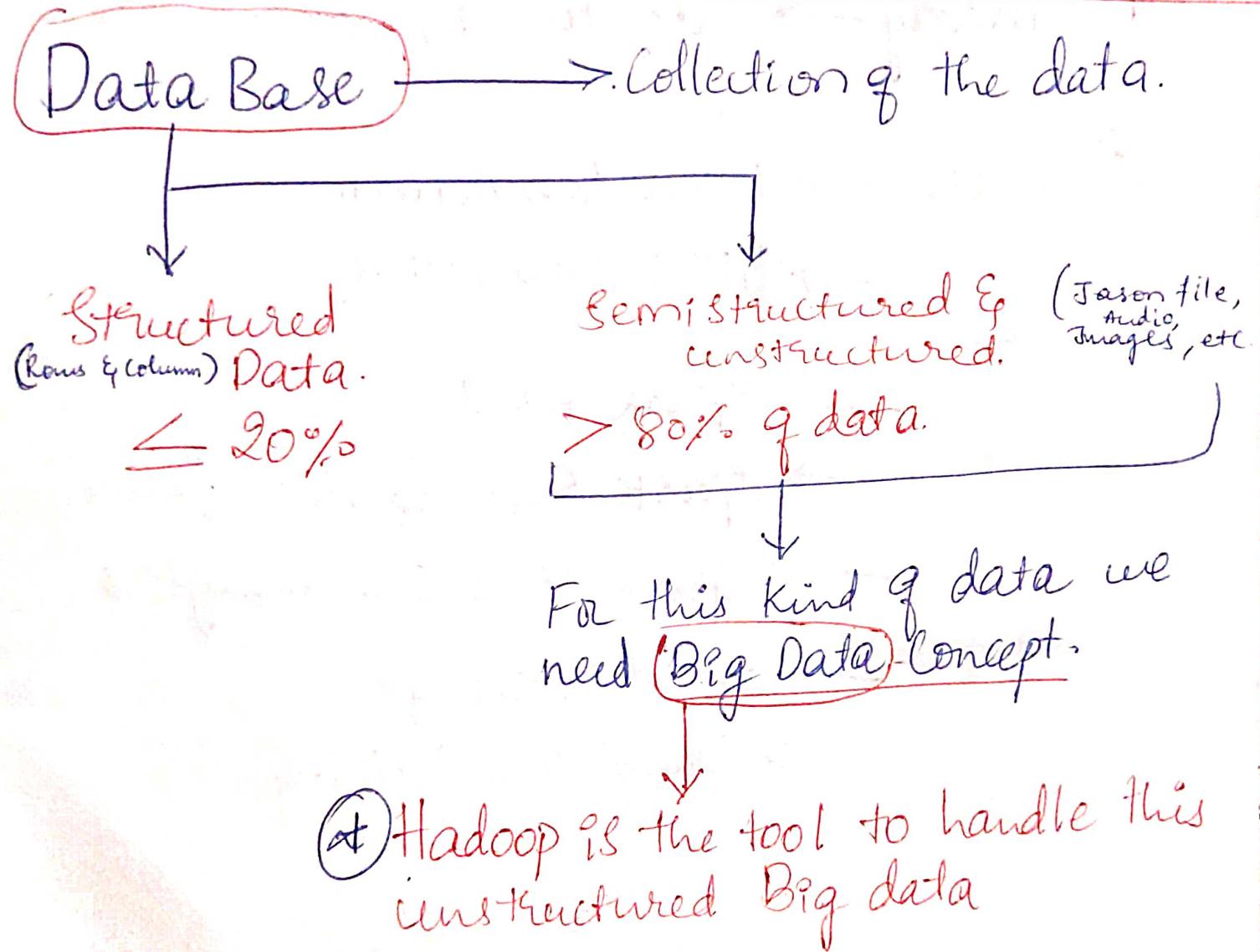
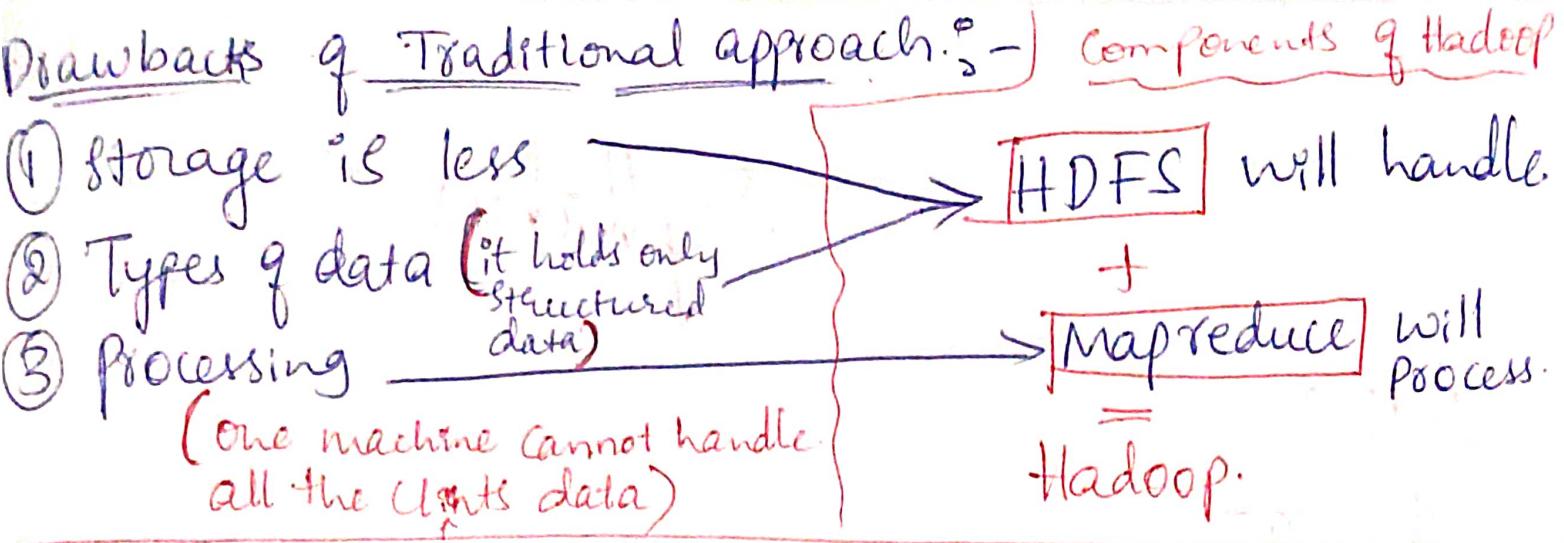


Syllabus

Content

- ① Theory → Into to BD, Hadoop, YARN, &
- ② Lab → Linux Lab Activity
→ HDFS Lab Activity
→ Spark RDD
→ Spark SQL & Data frame
→ Spark Streaming analytics
→ Hive Activity
→ Spark ML Activity.



Linux Activity :-

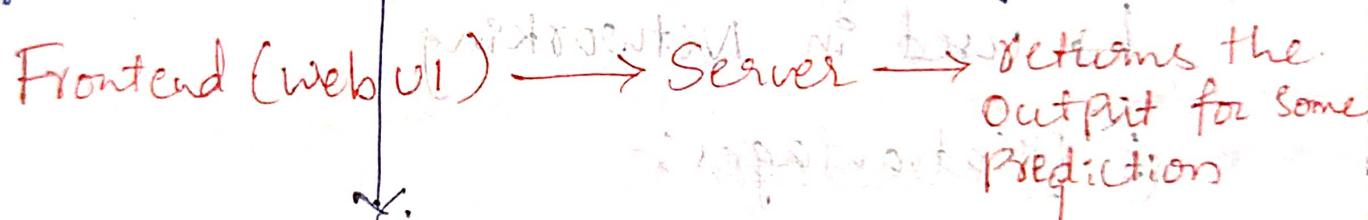
Typically in a project we work in a different environment.

Project life cycle

① Client (Gives requirement of project → ML/DE)

② Under standing the project and work with team

For ex:- To built the machine learning Application.



③ Deploy Application (Software Dev phase) :-

↳ Development Phase → In this phase diff engineers work.

↳ Testing Phase

↳ Production phase

Here in this phase, Machines used are the Linux OS.

Linux OS :-

1991 - Linus Torvalds
built it

*Advantages :-

- ↳ open source (Base Kernel was in C & C++)
- ↳ Security is high & stable
- ↳ Very light & high stability
- ↳ Don't need quick restart.
- ↳ High Performance.
- ↳ Preferred in Embedded Systems
- ↳ Many Linux distributions (Ubuntu, Fedora, etc.)
- ↳ It provides Shell Scripting
- ↳ Huge community support
- ↳ used in Networking

*Disadvantages :-

- ↳ CLI, Commands
- ↳ No. Standard edition of Linux.
- ↳ Learning Curve is very high.

We use Command Line Interface (CLI) to work with Linux.

