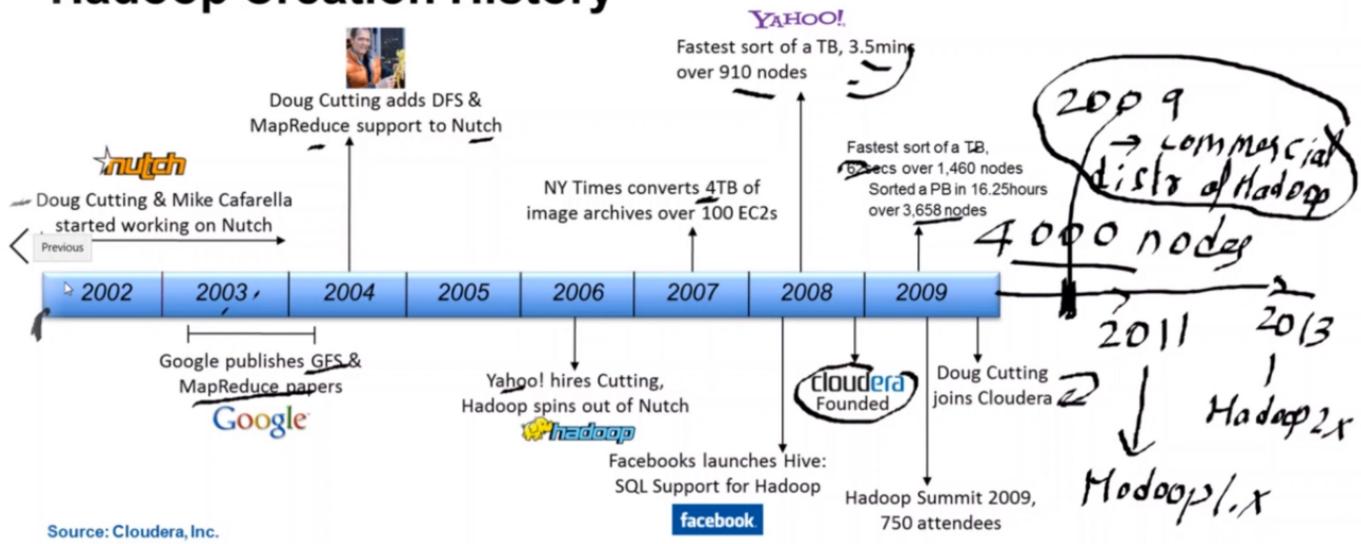


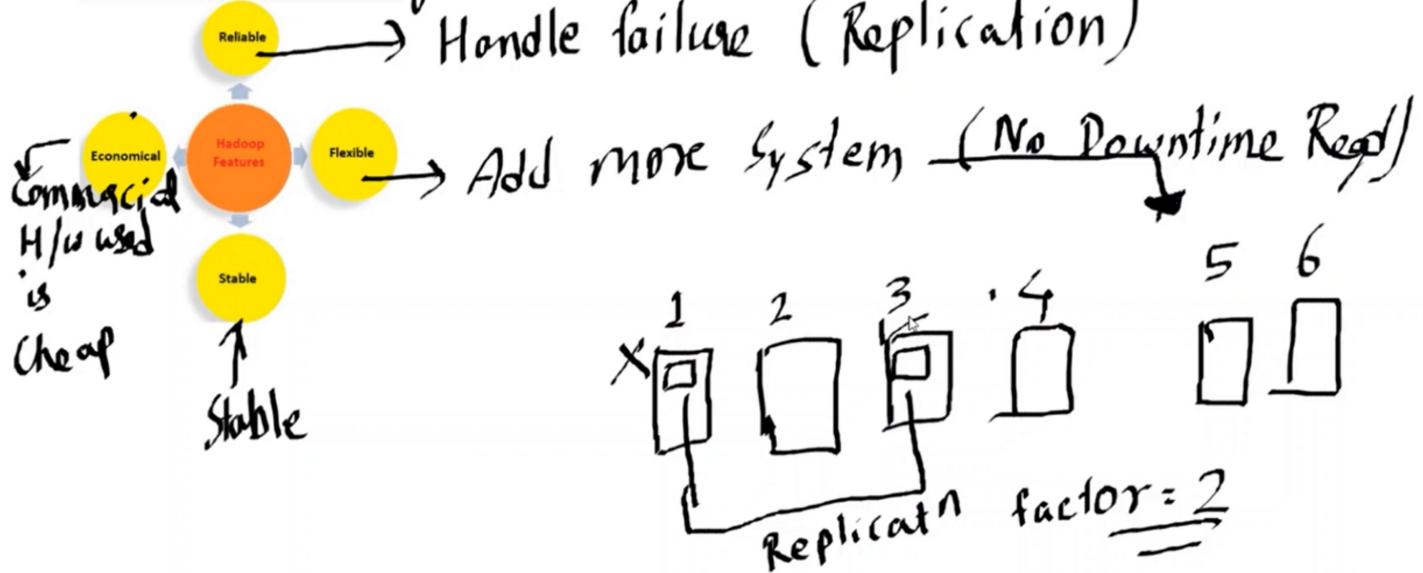
Hadoop

Sunday, October 17, 2021 3:56 PM

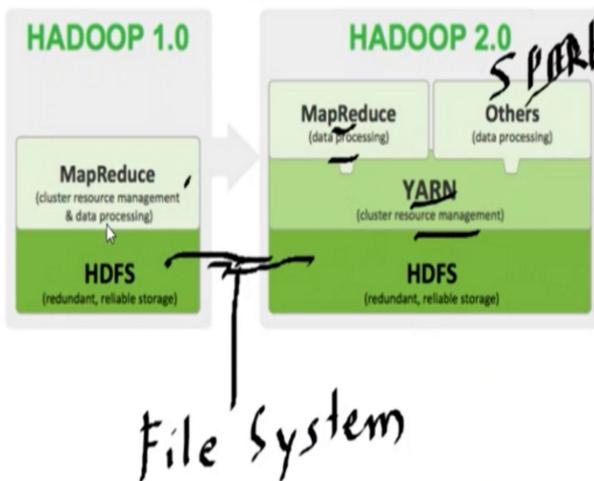
Hadoop Creation History



1.2 → Features of Hadoop



1.3 Major Components of Hadoop MapReduce



