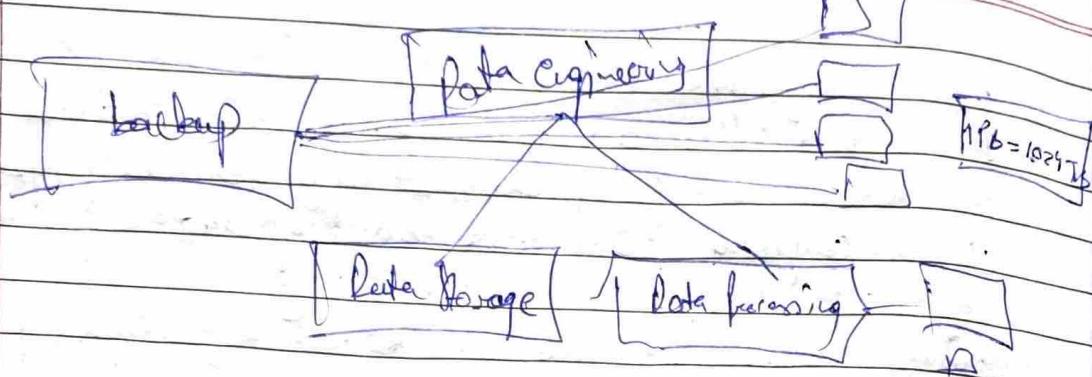
- It is basically an open-source, distributed computing system that provides a fast and general-purpose cluster-computing framework for big data processing.
- Mainly used for machine learning, graph processing and large-scale data processing.
- Spark Core key features :-
  - (i) Inclusion of essential I/O functionalities.
  - (ii) Important in observing the state of the spark cluster.
  - (iii) Task Dispatching
  - (iv) Fault Recovery
  - (v) It overcomes the frag of mapreduce by using in-memory.
- Spark Core is embedded with a special collection called RDD (Resilient distributed datasets).

RDDs basically handles data partitioning across all clusters.

2 operations performed on RDDs :-

Transformation

Action



## ⇒ Apache Spark SQL

- Spark SQL is a distributed framework for Structured data processing.
- It uses same execution engine while computing an output. It does not depend on API/ language to express the computation.
- It work to access Structured and semi-Structured Information.

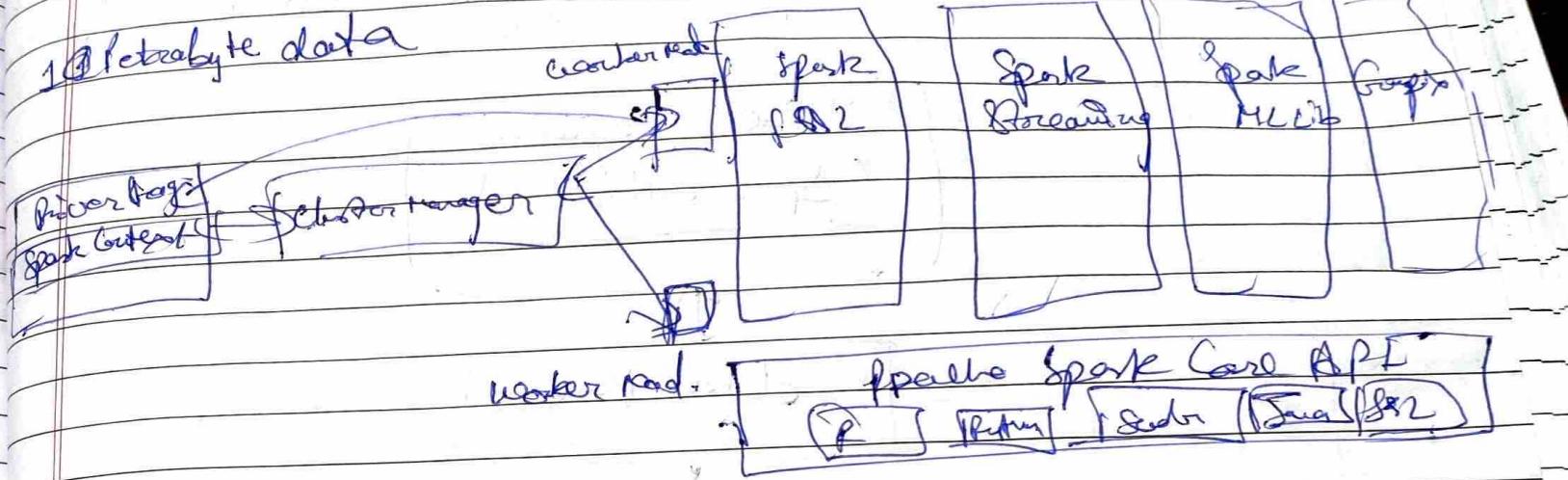
Features:-

- ① cost based optimiser
- ② full compatibility with existing HIVE data.
- ③ Data frames and SQL provide a common way to access a variety of data sources.

HIVE ⇒ Data warehousing and SQL like - query lang. system built on Hadoop.

