

## Why Big Data now?

---

- More data are being collected and stored
- Open source code
- Commodity hardware / Cloud

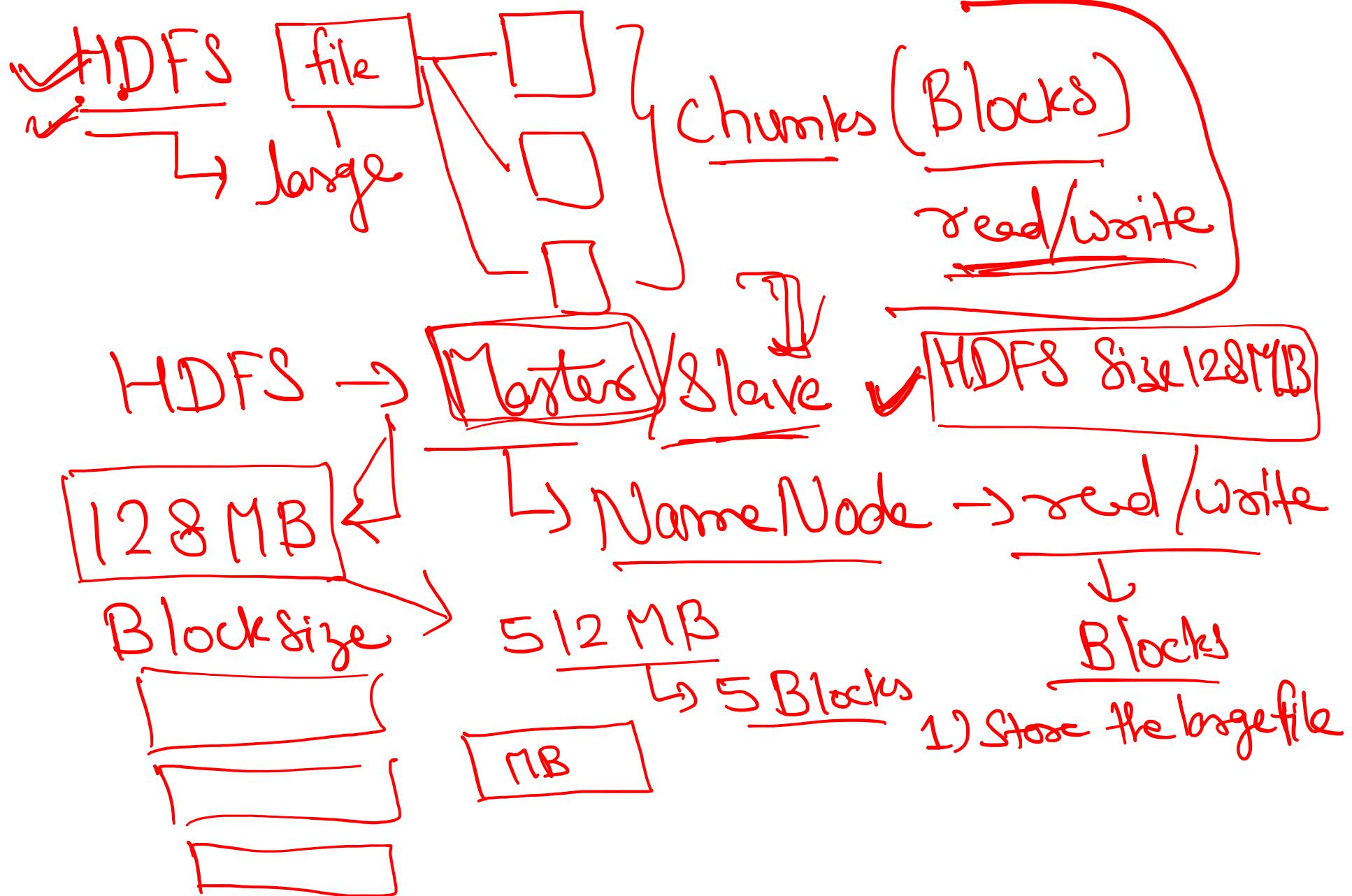
- 
- High-Volume
  - High-Velocity
  - High-Variety

