

### Spark Features :-

- 1) Spark is written in Scala Programming language, and runs in JVM
- 2) API - Scala, Java, Python, R
- 3) Interactive shell - Scala & Python
- 4) Data sources: SQL, NoSQL, S3, HDFS, Local file system
- 5) Good fit for iterative tasks like Machine Learning algos.

### Components of Spark :-



Data from various sources is collected and made to run in spark clusters (group of machines).

\*

## Data Engineering

### Data software storage

cheap storage      costly storage

cost of storage  
components is low

### Apache spark

- \* Spark is used for Data processing
- \* Apache Spark is general purpose

### cluster computing system.

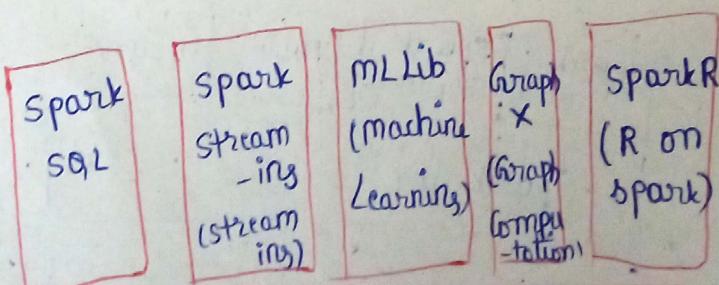
group of systems (virtual)  
working together

- \* It provides high-level API in Java,  
Scala, R, Python

- \* It has high-level tools for structured data processing,  
machine learning, graph processing and streaming

- \* Spark can run alone (or) on existing cluster manager

### Introduction to Apache Spark Ecosystem Components



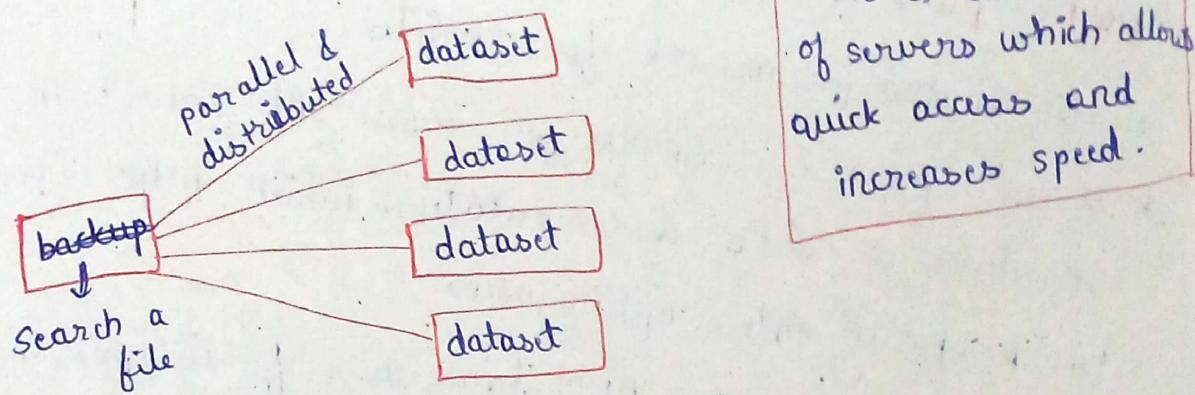
R    SQL    Python    Scala    Java

### Data processing

- depends on
- (i) RAM
  - (ii) storage
  - (iii) processor
  - (iv) I/O speed
  - (v) block speed
  - (vi) OS
  - (vii) hardware component
  - (viii) SSD
  - (ix) CPU

## 1) Apache Spark Core

- \* All functionalities provided by Apache Spark are built on top of Spark Core
- \* It delivers speed by providing in-memory computation
- \* capability
- \* Spark core is the foundation of parallel and distributed processing of huge dataset.



### Features of Spark core :-

- 1) control essential I/O functionalities
- 2) significant / important in programming
- 3) task dispatching (dividing tasks & assigning)
- 4) Fault recovery
- 5) overcomes MapReduce by using in-memory computation
- 6) observes role of spark cluster.

## RDD

- \* Spark core is embedded with special collection called RDD (Resilient distributed dataset)
- \* RDD is among the abstractions of Spark
- \* RDD handles partitioning data across all the nodes in a cluster.

- \* Memory is held in the memory pool of cluster as single unit

## Two operations

1) Transformation produces new RDDs from the existing RDDs.

2) Action Action is used when we want to work with actual data set.