→ Different users in Linux :-

- ① Root user :- Created by default in OS.
→ used to make changes in administration in OS.
- ② Regular user :-
 - ④ used for many functions
 - ④ It has username & password
- ③ Service Account :-
 - ④ used for some installation

→ We use Regular user on the Jupyter hub (Terminal)

① `pwd` → Present working directory.

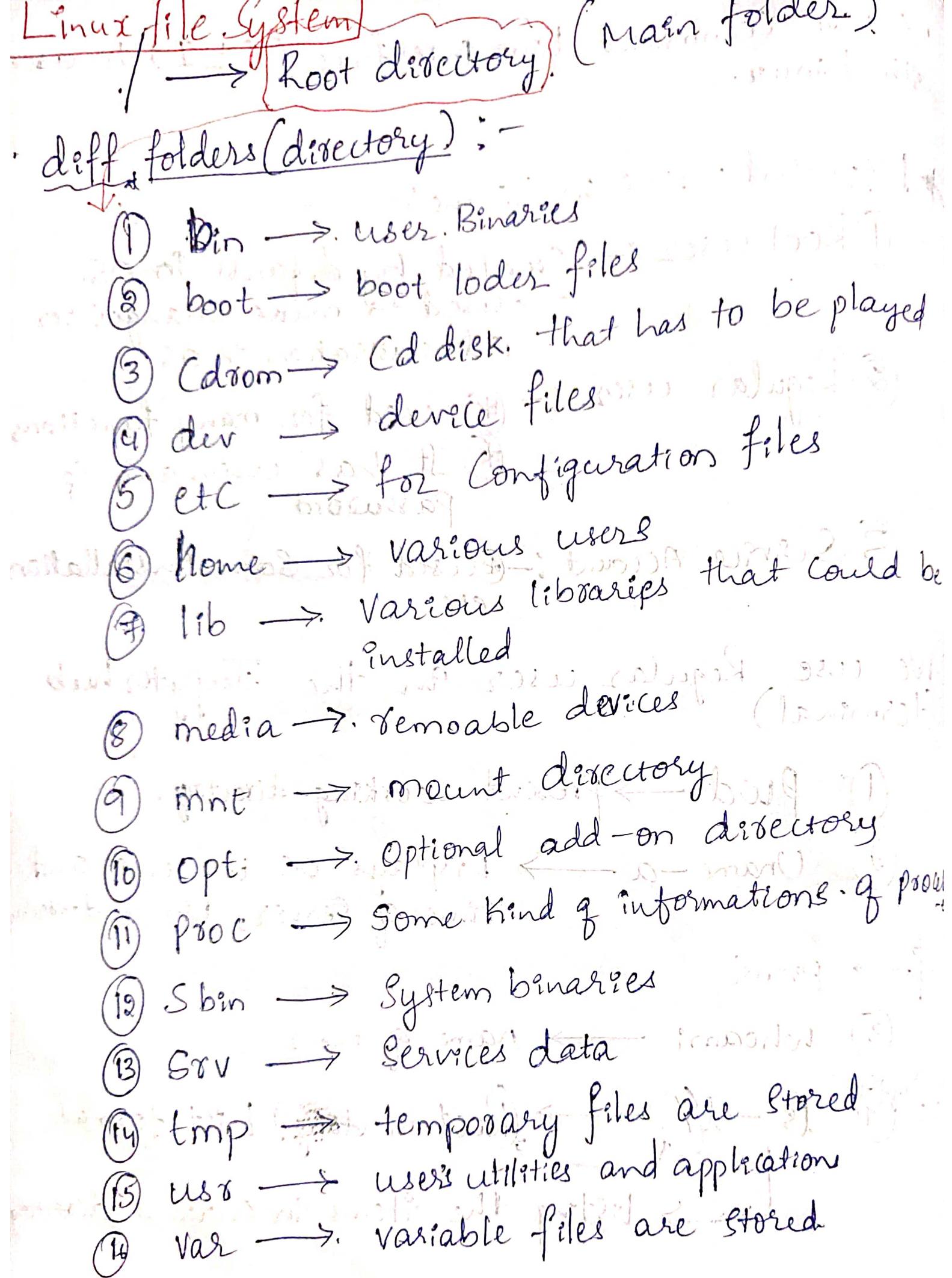
② `uname -a` → Displays OS we use and boot time of server, Linux distribution

\$ → prompt

③ `whoami` → name of user

ls -l / → (Root directory) List format

listing the items in current directory



④ How to create a directory (folder) in Linux.

↳ mkdir foldername

Ex :- xyz@fractal:~\$ mkdir 22jan.

⑤ How to enter the directory.

↳ ~~cd~~ Cd. 22jan. OR Cd/home/name/22jan

⑥ How to enter ~~previous~~ home. folder. (root directory)

↳ Cd.

⑦ Cd - → Take us to previous folder

⑧ Cd .. → Take to parent directory of folder

Permissions :-

d → directory

- → file

r → read. ($4 \rightarrow 2^2$)

w → write ($2 \rightarrow 2^1$)

x → execute. ($1 = 2^0$)

7 7 7
d rwx rwx rwx.
↓ ↓ ↓
user's group other
only

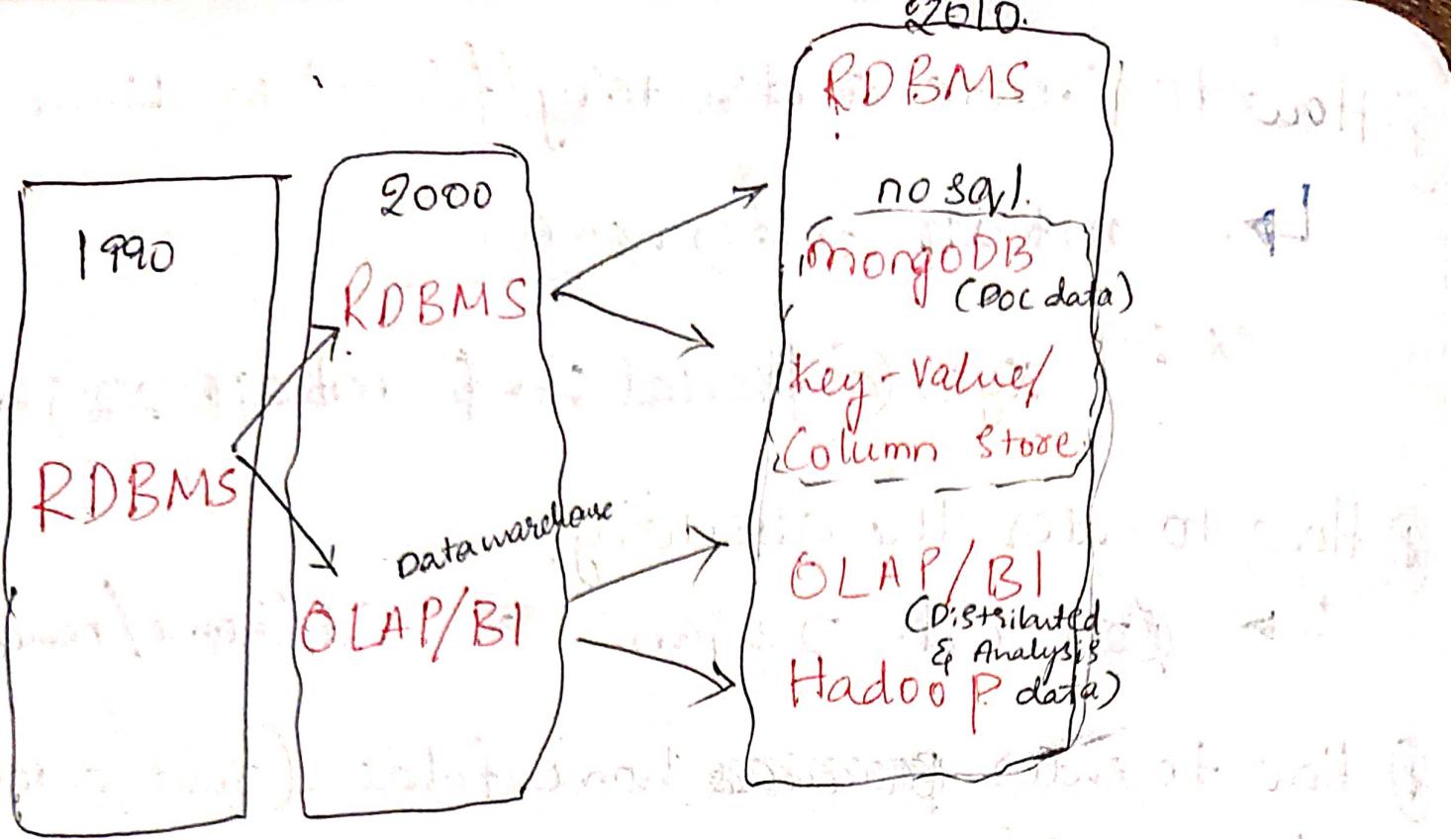
$$rw \rightarrow (4+2)=6$$

$$rx \rightarrow (4+1)=5$$

$$wx \rightarrow (2+1)=3$$

$$rwx \rightarrow (4+2+1)=7$$

Values.