TV's  
Value is big data

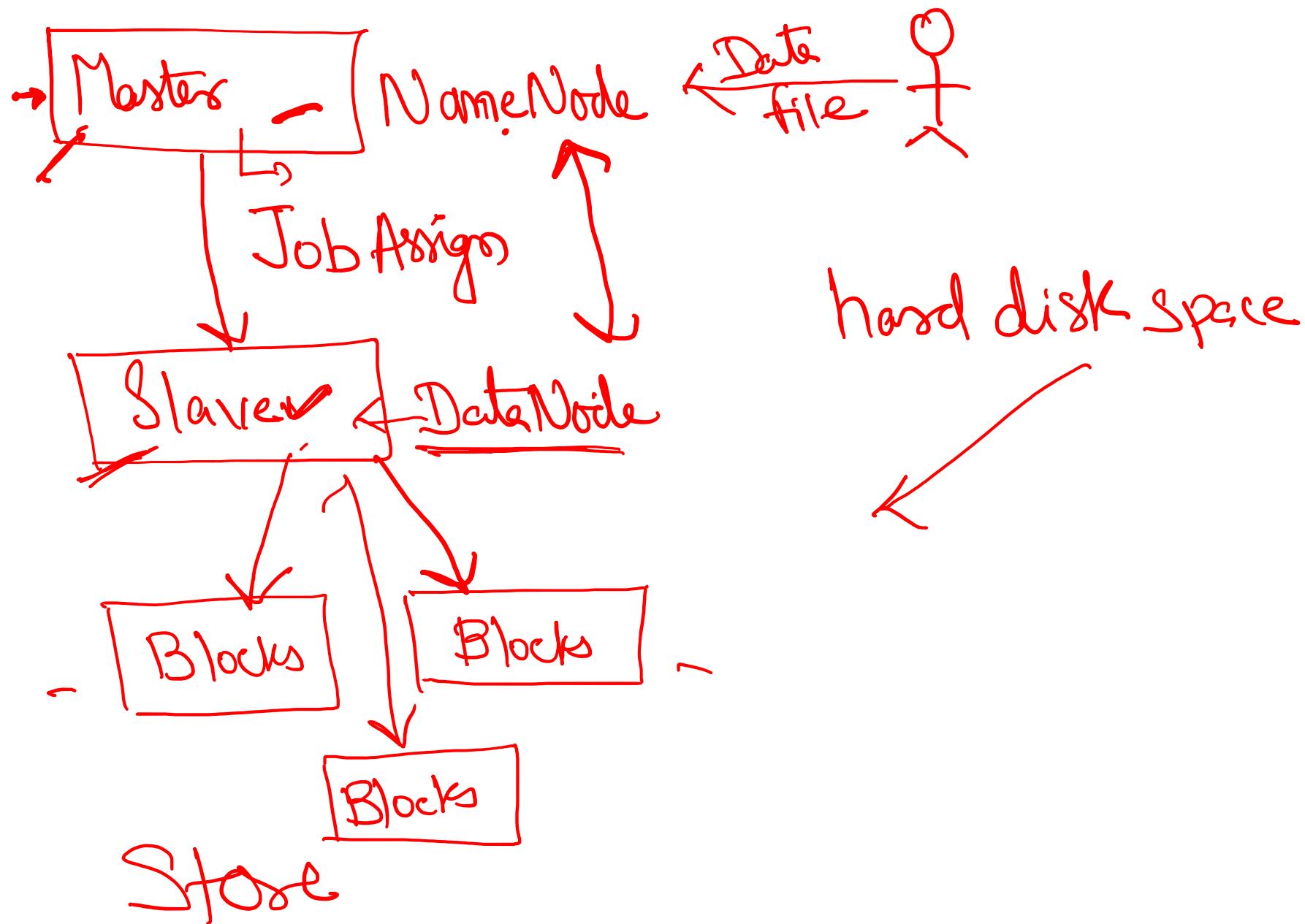
→ Artificial  
Intelligence

Hadoop Blocks – Namenode

~~Date~~ ✓



## Hadoop Blocks – Datanode ✓

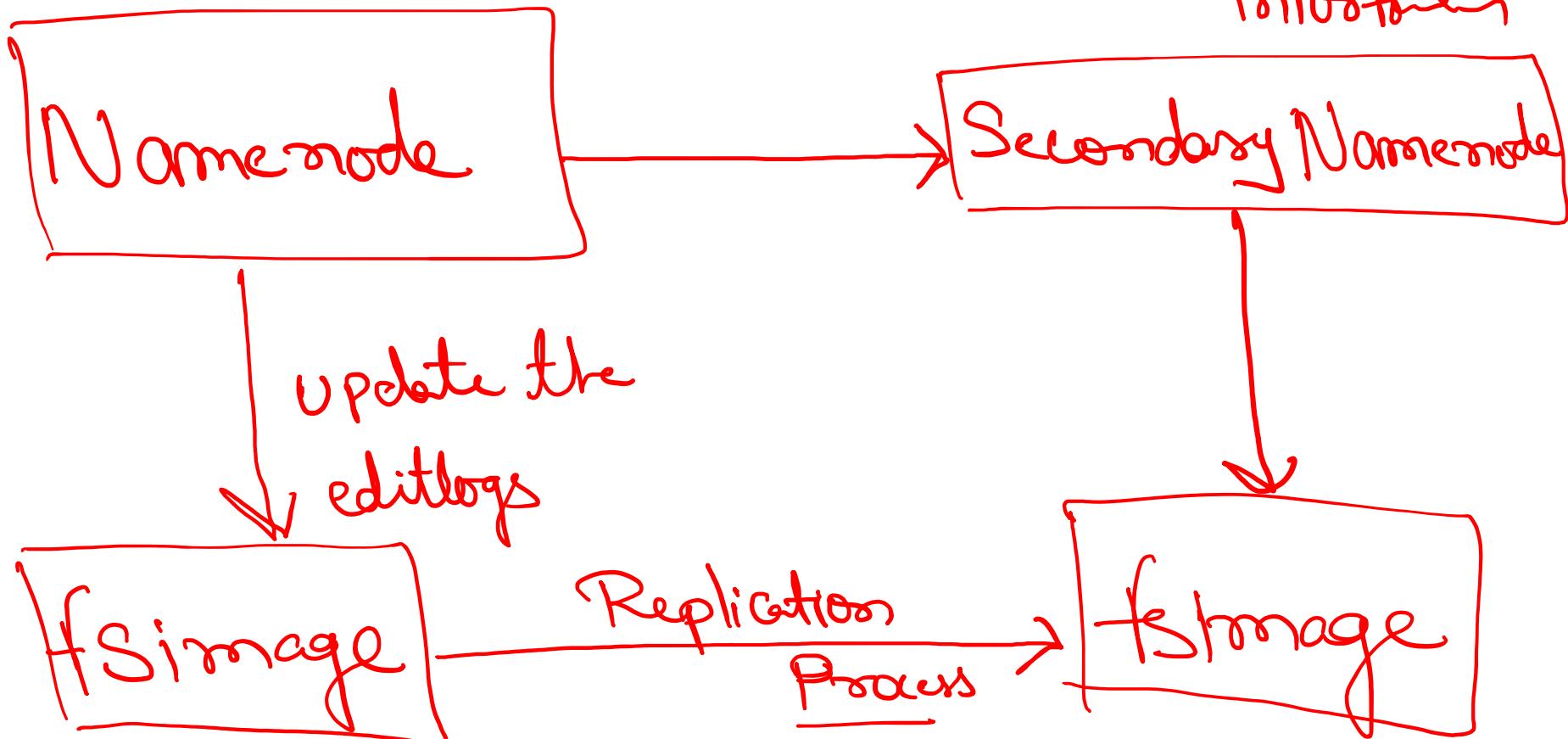


Hadoop Blocks – Secondary Namenode ✓✓

