1) Apache Spark (2009) is an open source distributed processing System · used for big data workloads.

History → 2009 · Project vc Berkley

2) Spark features :-

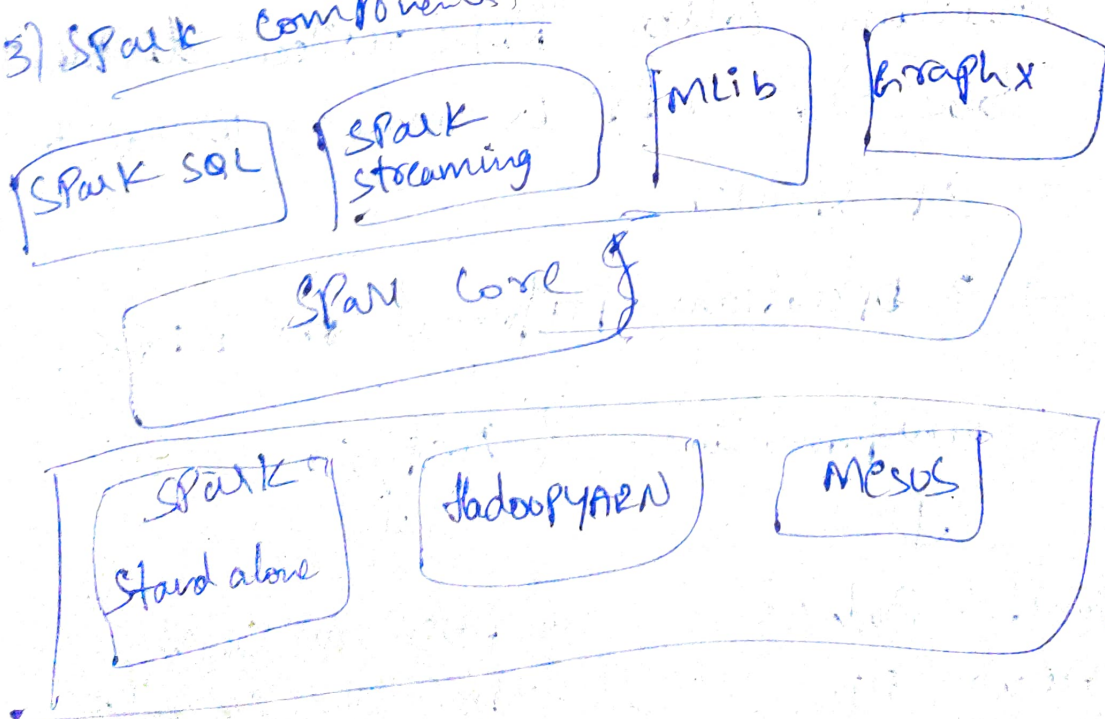→ written in Scala Programing language & runs in JVM.

→ API : Scala, Java, Python, R.

→ Interactive shell : Scala & Python

→ data Sources :- SQL, NOSQL, S3, HDFS, Local File system, Etc.

→ Good fit for iterative tasks like ML algorithm

3) Spark components :-

| Spark SQL | Spark streaming | MLib | Graphx |

| Spark core |

| Spark Standalone | Hadoop YARN | Mesos |

i) Spark Core:
- Functionalities are provided by Apache Spark are built on top of Spark core.
- delivers speed by providing in-memory computation.

## features

- Task dispatching
- fault recovery
- overcomes snag of mapreduce by using in-memory computation

* Spark core is Embedded with a special Collectop called RDD( resilient distributed dataset). RDD is among the abstraction of Spark.

→ 2 operations Performed on RDDS.

Transformation. & Action.

↓

function Produces new RDD from Existing RDD's

RDD's are created from Each other but when we want to work w actual data set, we use action

## ii) Spark SQL :

-) It is ~~distributed~~ distributed framework for structured data Processing.

-) using this Spark get more info about structure of data & computations.

-) with that info, it can perform Entire optimization.

**\* features :**

-) Cost based optimizer

-) Mid query fault - tolerance

-) Data frames.