Transition from databases to data warehouse to datalakes

What is Big Data?

* V's of BD

Volume, speed at which data is generated

Velocity, speed at which it is to be analyzed for actionable insights

Variety, Veracity

Structured
(in rows & columns)

Unstructured
(email, photo, etc)

Semi-structured
(in b/w)

Untrusted, Uncleaned

Big data is high volume, high Velocity, and/or high-variety information assets that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, & process automation.

Volume, velocity, variety, veracity

New Type of Data.

Sentiment	Clickstream.	Sensors	Geographic
understand how.	Capture and analyse (website visitors data trails etc) Optimize your web	Discovers the patterns	Analyse locations based data
Server logs	unstructured		
Research logs	web pages, emails, docs, etc		

Big Data. — model building Issues.

④

★ Hadoop :- (OS)

Hadoop is essentially something, where we can learn software on diff machines since we have huge amount of data and we have to process it.

★ Hadoop is an open source framework, capable of processing large amounts of heterogeneous data sets in a distributed fashion across cluster of commodity computers & hardware using a simplified programming model.

→ Hadoop framework is based on following principle:-

★ In pioneer days they used oxen for heavy pulling, when the load got heavier, they did not grow a larger ox, we should be trying for bigger computers, but for more systems of computers.

HDFS and MapReduce

HDFS → Reliable Shared Storage.
(Hadoop distributed File System)

Storage.

Distributed Storage

All Systems are interconnected to each other

MapReduce → Distributed Computation.

Hadoop

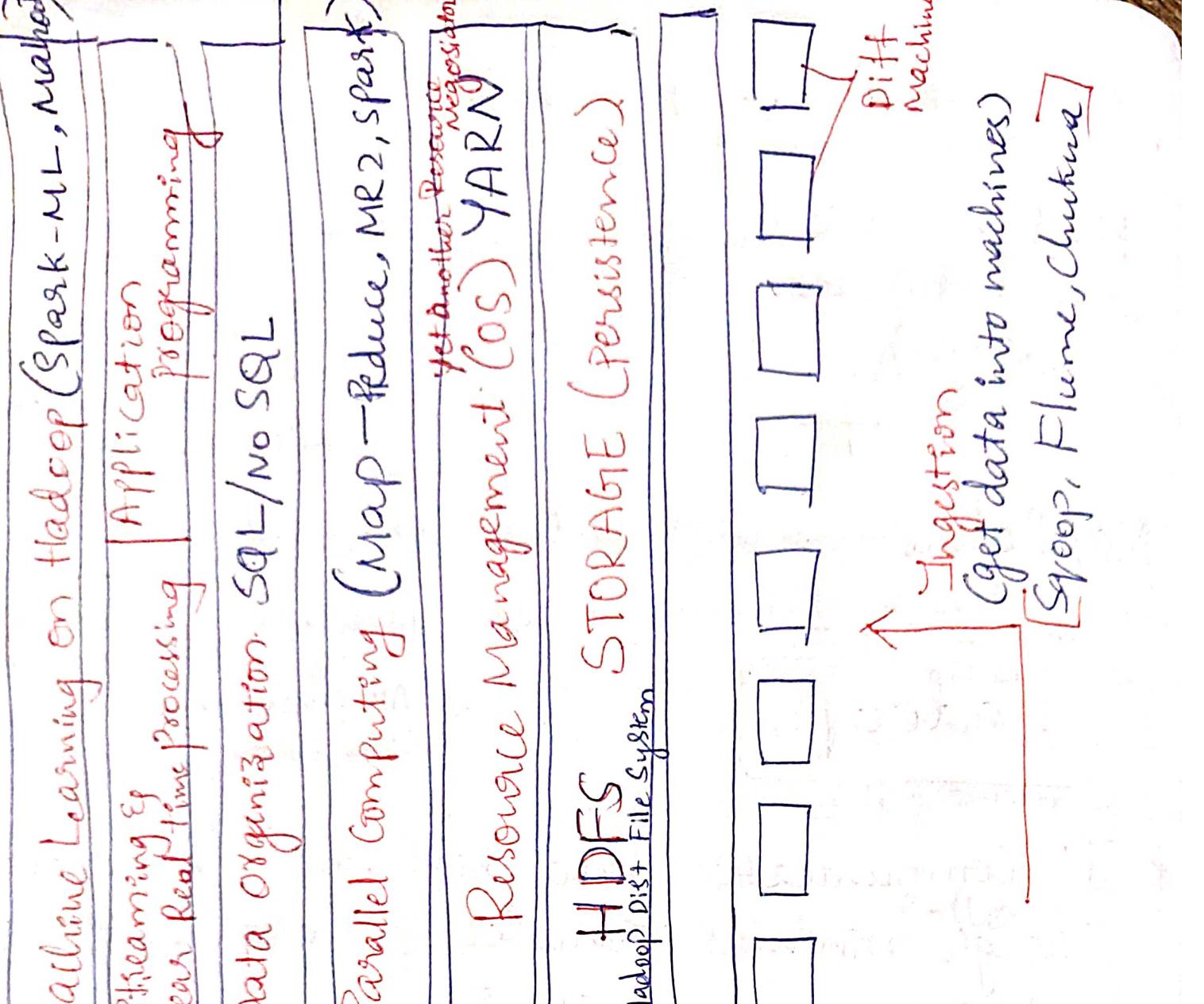
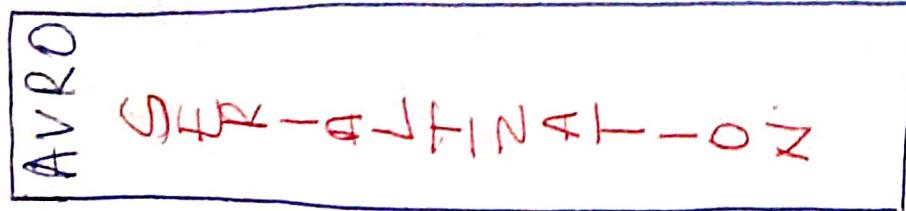
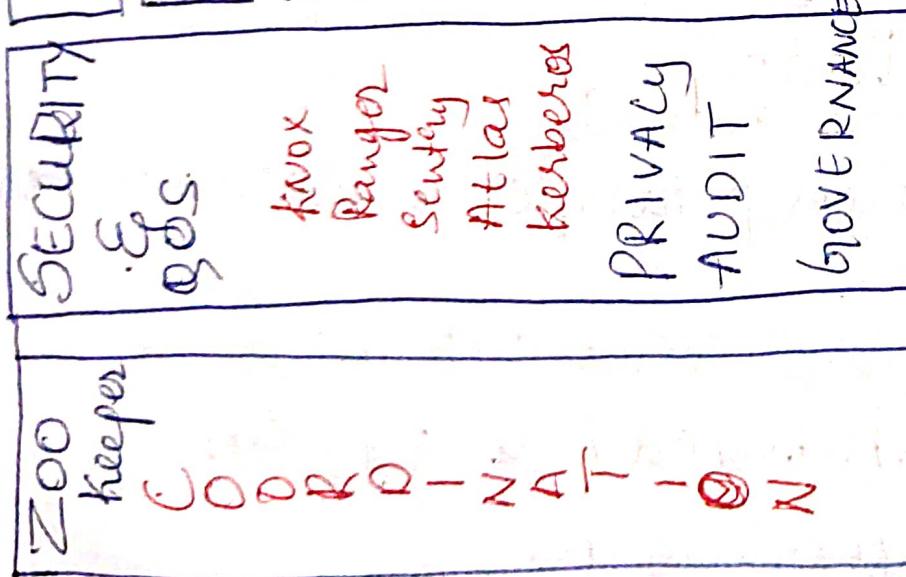
Parallel processing

All process happening in same time.

- ④ It communicates b/w diff nodes (machines) through minimal transfer of data and high level programming is done.

Why do we need hadoop :-

- ④ We want to process petabyte of datasets very quickly
- ④ Data may not have strict Schema.
- ④ Very Large distributed File System.
- ④ Runs on heterogeneous OS.
- ④ Fault tolerant.