## Apache Spark Ecosystem.

10 Petabyte data

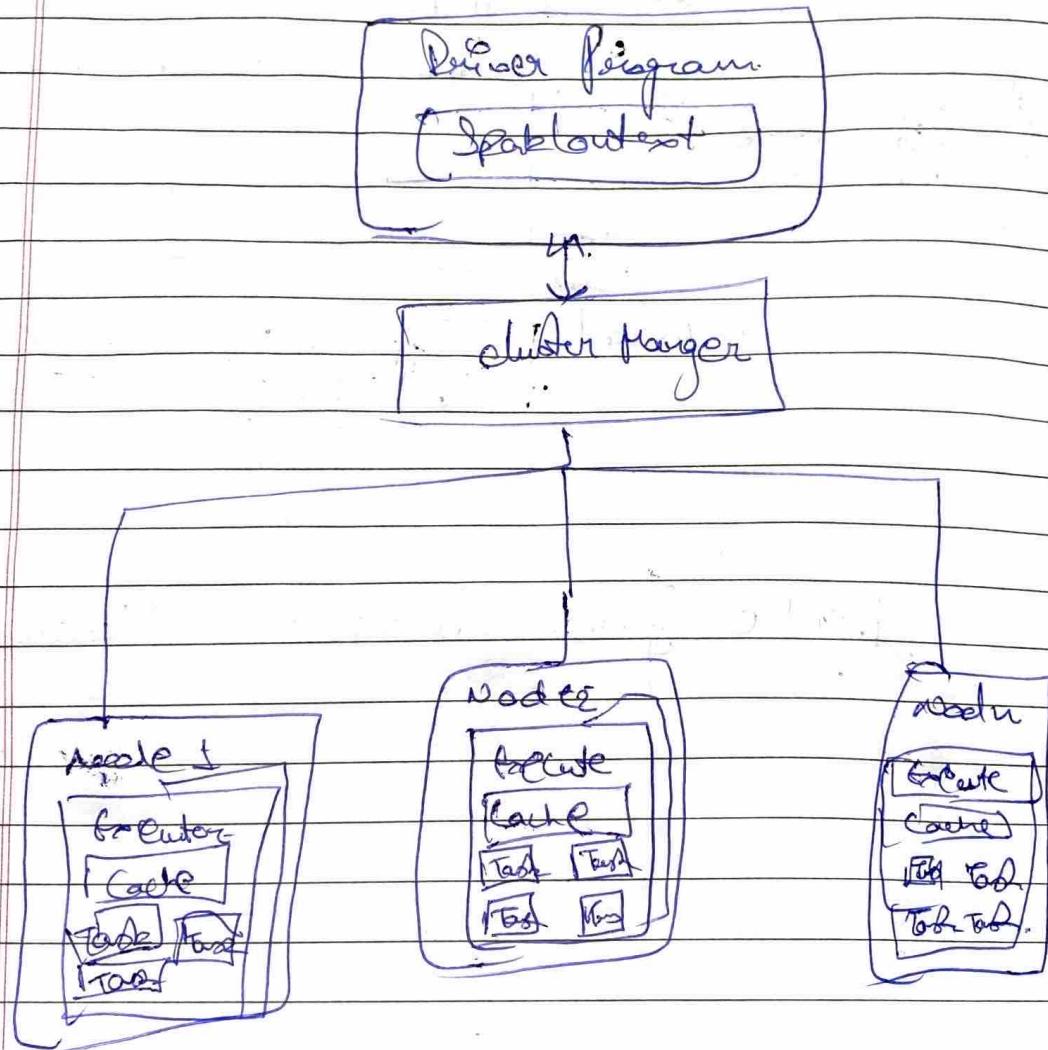


## Apache Spark Cluster Architecture.

## → Apache Spark Streaming -

Cluster → Group of Machines

- Can create in databricks or spark RDDs
- Work together to perform distributed processing



Spark Information will load on cluster, so it becomes Spark cluster.

like

Streaming occurs on platforms like netflix, youtube, etc etc.

2 types of data :- Streaming - live data storage + already stored.

## → Apache Spark Streaming

- It is an add on to spark API which allows scalable, high throughput, fault tolerant stream processing of live data streams. Spark can process data from sources like netflix, youtube etc.

## → Working of spark Streaming

1) Gathering :- Gathering of data from diff sources like kafka, flume, busin, file systems and socket connections etc

2) Processing :- Gathered data processed using complex algorithms

3) Data Storage :- Processed data is pushed out to file systems, databases and in dashboard

4) Data Stream :- It signifies continuous flow of data. It is basically a sequence of files.

## ⇒ Apache Spark MLlib

- It is a Scalable ML library that focused on high quality algos. and high speed.
- Purpose for MLlib creation is to make ML scalable and easy.
- Helps in implementation of clustering, regression, classification and collaborative filtering.

## ⇒ Core Concepts

- ① Job :- piece of code which takes I/O from HDFS or local, performs some computation on the data.
- ② Stages :- Jobs are divided into two stages:-  
Map & Reduce Stages.
- ③ Tasks :- Each Stage has some task one task per partition.
- ④ DAG :- Directed Acyclic Graph.
- ⑤ Executor :- Process responsible for executing my Task
- ⑥ Master :- The machine on which the other progs runs.
- ⑦ Slave :- machine on which Executor runs.

- ① Spark Context :- Represents connection to a spark cluster and can be used to create RDDs, etc.
- ② DAG Scheduler :- Compiles a DAG of stages for each jobs and submits them to Task Scheduler along with preferred location for tasks.
- ③ DAG - ~~Directed Acyclic Graph~~ - Collecting Tasks together.
- ④ Task Scheduler :- Responsible for sending tasks to cluster, running them and retrying in case of failure.
- ⑤ Scheduler Backend - Backend interface for scheduling system that allows plugging in different implementations (HDFS, YARN, Local).

- ① How Spark (Works)?  
Every layer is an interpreter. Spark uses Scala Interpreter.
- ② Code entered in spark console creates a operator graph.
- ③ DAG scheduler divides operator graph into map and reduce stages.  
A task stage is comprised of tasks based on partition of I/O data.
- ④ DAG scheduler pipelines operator together to optimise the graph.
- ⑤ Finally the DAG scheduler passes the stages to the task scheduler which is launched via cluster management (spark standalone / yarn based).