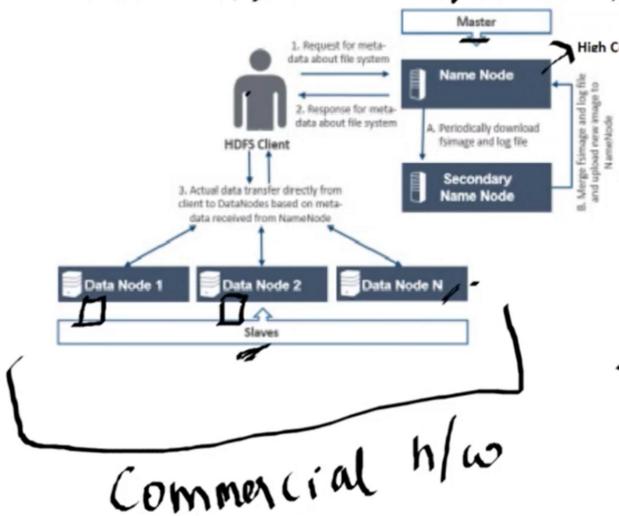
SPIKE

- ① Map → Distributes the Queries
- Reduce → Gathers the Result
- ② Resource Management

2. HDFS

- 2.1 → what is HDFS?
- Hadoop Distributed File System
- Core Component
- like Regular file system (Txt, excel, image, etc)
- TB & PB