fstorage- Stores the snapshot of its info of all

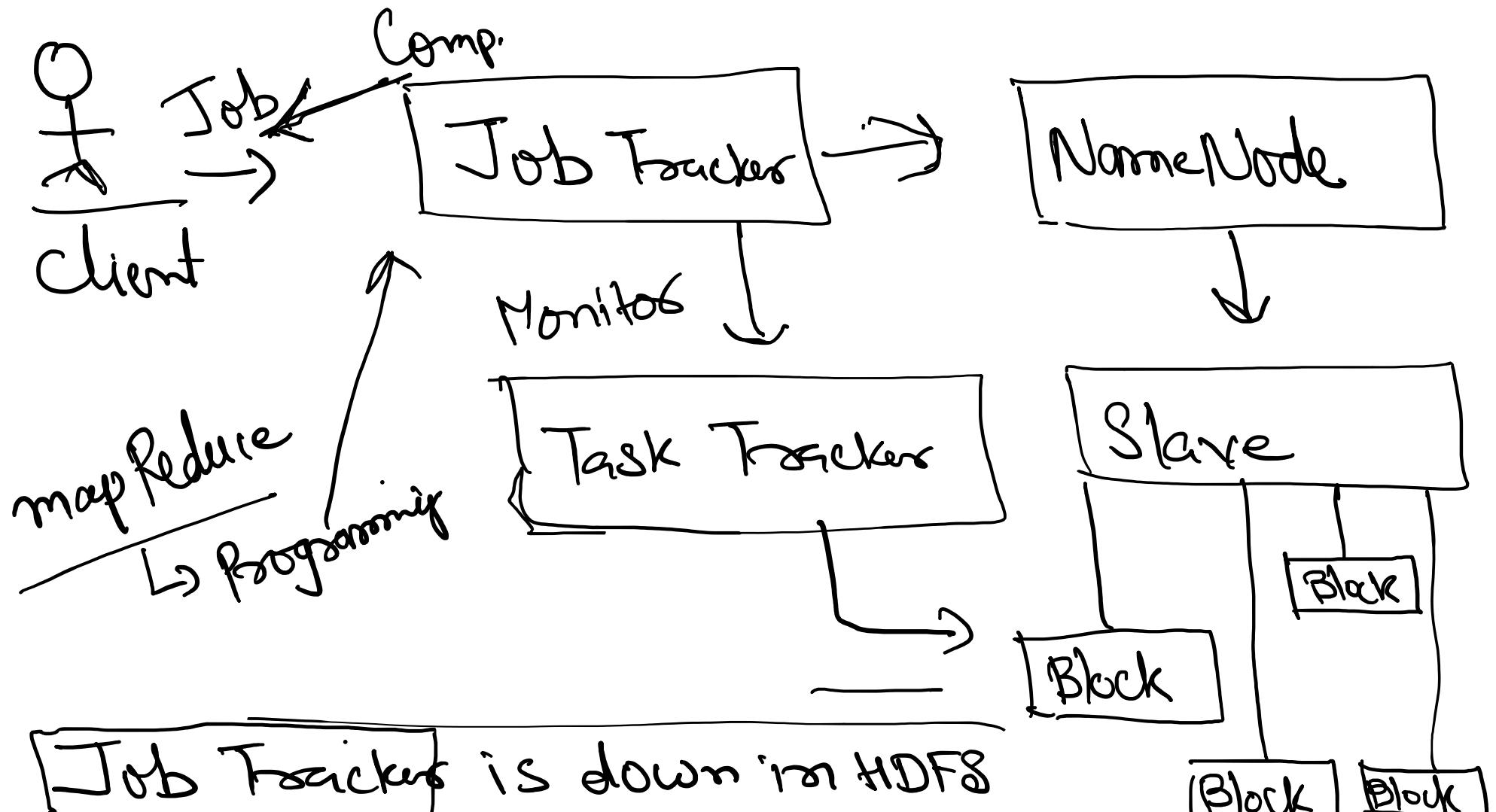
Namenode - holds metadata for HDFS in blocks

info of all



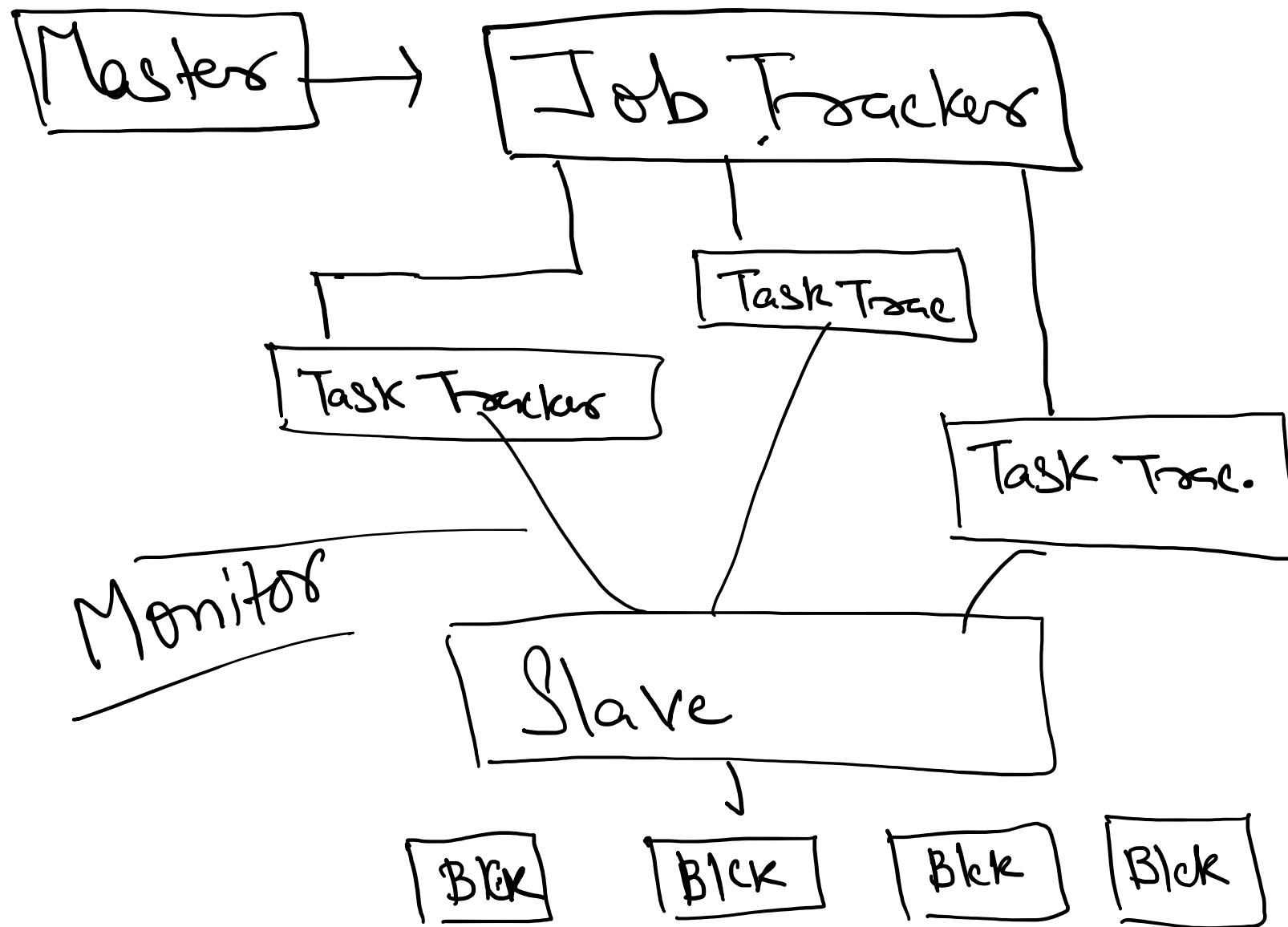
Editlogs → keeps track of every change of HDFS.