Major components of Hadoop :-

① The Storage Layers

(GFS → Google File System)

HDFS → Hadoop Distributed File System

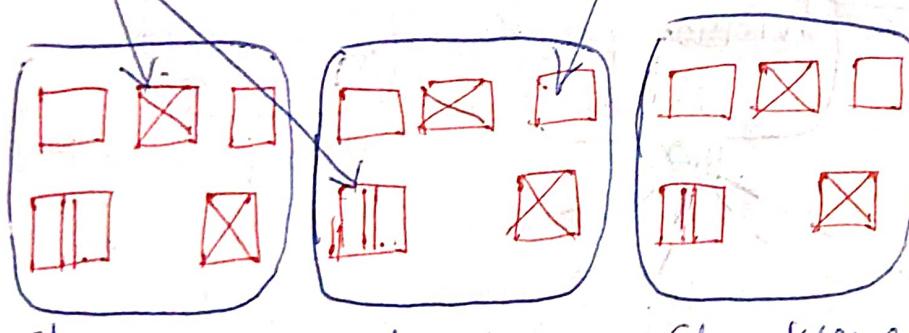


file

A file....

is made of 64 mb
Chunks

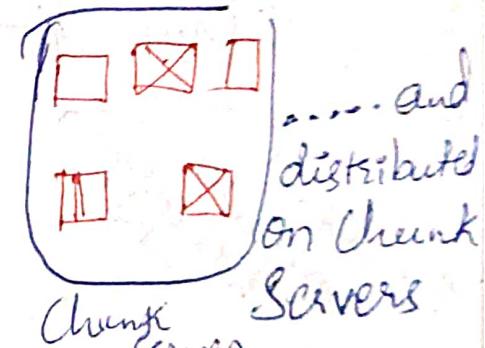
... that are replicated
for fault tolerance



chunkserver

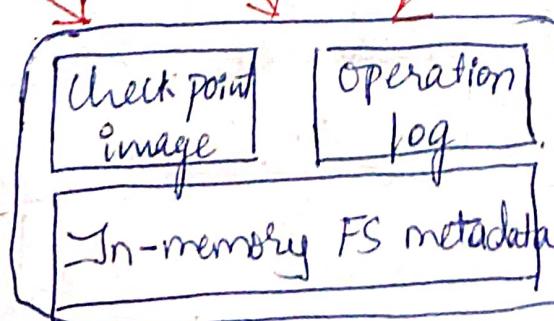
chunkserver

chunkserver



chunk
server

Master

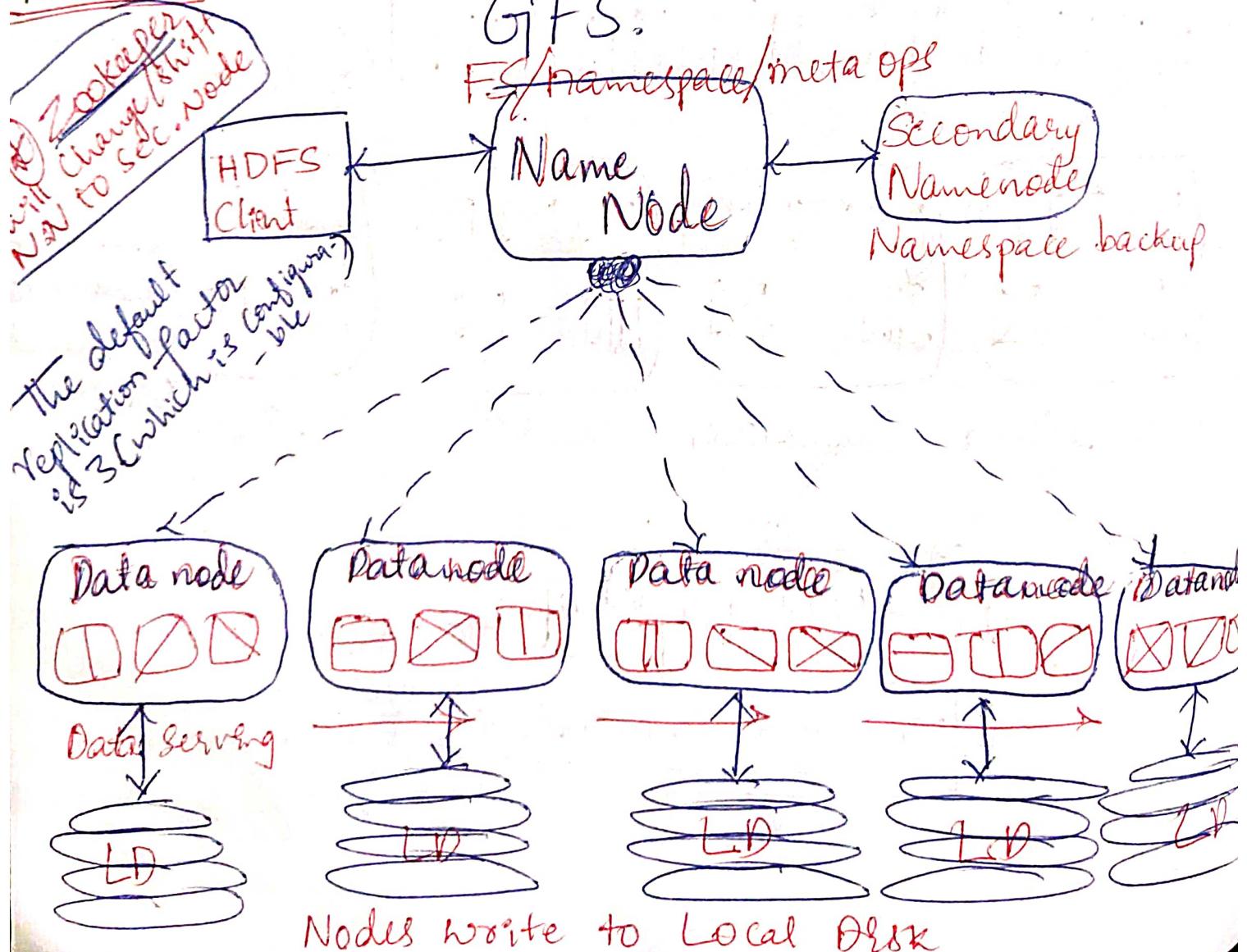


The master
manages the
file system
name space.

GFS Master Responsibilities

- ★ Metadata Storage (Information of the data in DataNodes)
- ★ Namespace management / locking
- ★ Periodic communication with Chunk Servers
- ★ Chunk Creation, re-replication, rebalancing
- ★ Garbage Collection.
- ★ "stale replica" deletion.

HDFS || CDH3: Open-Source Implementation of GFS.



④ Name node manages overall file system meta data
↳ It can be a single point of failure in the system

④ Datanodes (one per machine) stores and retrieves actual data
↳ Reports to ~~data~~ name node.

④ Both datanode & Namenode include a webserver, so node status can be easily checked.

④ Secondary namenode: does housekeeping & acts as a memory backup.

④ Namenode is "rack-aware" → Knows how machines are arranged.

④ Namenode does not directly read & write data.

↳ Clients gets data location from namenode
↳ Clients interacts directly with datanode to read/write data.

④ Name node keeps all block metadata in fast memory.

④ For every 3 seconds every datanode sends the signal to namenode with the data intact, this is called Heartbeat.

Rack Awareness :-

(in-built Rack Awareness Algorithm.)

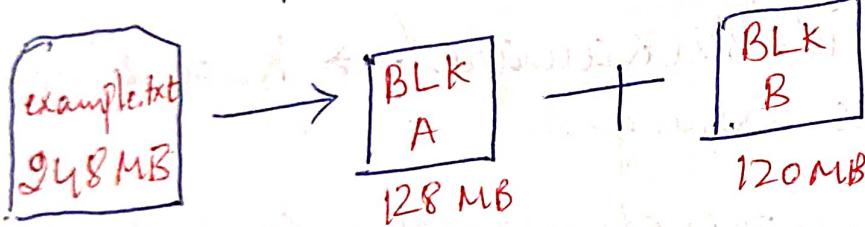
④ Hadoop has the Rack Awareness Algorithm in it.

↳ It improves the network performance

↳ To prevent loss of Data

HDFS write Architecture :-

④ Write once, Read multiple times
Storing the data



④ At first, the HDFS Client will reach out to the name node for a write request against the two blocks, say Block A & Block B

④ The name node will then grant the client the write permission and will provide the IP addresses of the Data Nodes where the file blocks will be copied eventually.

For Block A, List A = {IP of datanode 1, IP of DN 4, IP of DN 6}
For Block B, Set B = {IP of dn 3, IP of DN 7, IP of DN 9}

There are following steps to undergo this:

- ① Setup of Pipeline (Path).
- ② Data streaming and replication
- ③ Shutdown of pipeline (Acknowledge^{new stage})

- ④ Setup of pipeline means the NameNode will give the IP address of the datanodes to the Client for the storage.
- ⑤ Data streaming and replication includes the writing of the data in the specified datanodes.
- ⑥ Shutdown of pipeline means intimating the NameNode that the writing process has been completed and then the NameNode stores it in metadata.

HDFS Read Architecture:- It is very simple, the Client must specify the blocks (data) to be read, the NameNode will find the location of the blocks (if they are in one Datanode or not), then sends the IP address of the datanodes to read the data to clients.