2.2 - HDFS Architecture



1. Master - Share arch

2. Name Node

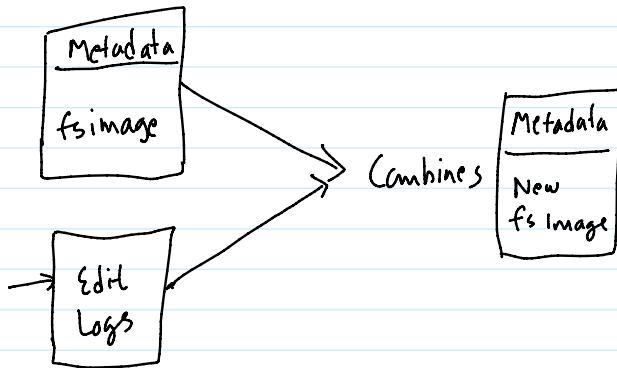
→ Manage file system

→ 300MB

Blocks: 128MB + 128MB + 44MB

→ Daemon Runs on Name Node & Data Node Machine

* Secondary NameNode



Hadoop 1.x

→ No Standby NameNode

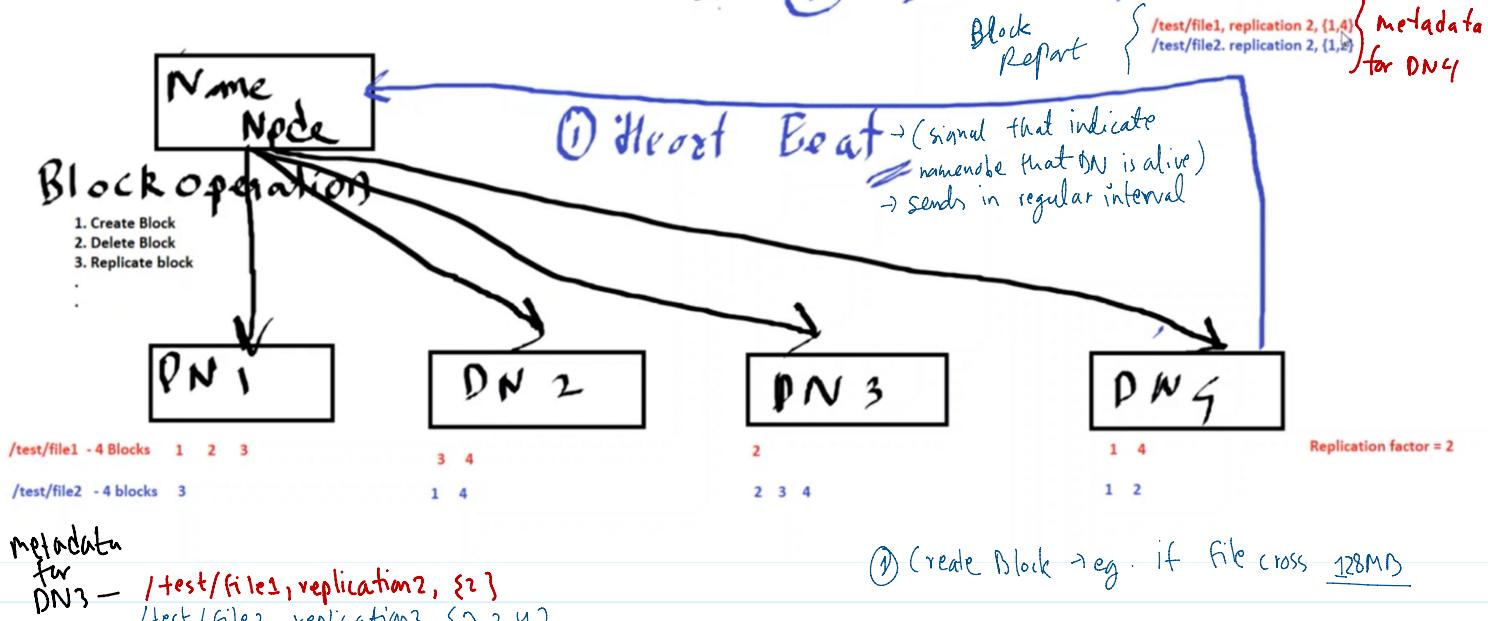
→ Can be called single point of failure
(if namenode fails)