## Hadoop Blocks – Job Tracker ✓



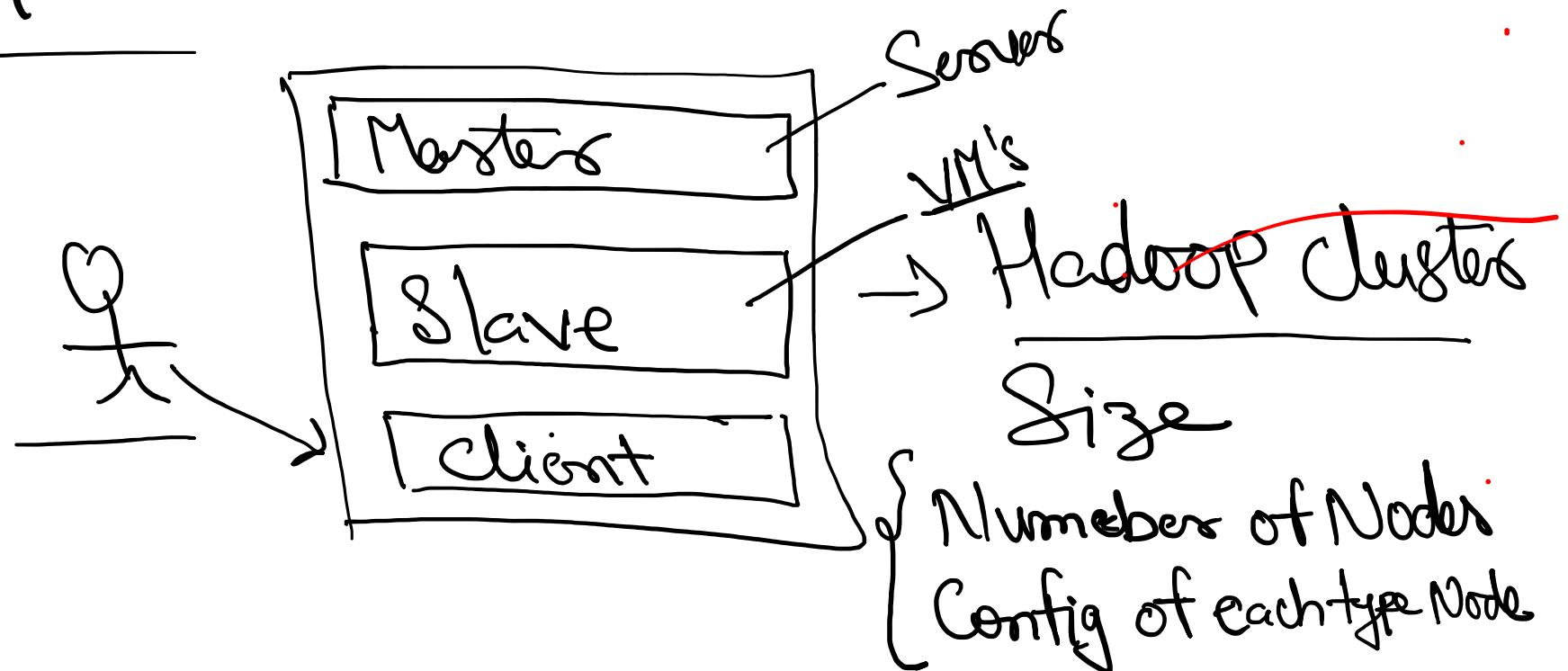
mapReduce → halt

## Hadoop Blocks – Task Tracker



## Hadoop Cluster - Introduction ✓

Apache ⇒ pig data analytics, process tasks to be broken down into smaller tasks that can be performed into parallel by using any algo like mapReduce.