MapReduce

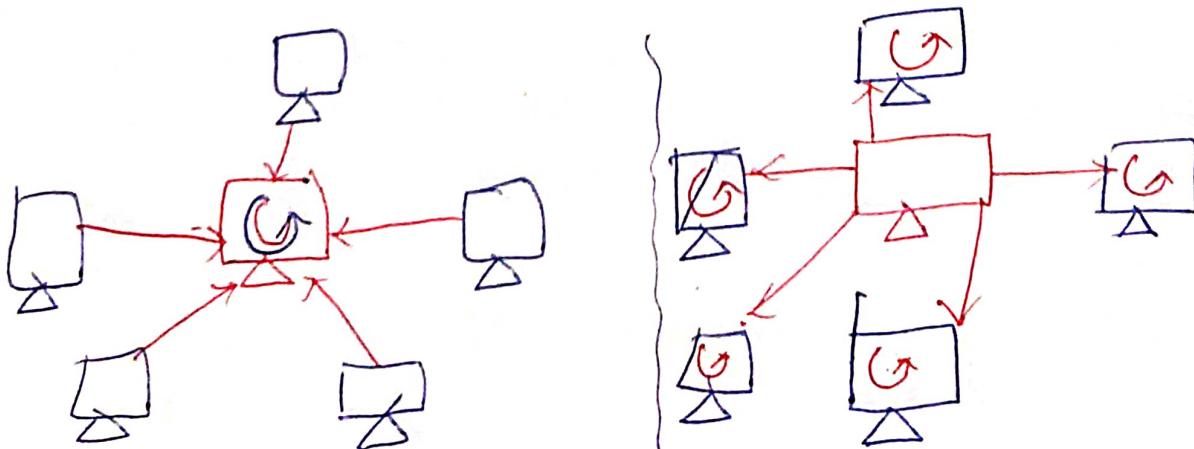
(F) MapReduce is used for batch processing (historical data)

Drawback

- (A) MapReduce is a programming model for processing large data sets with a parallel and distributed algorithm on a cluster.
- (B) MapReduce when coupled with HDFS can be used to handle big data. The fundamental of this HDFS - MapReduce system, which is commonly referred as Hadoop.
- (C) The basic unit of information, used in MapReduce is a (key, value) pair. All type of structured and unstructured data need to be translated to this basic unit, before feeding the data to MapReduce model.

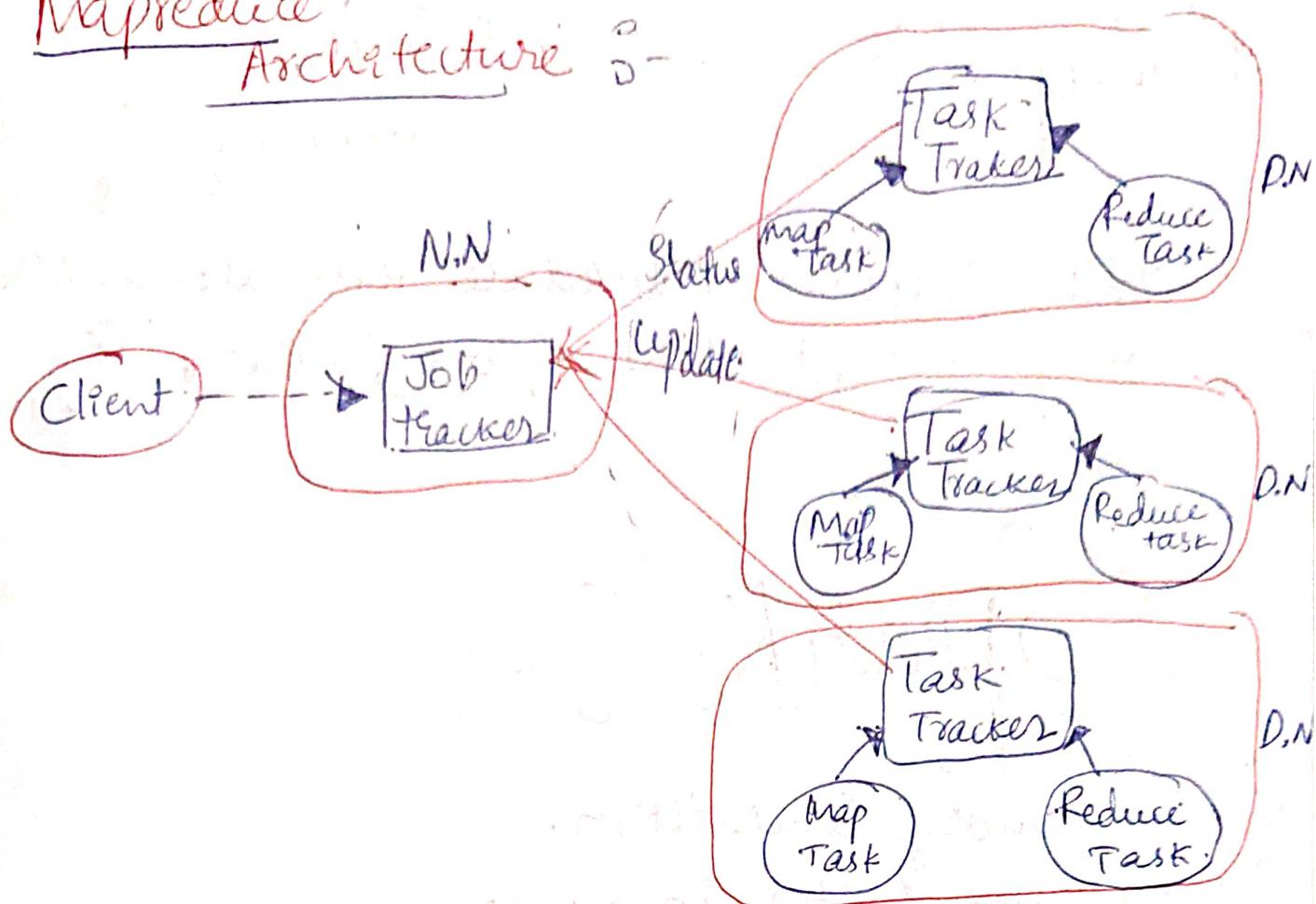
Advantages of MapReduce :-

① Parallel Processing :-



Moving data to Processing unit, moving Processing unit to the data
(Traditional Approach) (MapReduce Approach)

Mapreduce Architecture :-



Drawbacks

- ④ Batch Processing
- ④ not suitable for Real time Data / Data streaming
- ④ It supports upto 4000 Nodes per cluster
- ④ It has a Single component: Job tracker to perform many activities like Resource management, Job Scheduling, Job Monitoring, Re-scheduling Jobs etc
- ④ Job tracker is single point of Failure
- ④ It runs only map/reduce jobs
- ④ It follows slots concept in HDFS to allocate resources.

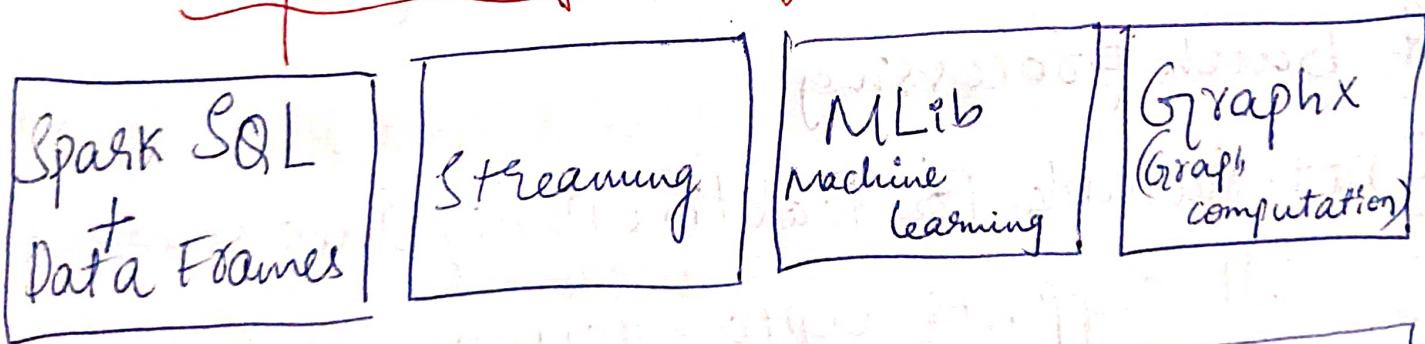
Hadoop 2.X. → YARN
(Yet Another Resource Negotiator)

- Yarn performs Mapreduce jobs along with many other jobs like Spark, HBase, Tez etc