Hadoop 2.x

→ There is secondary namenode

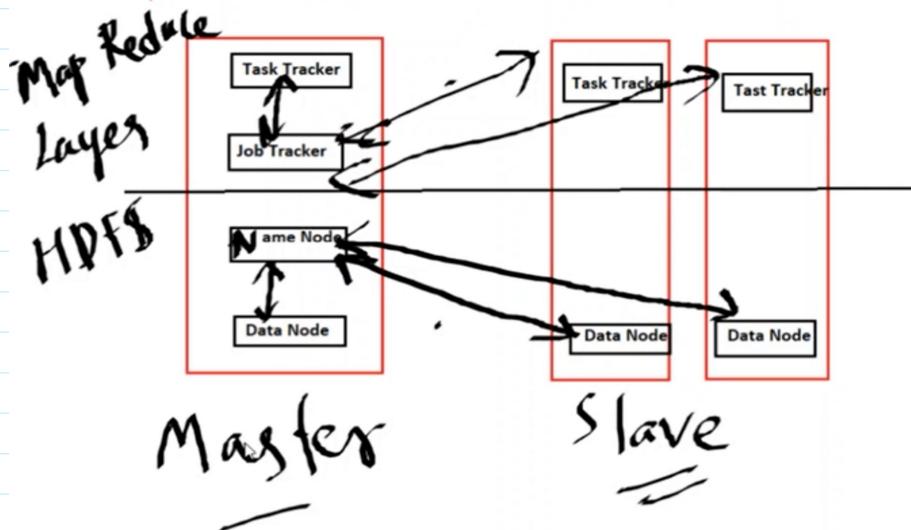
→ Not single point of failure.
(recover by standby namenode)

2.2 - HDFS Architecture - ② Block Report

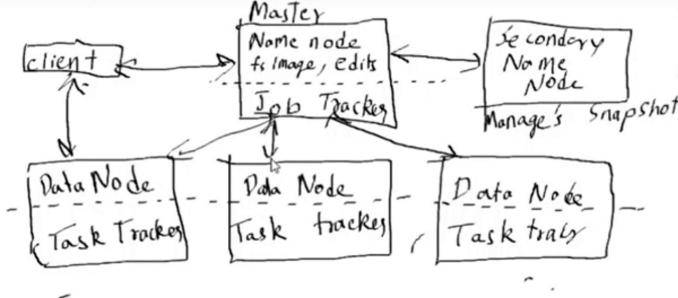


2.2 - HDFS Architecture..