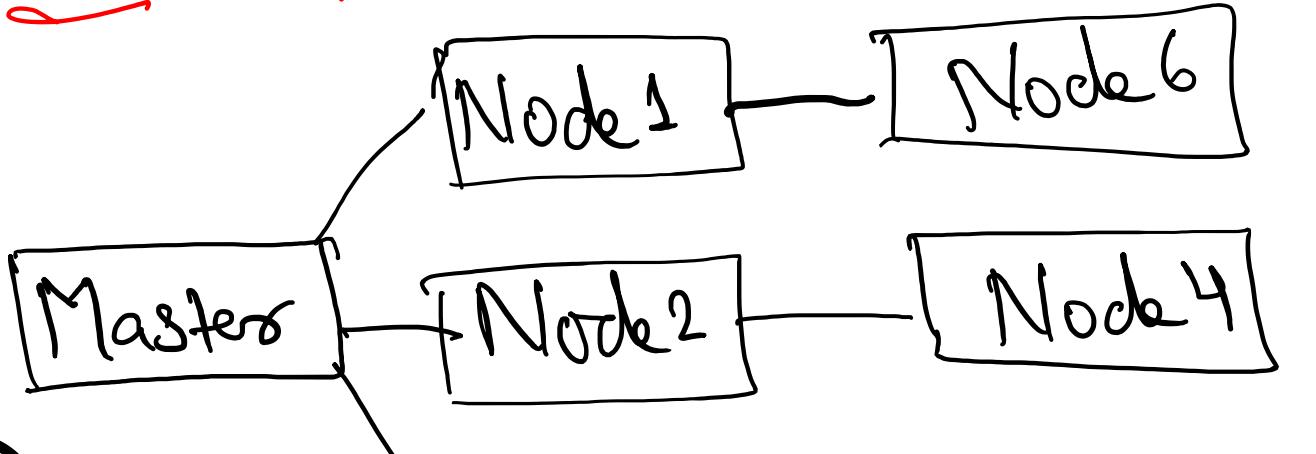


## Hadoop Cluster - Configuration

### Web Portal - Smap

Data Store - Struct, Semi-Struct & Unstruct Data

Schemas =>



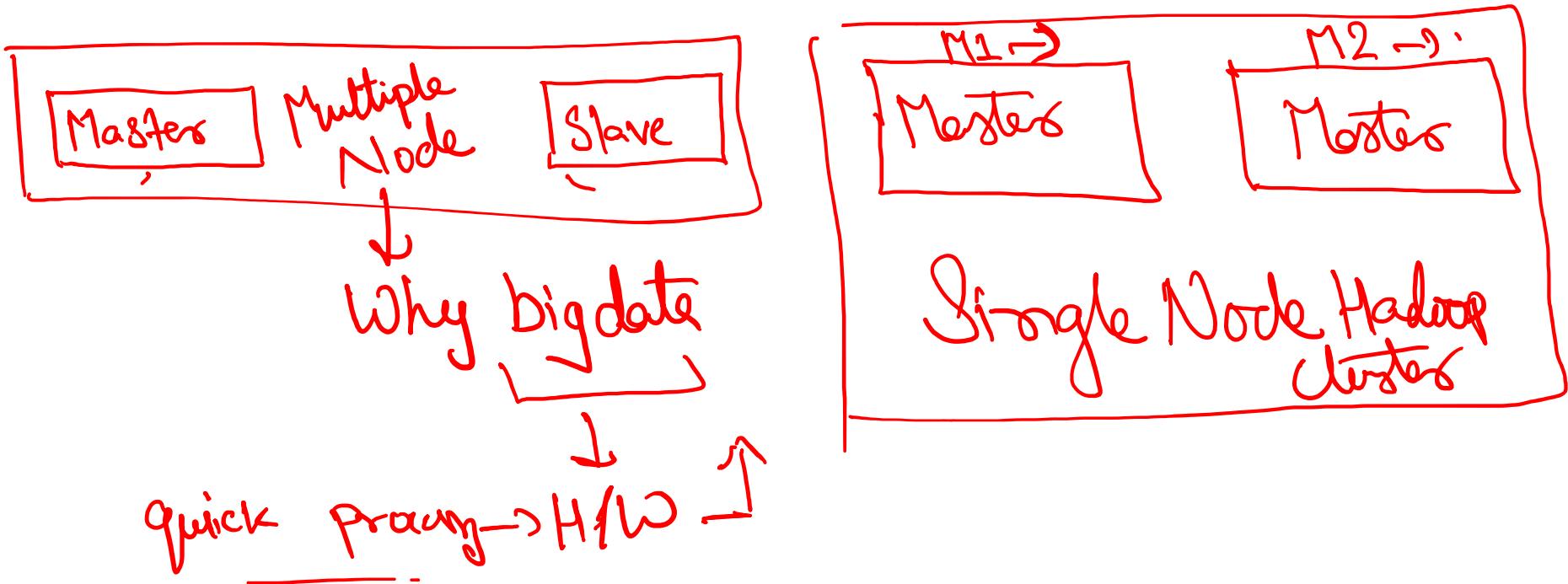
Properties ✓

- Scalability
- flexibility
- Speed
- No data loss
- Economical

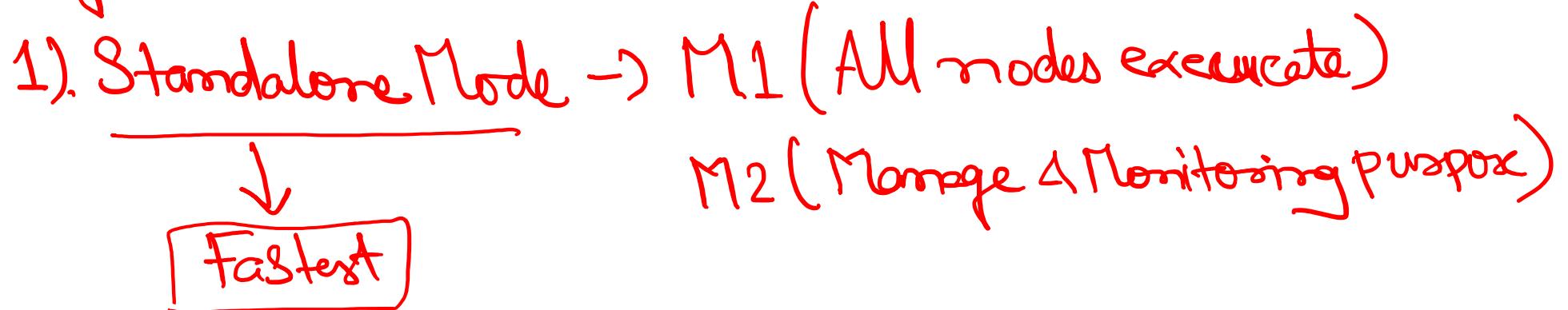
Slaves

## Type of Hadoop Cluster

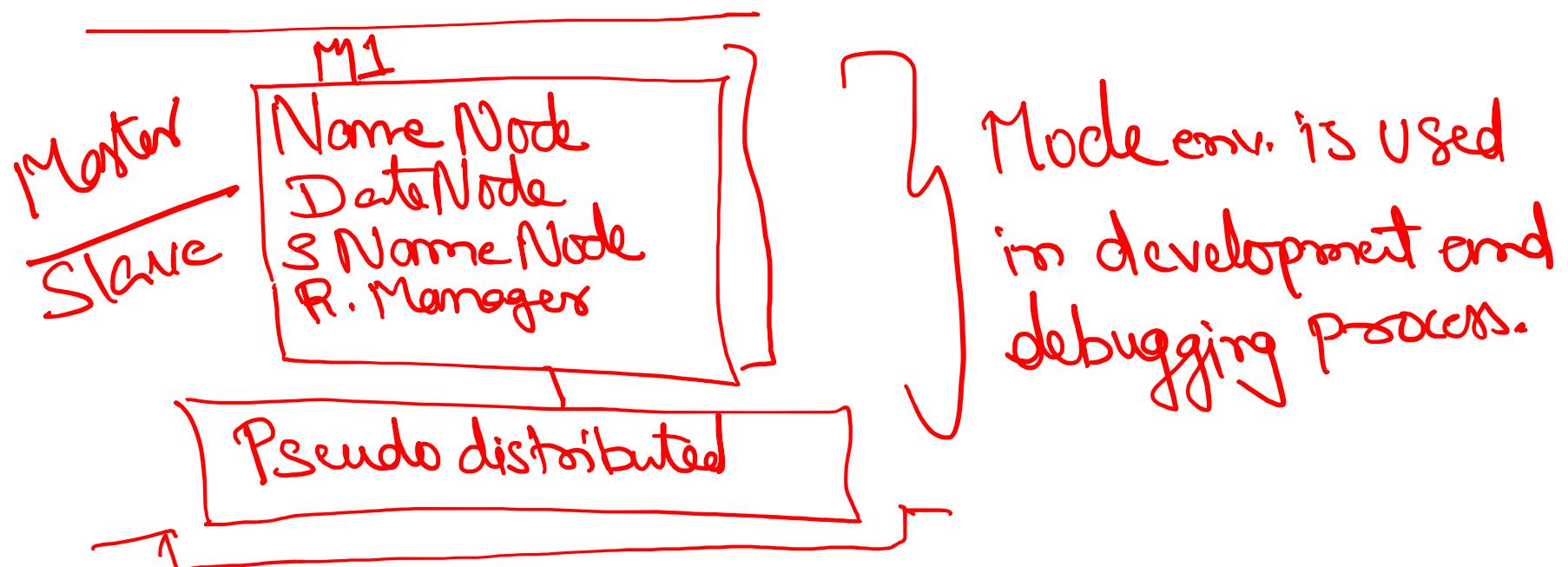
- 1) Single Node Hadoop | 2) Multiple Node