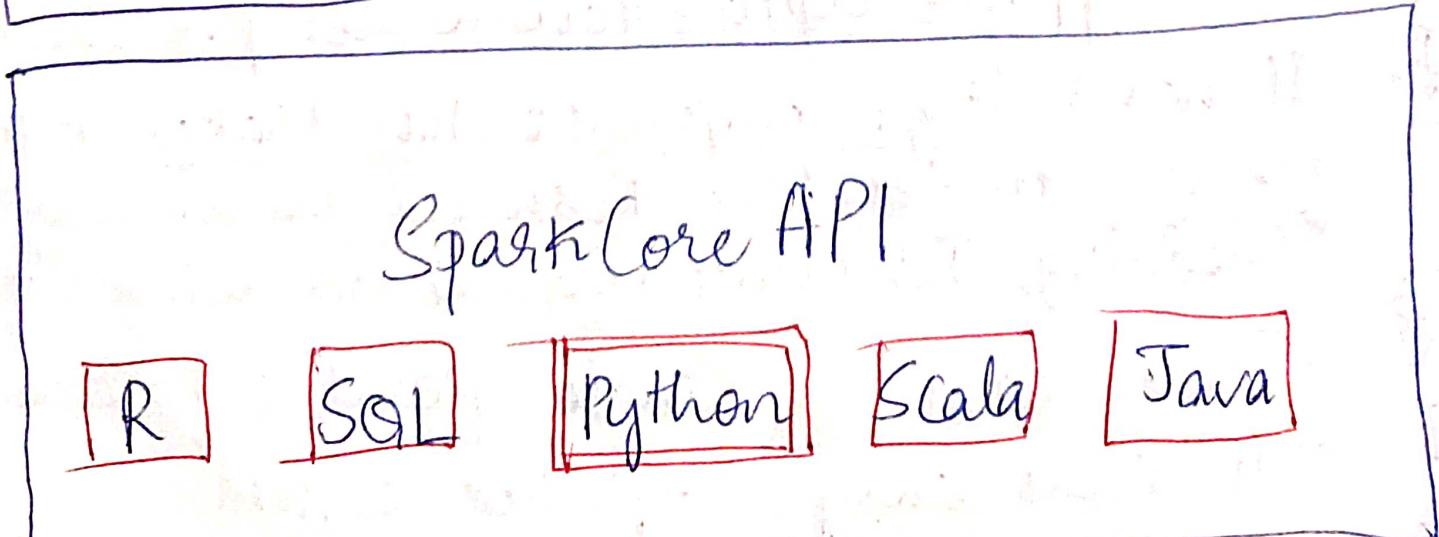
Spark

- In-memory computation.
- 100 times faster in memory
- Supports → ML, SQL, streaming, Graph algorithms

Spark Ecosystem

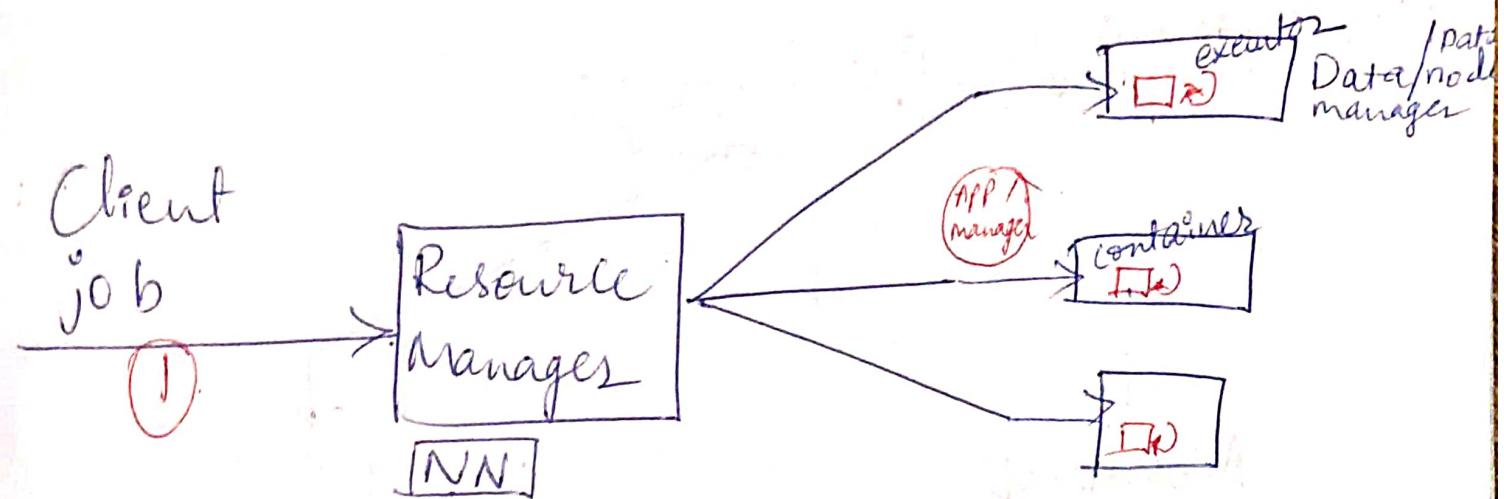


Spark Core API



Spark & YARN

Architecture.



Steps

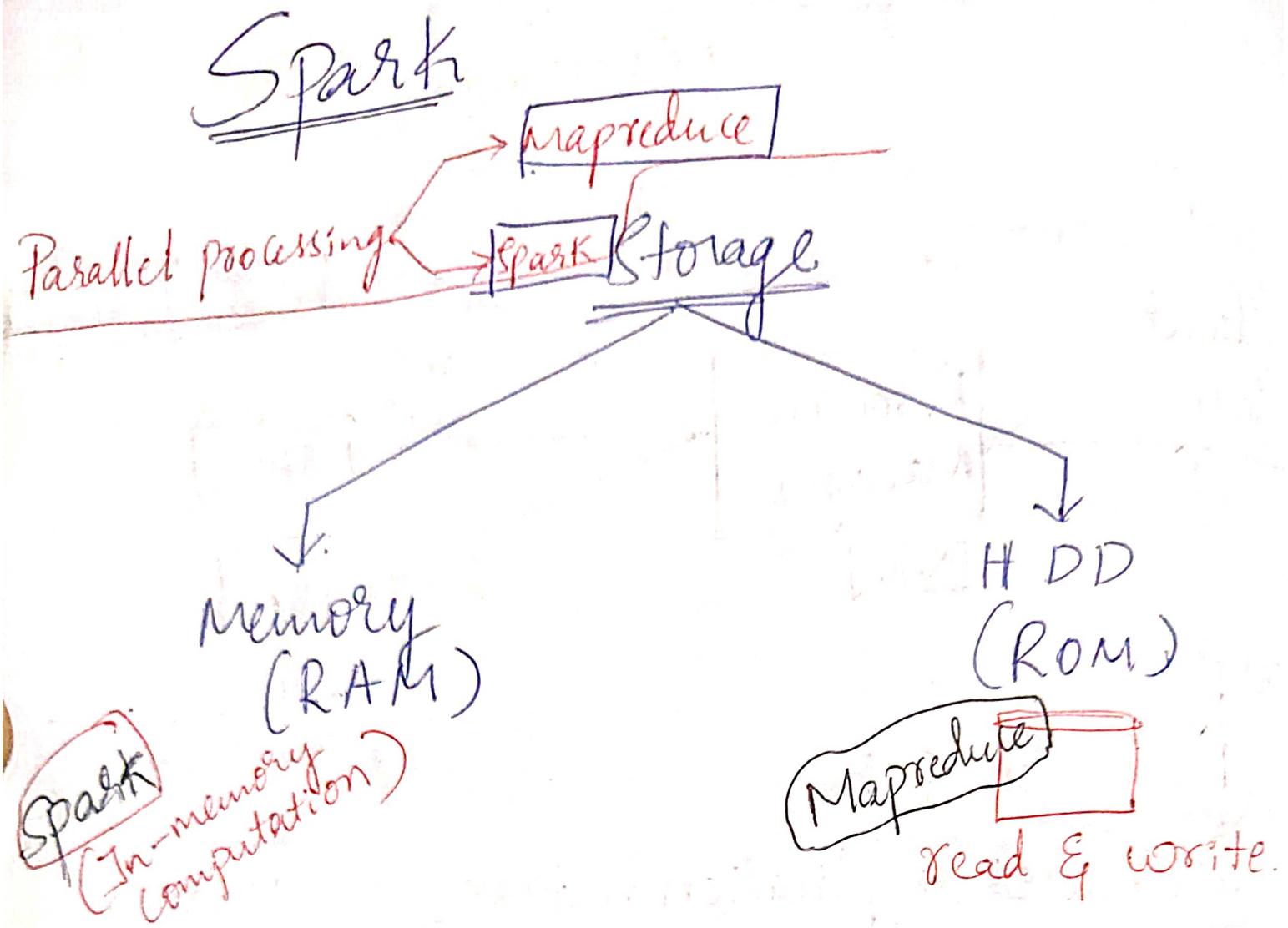
- ① Client job
- ② Creates application master
- ③ R.M creates Container(memory location)
- ④ Application master sends to the job to Container & executes the job
- ⑤ And the job is sent back to Client

Yarn

- ① It will handle Resource manager
- ② Allocates memory location, Application manager etc.

Spark

- ④ Spark says "what job to be done!"
- ⑤ It is in memory computation.
- ⑥ It executes the job in container



④ Spark-execution :-

Spark → R.DD $\xrightarrow{\text{stores it to}}$ DF \longrightarrow SQL.

There are 2 operation

① Transformation.

② Action.

This submits the job to YARN.

④ It will assign executor to solve this job using application manager.