Typical Hadoop Cluster



2.2 - HDFS Architecture



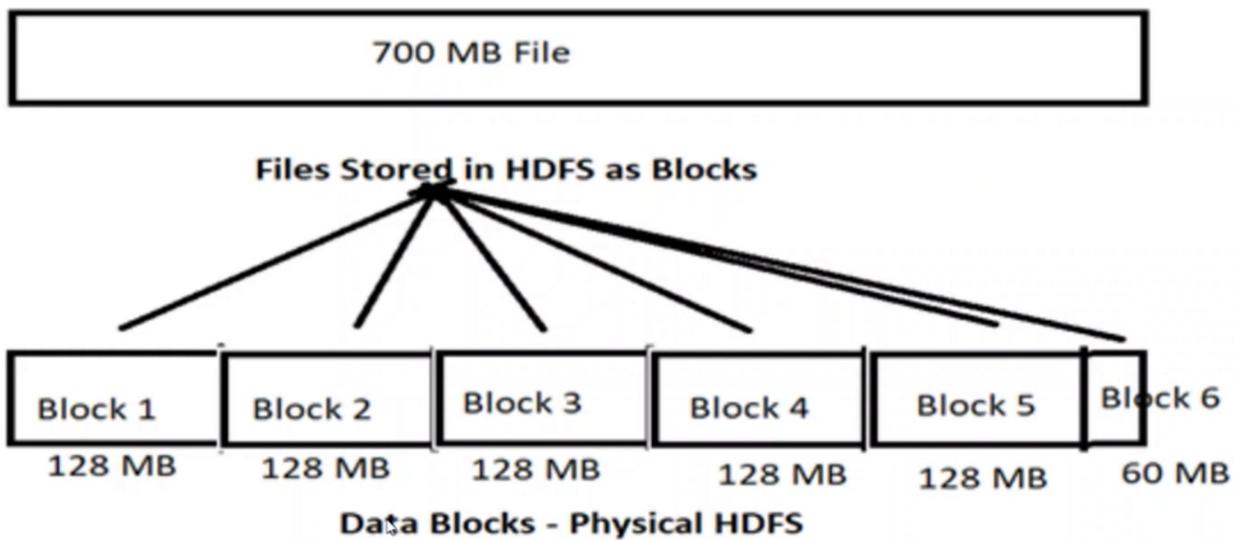
Job tracker → 1 instance of JT in Master

MR job is submitted to JT

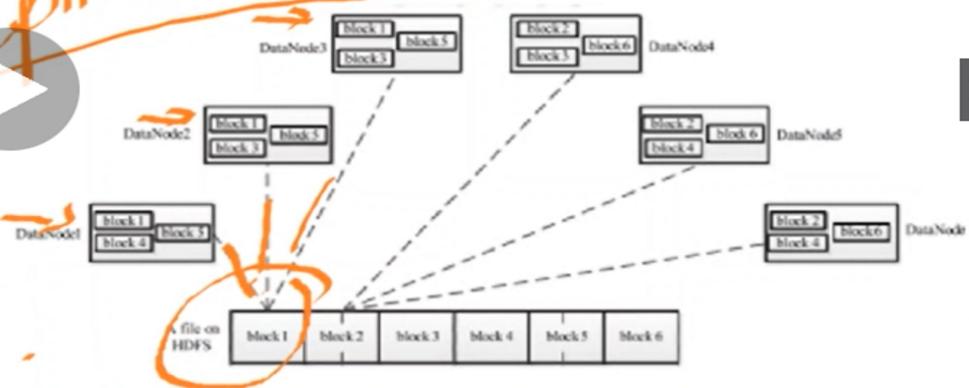
Initiate task in nodes where file is located

Task tracker ← multiple instances of TT in Slave

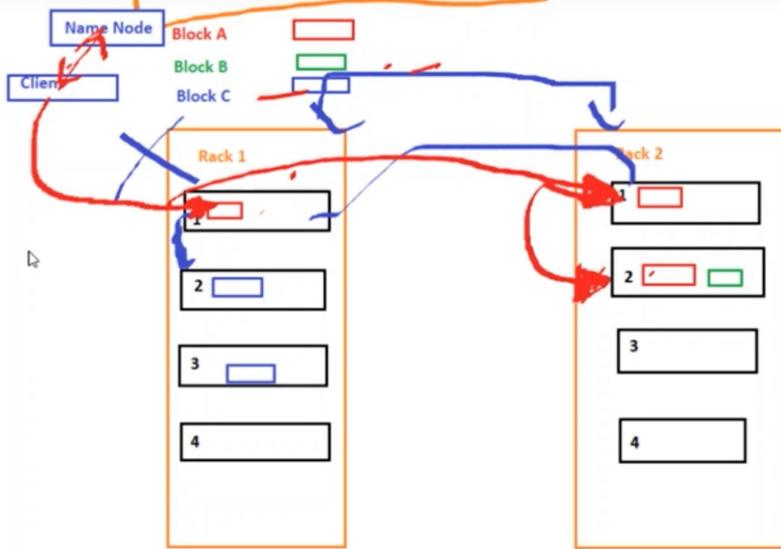
- ↳ Receives information from job tracker.
- ↳ Heart beats to job tracker.