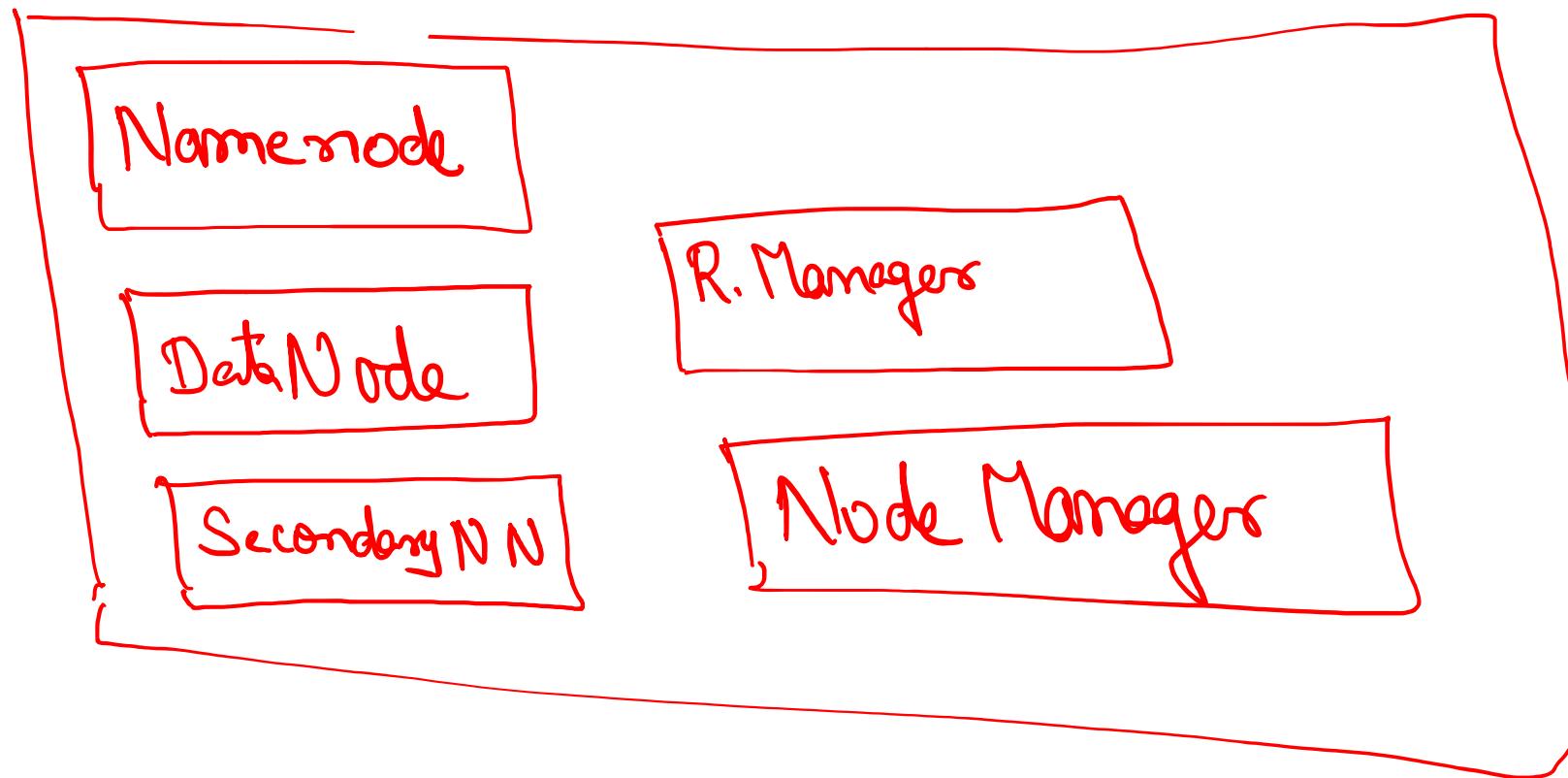
# Types of Mode of Operation



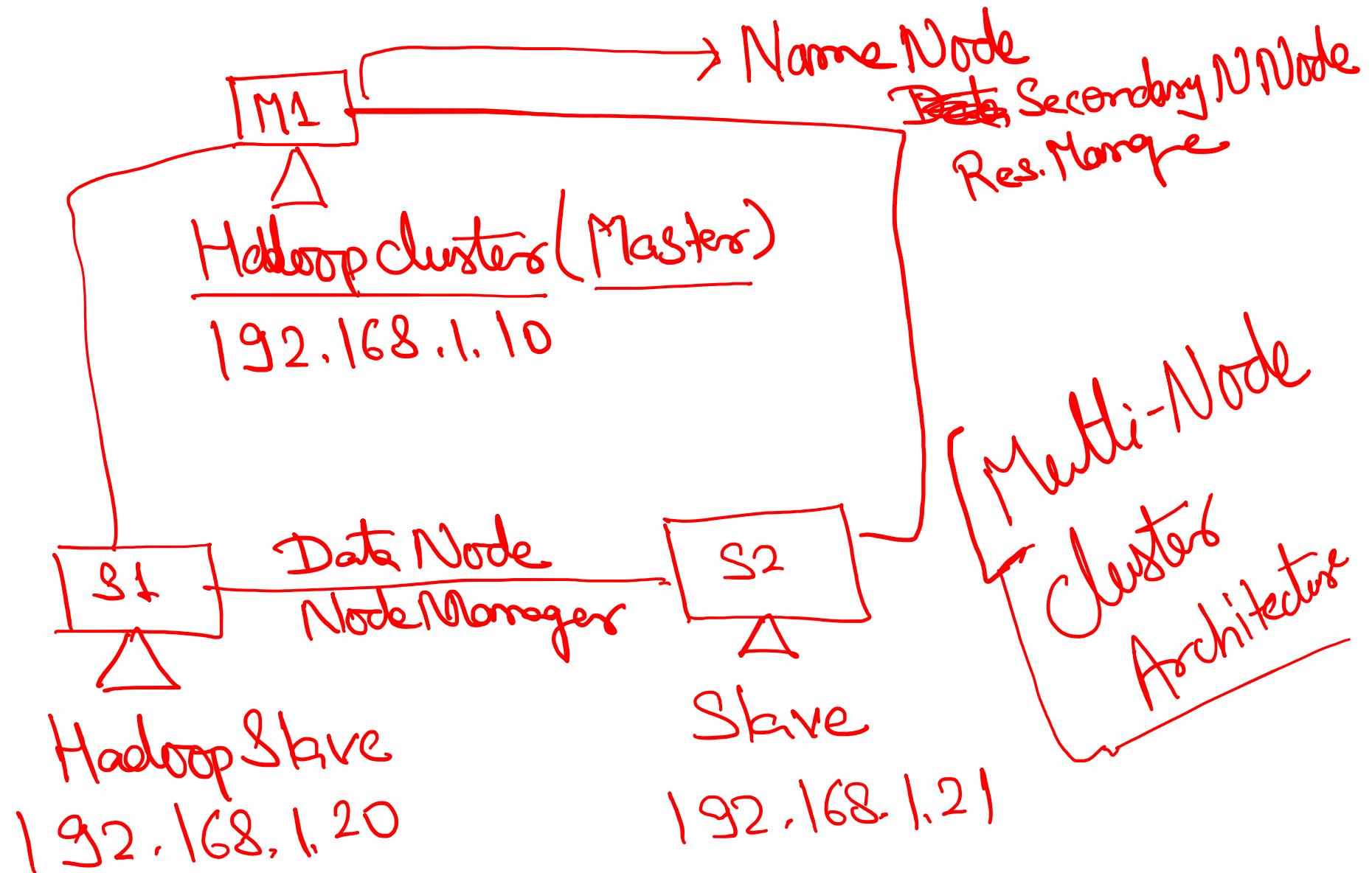
## 2) Pseudo Distributed Mode → ✓



### 3). Fully Distributed – Multi Node clusters(☞)



## fully Distributed Mode



Config file - [Imp] → Apache Hadoop docs.  
[SURRY]