## iii) Spark Streaming :

=) It allows Scalbi, high-throughput, fault-tolerant stream. Processing of live data streams.

=) 3 phase of spark streaming.

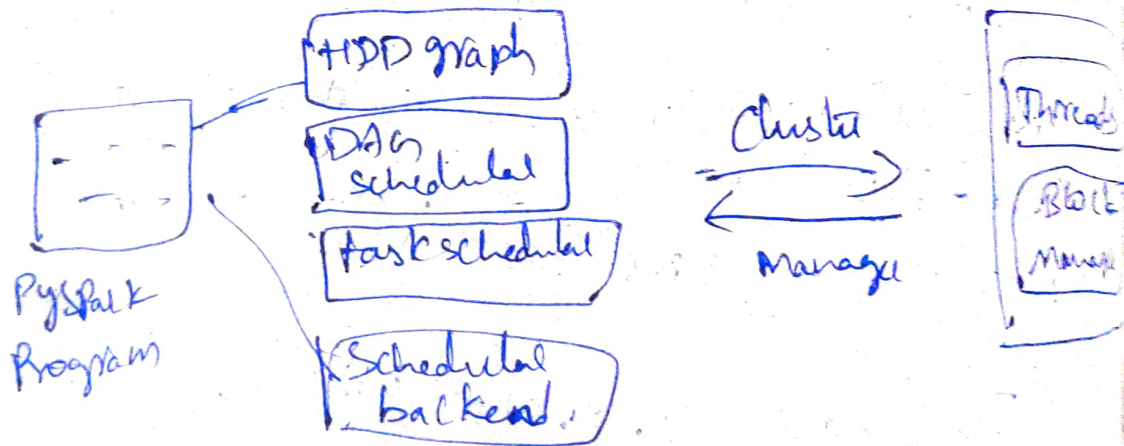| a) Gathering | b) Processing | c) Data storage |
|---|---|---|
| ↓ | · Gathered data | Procesed data is |
| 2 categories | is Processed | Pushed out to |
| basic Source :- Sources available for streaming content Ent Socket connection | using complex | file systems, databases & dashboards. |
| Advanced source :- Kafka, etc. | Algorithms | -) know as discretized stream |

iv) Spark MLlib

→ MLlib is scalable Machine learng library the
    discurses both high-quality algorithms
    high speed.

→ motive is make ~~neable~~ ML scalable &
    Easy.

→ Consists of clustering, regression, dassi,
    & collaborative filtering.

Spark component :

    DAG scheduler
    task scheduler
    schedule Backend.


PySpark Program — HDD graph, DAG Scheduler, task scheduler, Scheduler backend — Cluster Manager — Threads, Block Manager

→ Spark content → Represents connection
    to spark cluster.

→ DAG scheduler → compute a DAG
    of stages for Each Job

-) **Task scheduler** → responsible for sending the tasks to cluster, running them, retrying failures & to run the Jobs.

-) **Schedule Backend** → for scheduling system that allows plugging in diff implementations (Mesos, YARN, Stand alone, Local).

3) **How SPARK works:-** Spark has a Small Code base and system is divided in various layers

8 Layers — Interpreter
spark create a Operator Graph when it runs an action. Graph submitted to DAG scheduler the DAG divide graph into stages & also optimised the graph.

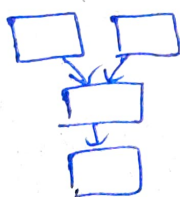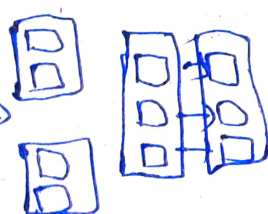-) stages are passed on task scheduler.

RDDObjects      DAG scheduler      task scheduler      worker



DAG →     taskset →     Cluster Manager →     task (threads)
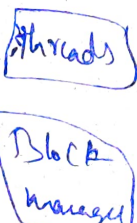
Block manager

build operator DAG      split graph into stages of tasks      launch tasks via cluster manager      Execute tasks store & serve blocks