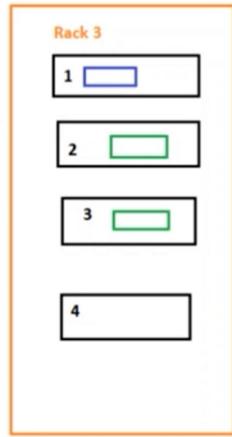
Replicat'n factor = 3



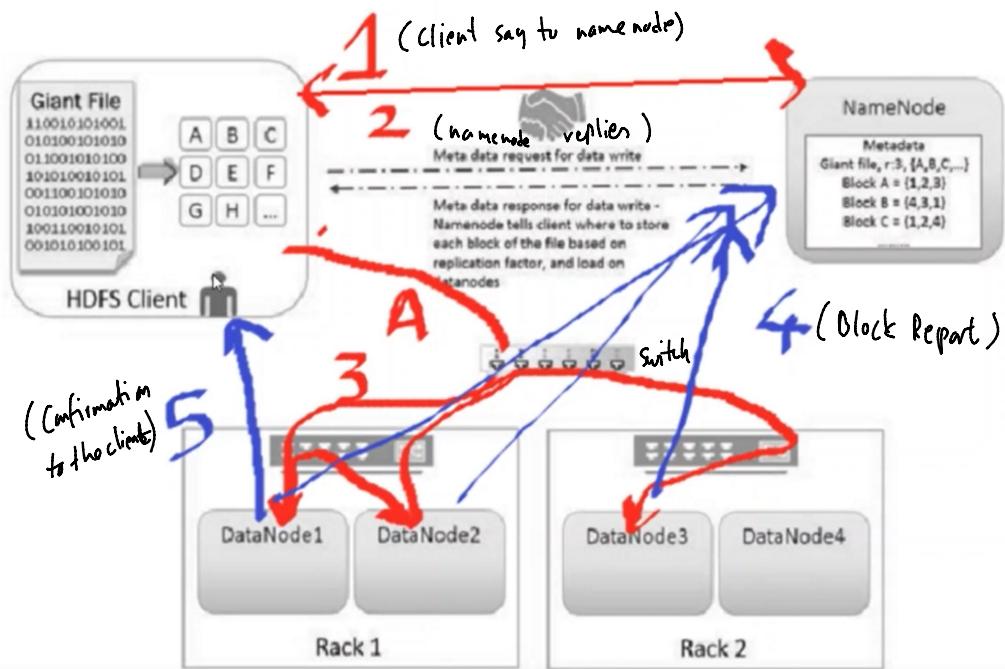
2.4 → Rack awareness



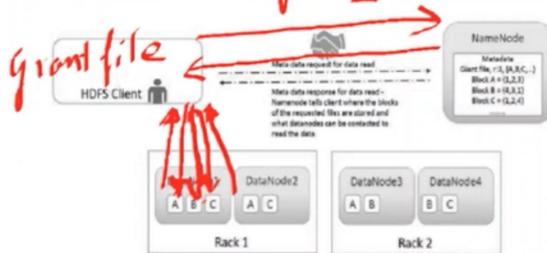
1. First Replica in local Node
2. Second Replica in remote rack
3. Third replica in same remote rack
4. Additional replicas are randomly placed