1). HADOOP-ENV.sh → JDK

[\${JAVA\_HOME}]

2) CORE-SITE.XML → Cluster Info

↳ NameNode - IP, Port#,

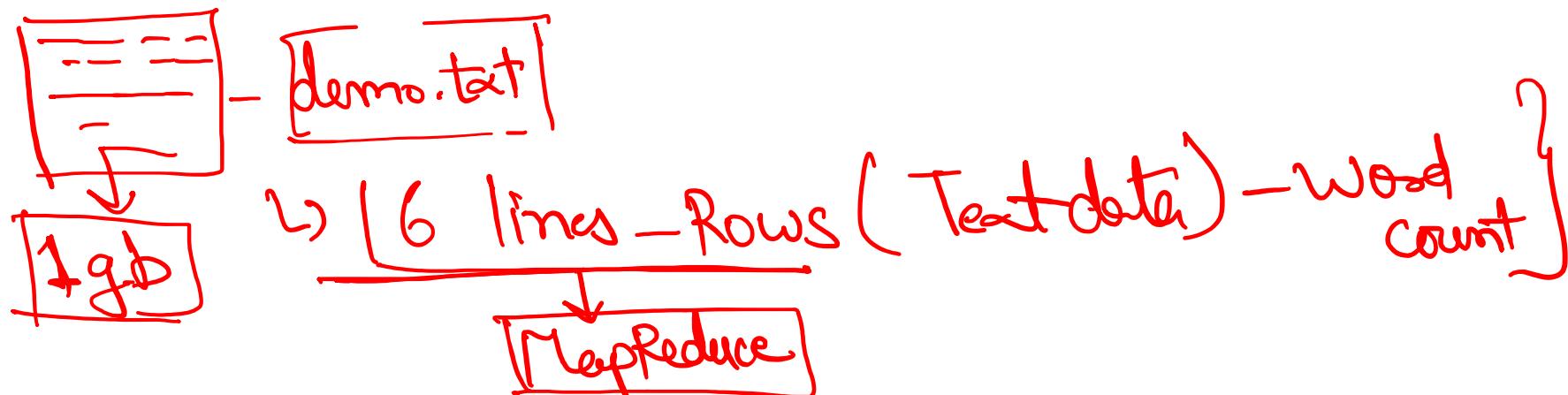
3) HDFS-SITE.XML → Replication of Block

4) MAPRED SITE.XML → MapReduce Program Settings

5) Masters -

6) Slave - Masters Node info a list of hosts, one per line  
IP address of Slave nodes.

MapReduce  $\Rightarrow$  it is responsible for processing the  
data file in hadoop ecosystem



## Map Reduce

- Map Task
- Reduce Task

X BB  
C BA  
X AC

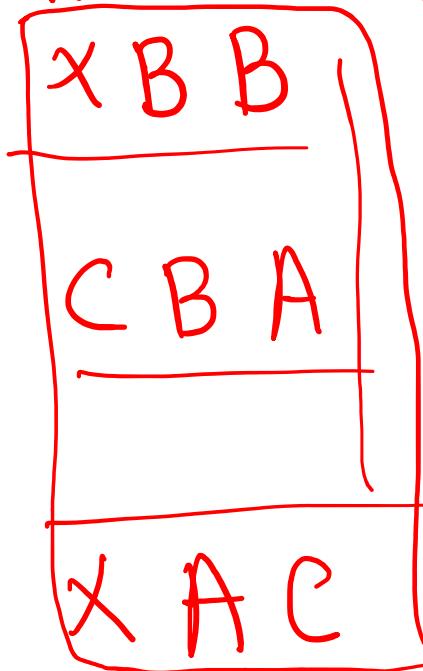
X BB

Example of MapReduce

Program →

Input

file



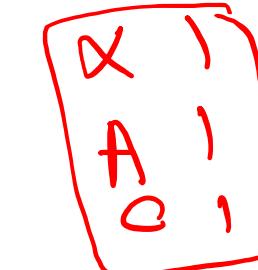
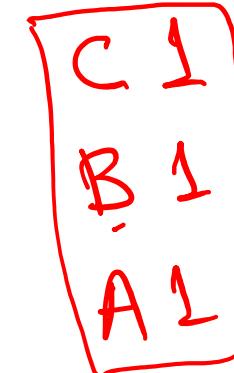
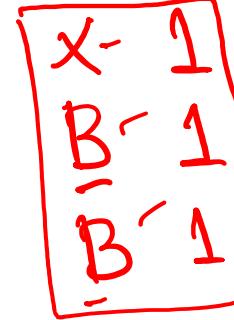
SPLIT