YARN → Resource manager layer

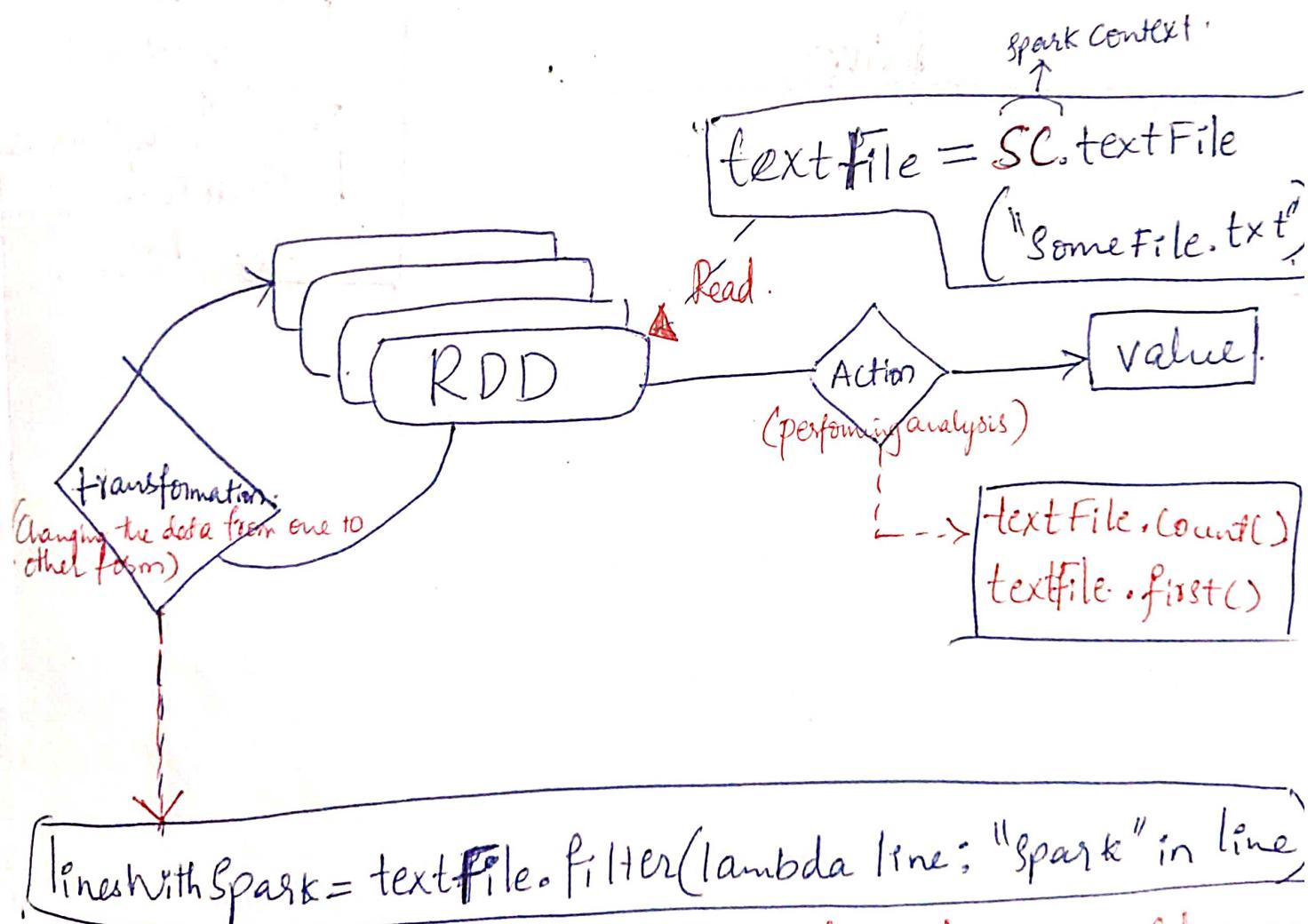
→ Resilient Distributed Dataset (RDD) - Key Spark Construct

↳ RDDs represent data or transformations on data.

↓
Sum, Percentage
or Avg, etc

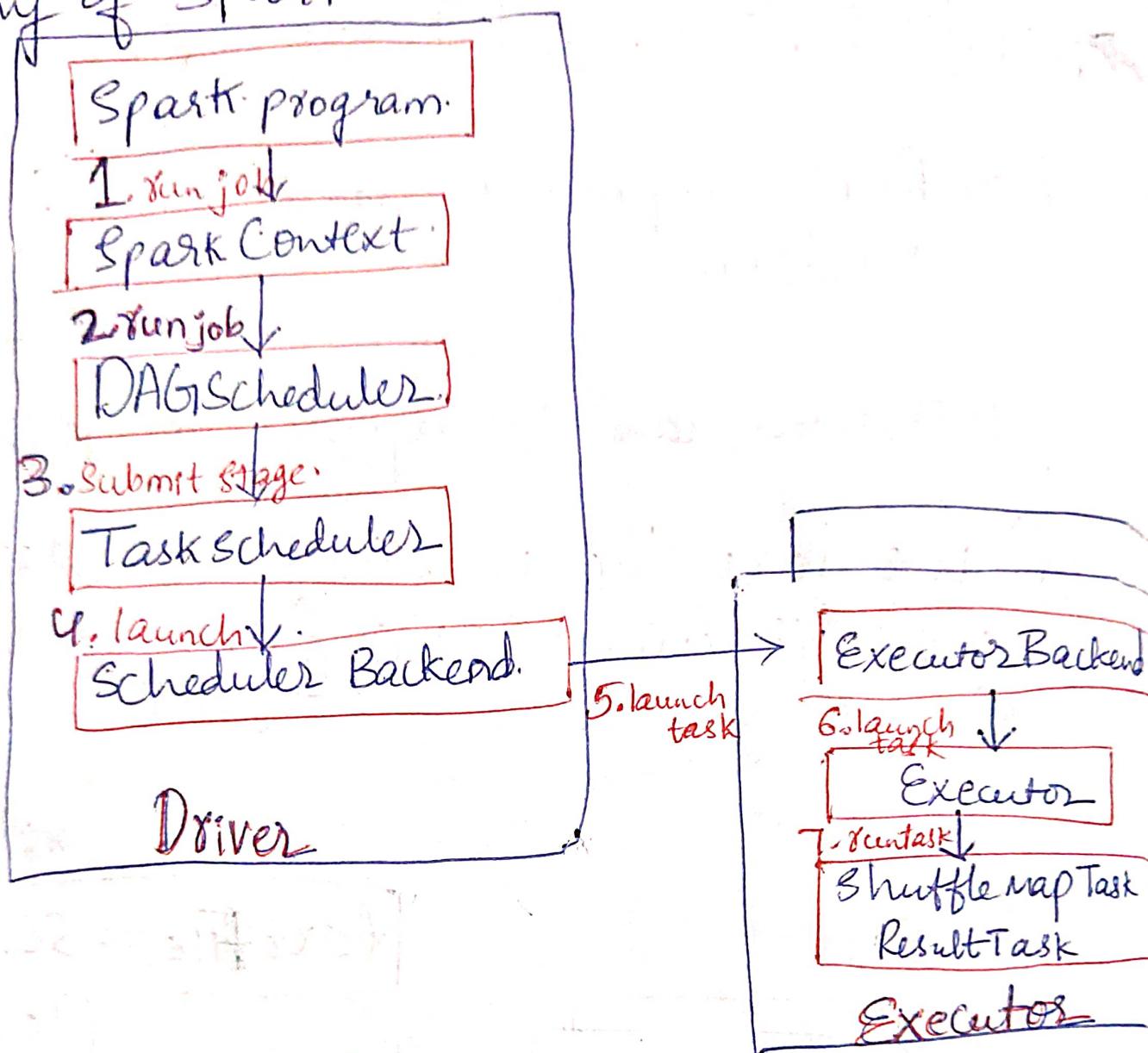
↳ Actions can be applied to RDD.