## Apache Spark SQL :-

- \* distributed framework for structured data processing.
- \* Using sparkSQL, spark gets more info about the structure of data & computation which helps spark to perform extra optimization.
- \* doesn't depend on API language for computation but uses same execution engine while computing OLP.
- \* sparkSQL works to access structured & semi-structured info.
- \* It enables powerful, interactive, analytical appln across both streaming & historical data.
- \* It is a spark module for structured data processing.
- \* acts as distributed SQL query engine

## Features

- 1) cost based optimizer
- 2) mid query fault tolerance

- 3) Full compatibility with existing Hive data
- 4) Dataframes and SQL provide a way to access a variety of data sources like Hive, Avro, ORC, JSON, parquet, JDBC
- 5) allows to carry structured data inside spark programs using SQL or other

### Apache Spark Streaming

- \* It allows scalable, high-throughput, fault-tolerant stream processing of live data streams.
- \* Spark can access data from sources like kafka, Flume, Kinesis or TCP socket.
- \* It can operate using various algos.
- \* Data received is given to file streaming, databases
- \* Spark uses micro batches for real-time streaming
  - allows a process or task to treat a stream as sequence of small batches of data.
- \* Hence, spark streaming groups live data into small batches.
- \* Now sent for processing

### Phases of Spark Streaming :-

#### (a) Gathering

built-in stream sources

Basic sources → available in  
streams Context API

Ex:- file systems

Advanced sources

Ex:- kafka, Flume, kinesis

- \* Spark access data from different sources like kafka, flume, kinesis or TCP sockets.

### b) Processing

gathered data is processed using complex algos

Ex:- map, reduce, join and window

### c) Data Storage

processed data is pushed out to file systems, databases

DStream :- High level abstraction provided by spark  
streaming is Dstream (or) Discretized stream.

→ signifies continuous stream of data

#### Formation of Dstream

sources like kafka, Flume & kinesis

by high-level operations on other Dstreams.

\* Dstream  $\Rightarrow$  Sequence of RDDs.

### ④ Apache Spark MLlib (Machine Learning Library)

MLlib  $\rightarrow$  scalable ML library  $\begin{cases} \text{high-quality algo} \\ \text{high speed} \end{cases}$

\* Main aim  $\rightarrow$  make ML scalable & easy

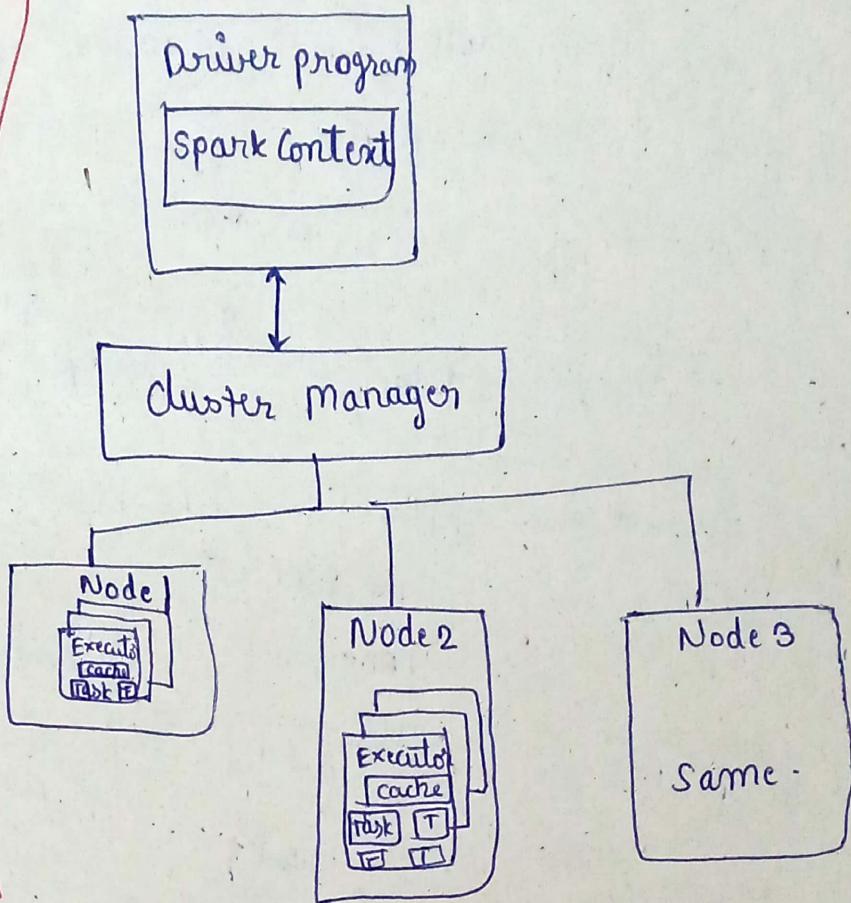
\* Contains many ML libraries  $\rightarrow$  has various algos like clustering, regression and also low level algos are present.

\* In 2.0 version, Dataframe based API is primary ML API for spark (bcz more user friendly)

\* some more benefits  $\rightarrow$  Data sources, SQL DF queries, uniform APIs.