X BB

CBA

X AC

MAP



COMBINE

A 1  
A 1

B 1  
B 1  
B 1

C 1  
C 1

X 1  
X 1

Reduce

A 2  
B 3  
C 2  
X 2

## Map Reduce flow

Input

Map ✓

Shuffle / Sorting ✓

Reduce / Output



client

Job

Map Reducer  
Clusters

Job Parts

Job Parts

Map

Reduce

Output

Javatpoint.com

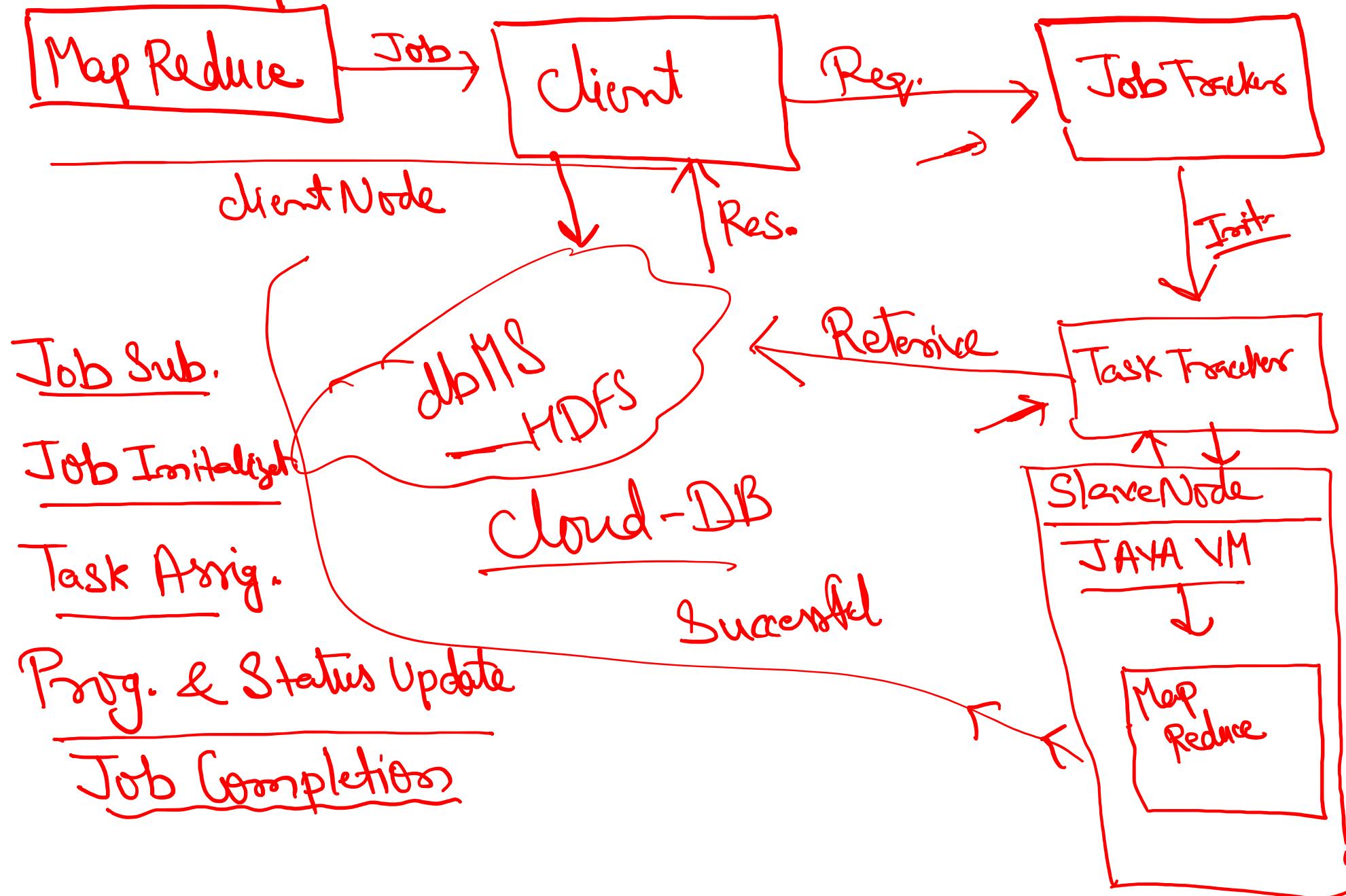
Map Reduce

Weather Dataset

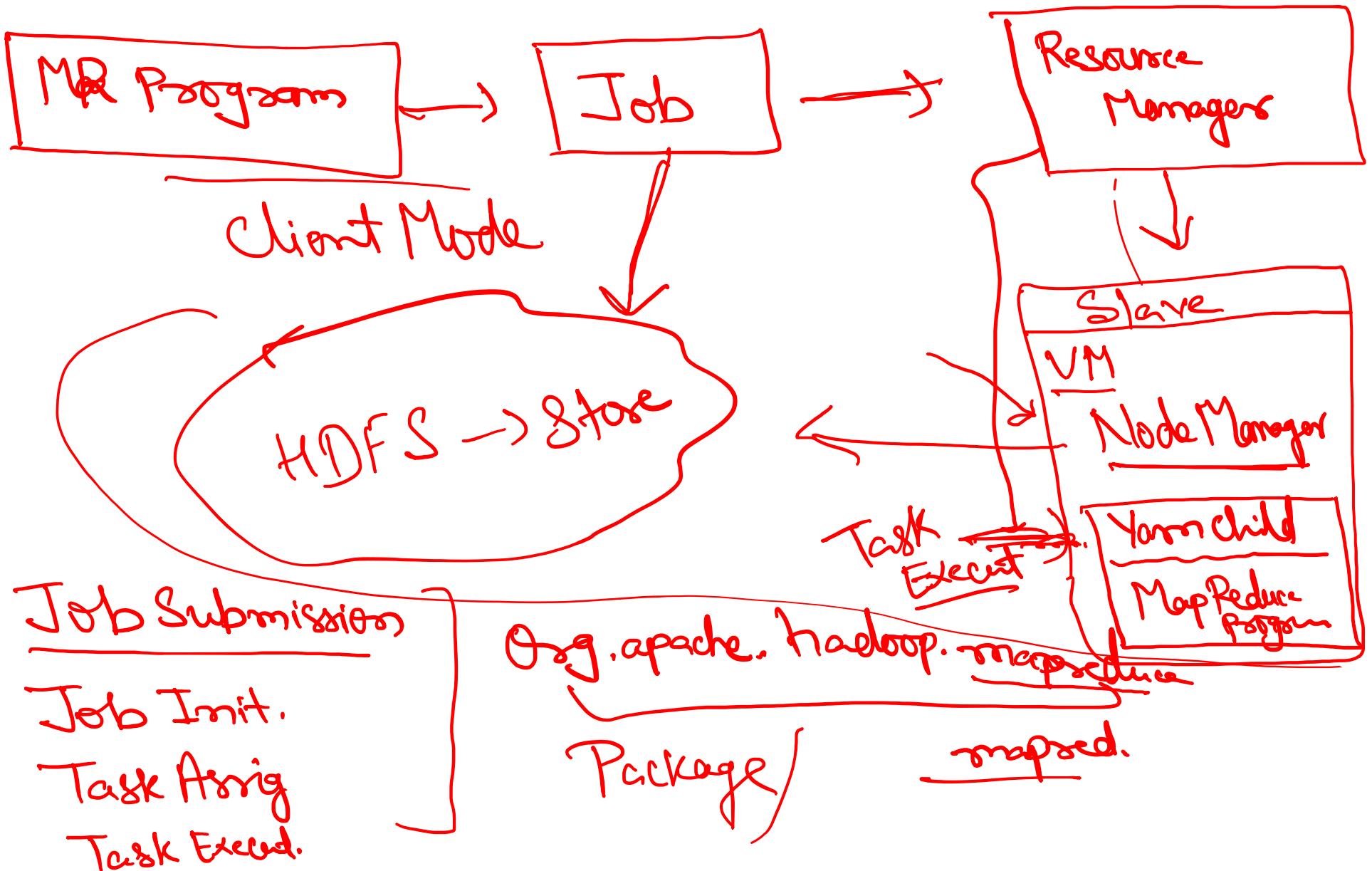
NCDC

Old & New - MapReduce framework

## Old Map Reduce (Classic)



# YARN → (New MapReduce)



Driver Code

Mappers Code

Reducer Code

Recorders Code

Combiner Code