2.5 Write file



2.6 → Read file



Checksum

Hidden

A + B + C + ... + t = Giant file

calculate checksum

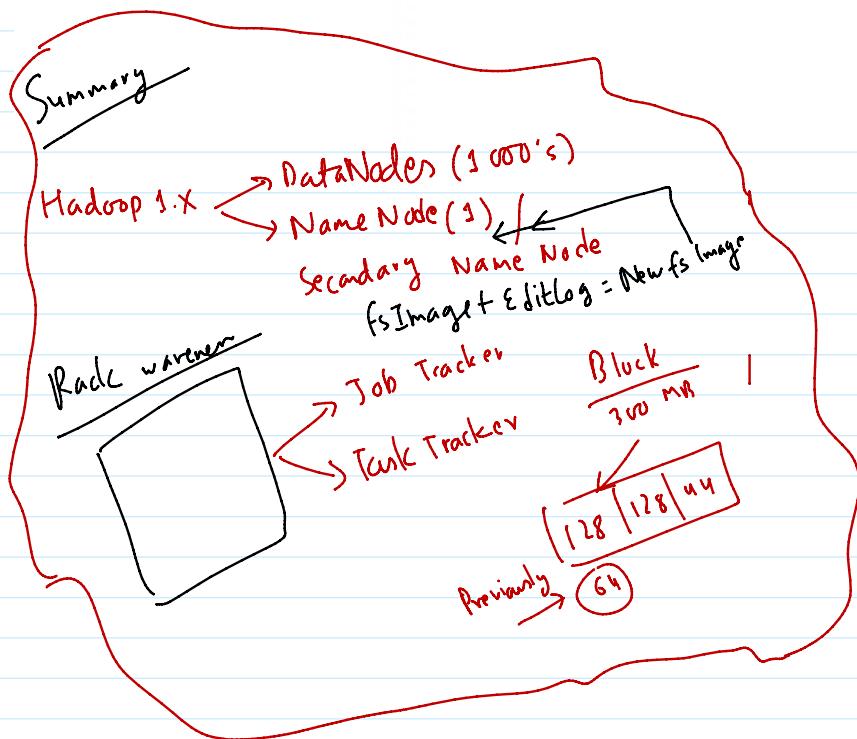
compose

with checksum from Name Node

4016 3601 0005 2007

Data - Checksum

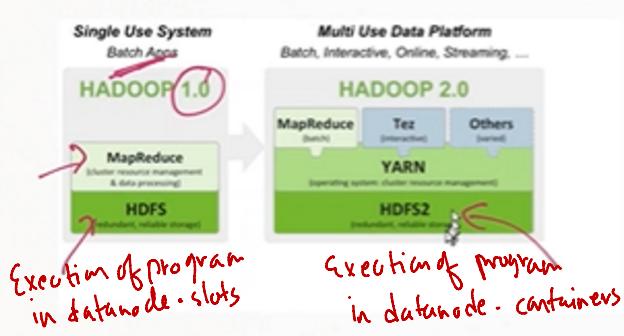
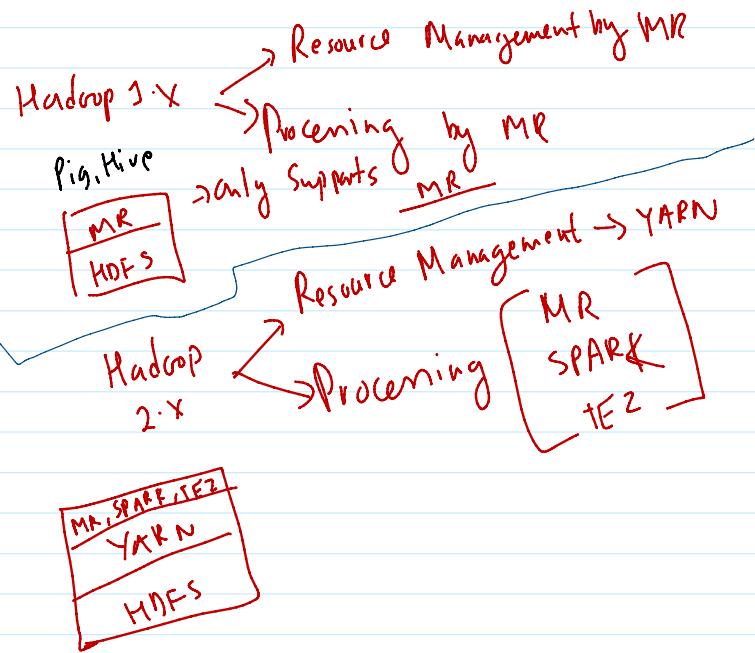
- 3 - Limitation of Hadoop 1.X
- SPOF → Name Node. → Stand By
 - 1 Job tracker → Bottleneck
 - 1000's Task trackers
 - ~ 4000 Nodes → 10,000's
 - (upto) Hadoop 2.X