## Architecture      Components (Architecture)

→ ML lib  
was linear  
algebra  
package  
**Breeze**  
↓  
collection  
of libraries  
for numerical  
computations  
& ML



Master :- The machine on which driver program runs

slave :- The machine on which executor program runs

Executor :- the process responsible for executing a

task

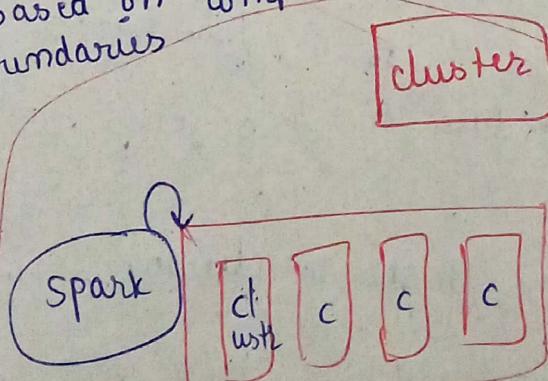
DAG :- Directed Acyclic Graph

job :- piece of code which reads some I/O from local,  
performs some computation on data & writes O/P.

Task stages :- jobs are divided into stages.

based on computational  
boundaries

group of virtual machines

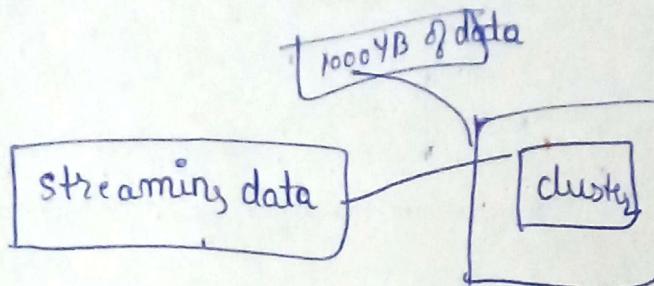


Jobs - tasks :- Each stage  
has some tasks.

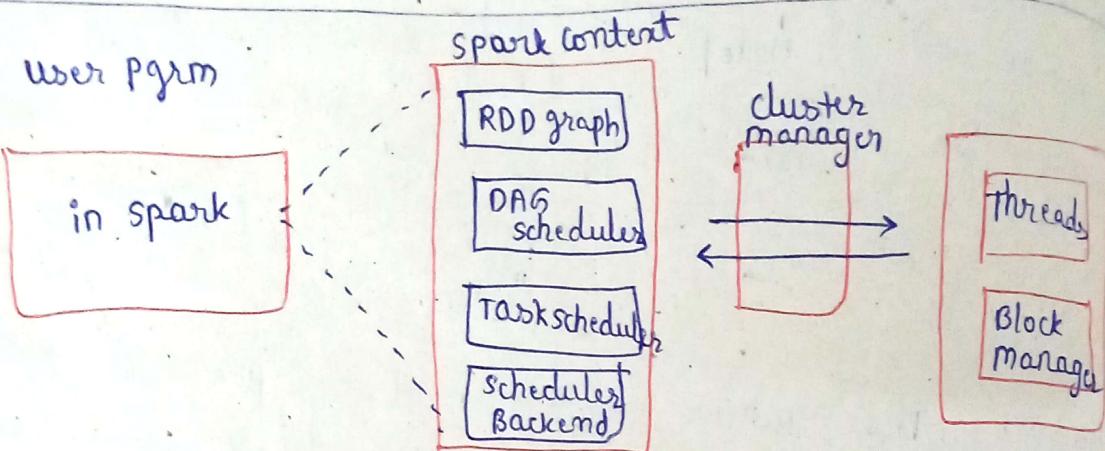
\* one task is executed on  
one partition of data  
on one executor.

- \* Spark SQL is preferred over SQL bcoz spark SQL deals with unstructured data, faster processing
- \* PySpark  $\Rightarrow$  spark + python

\*



- \* Used in various applications like banking



### Spark Context:-

- \* represents connection to spark cluster
- \* used to create RDDs, accumulators
- \* broadcast variables on that cluster.

DAG Scheduler :- , collects tasks together, organized

- \* computes a DAG of stages for each job and submits them to Task Scheduler
- \* TS determines preferred locations for tasks and finds minimum schedule to run jobs

### Task Scheduler :-

- \* responsible for sending tasks to cluster, running them, retrying them if there are failures

stand alone cluster  $\Rightarrow$  spark is installed in our pc  
yarn  $\rightarrow$  yet another resource navigator

M模  $\rightarrow$

Block Manager Scheduler Backend  $\rightarrow$  backend interface for scheduling systems that allows pluggins in different implementations.

How Spark Works?

- \* Spark has small code base and system is divided into various layers.
- \* Each layer has some responsibilities and these layers are independent of each other.

Layer-1 :- [Interpreters]

spark uses Scala interpreter by some modifications

$\downarrow$   
Enter code in spark console (like creating RDD's, apply operators)

$\downarrow$   
Spark creates operator graph

$\downarrow$   
when user runs action, graph submitted to DAG scheduler

$\downarrow$   
DAG scheduler divides op graph into (map & reduce) stages.

Stage :- comprised of tasks based on partitions of input data

$\downarrow$   
DAG scheduler pipelines operators together to optimize graph.

Final result of DAG scheduler is set of stages.

↓  
stages are passed to Task scheduler

↓  
TS launches tasks via cluster Manager (spark  
standalone / Yarn / Mesos)

↓  
Worker executes the tasks