④ we load the unstructured data and get structured data (pre-processing)

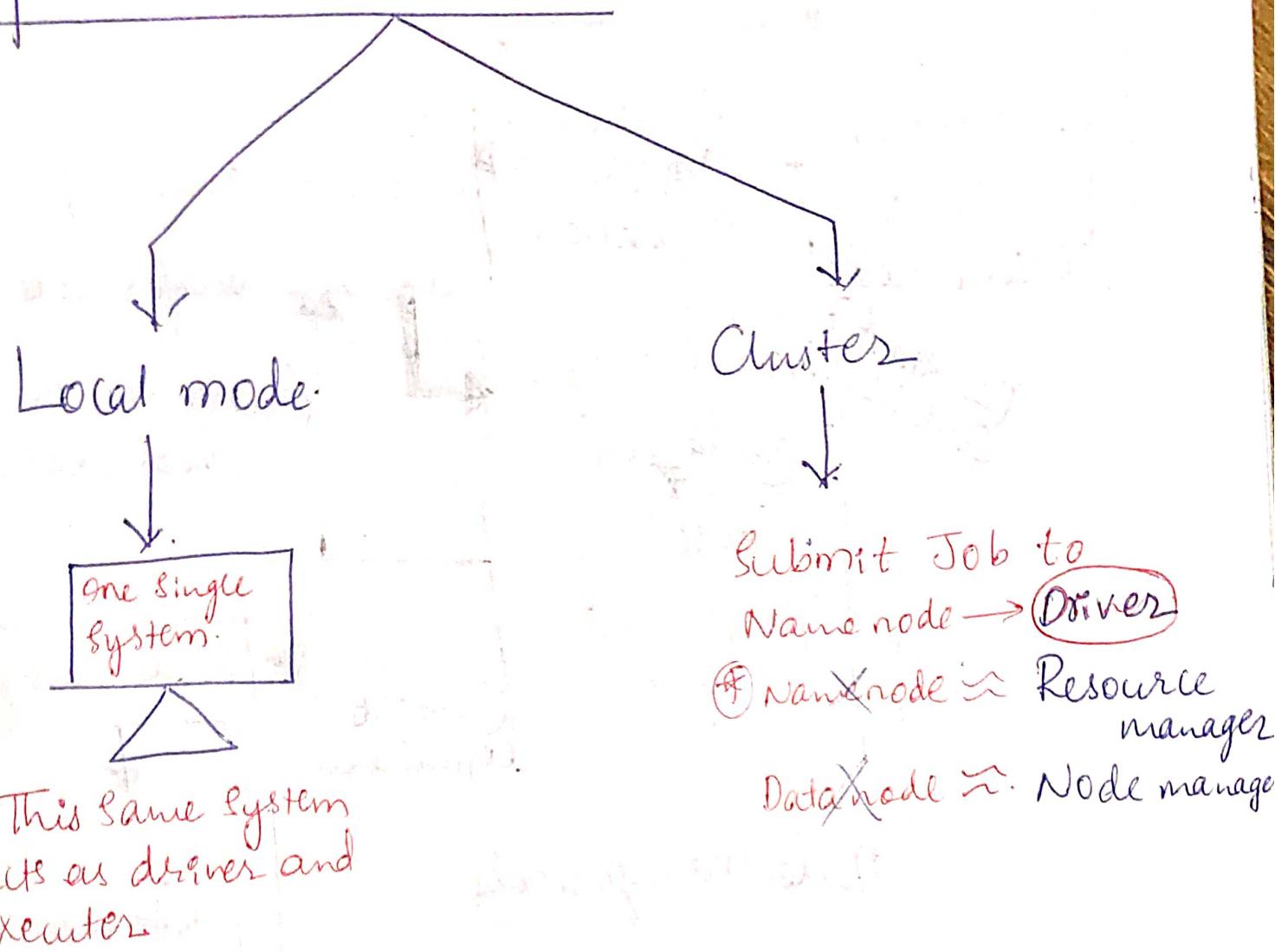


④ we can modify RDD, therefore for any change we do transformation and store as new RDD.

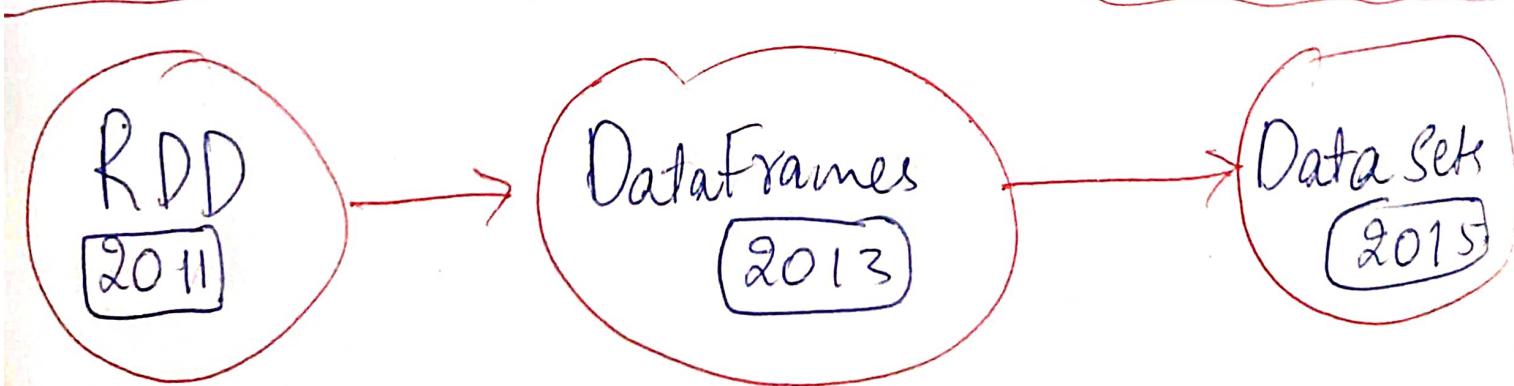
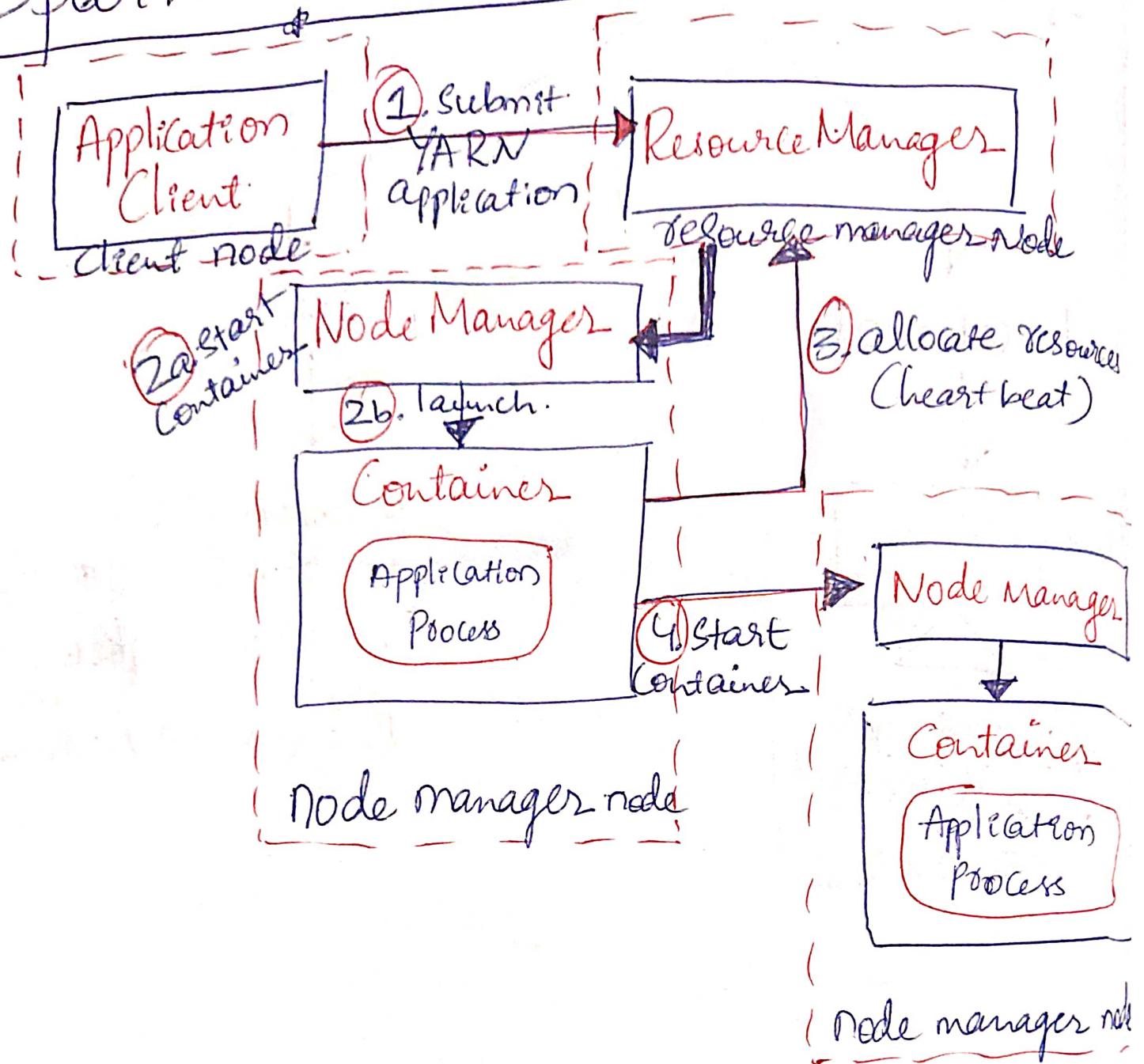
Anatomy of Spark Job Run :-



Spark Job Execution:



Spark on Yarn



Spark+ML :-

Dataset (Given)



We need to build the model.

① preprocessing + ML model.
 (Spark)
∴ Spark ML.