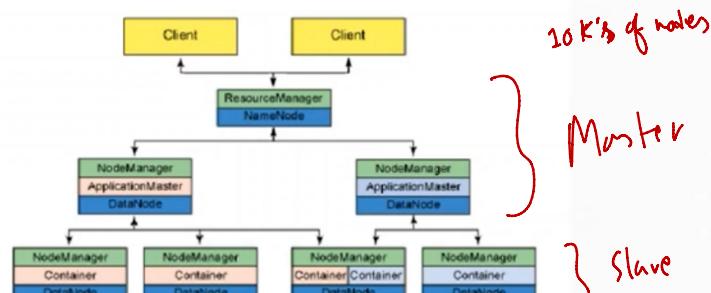
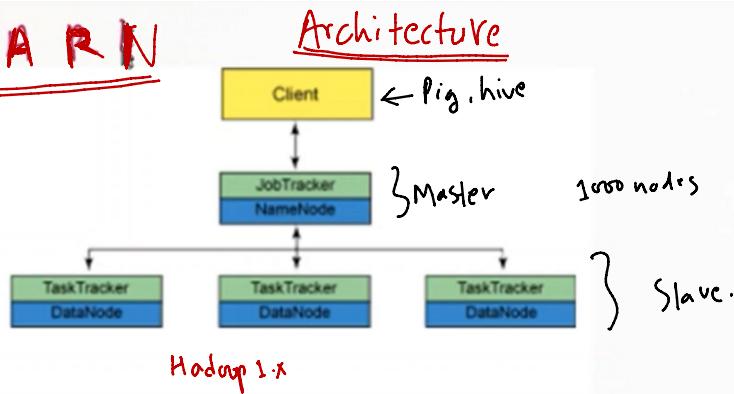
- YARN
- Yet Another Resource Negotiator
- Managing Resource Cluster
 - Introduced from Hadoop 2.0
 - Previously done by MapReduce API

Why YARN?

Hadoop 1.X → Resource Management by MR

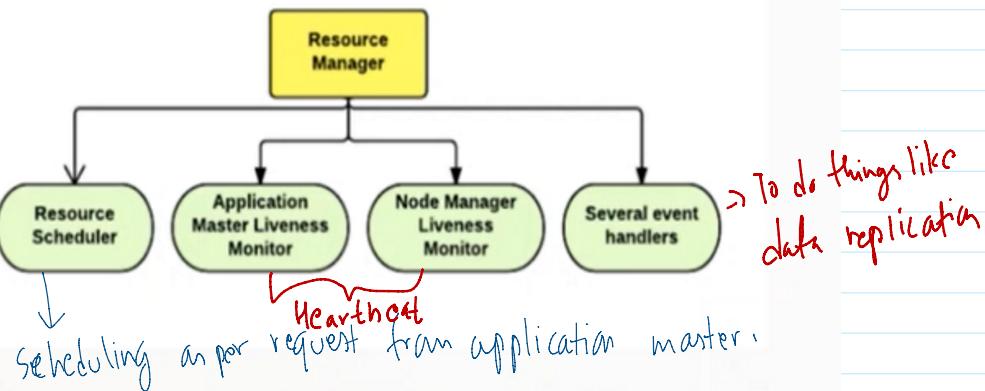


YARN



* There can be multiple slots/Containers inside a single node.

Resource Manager



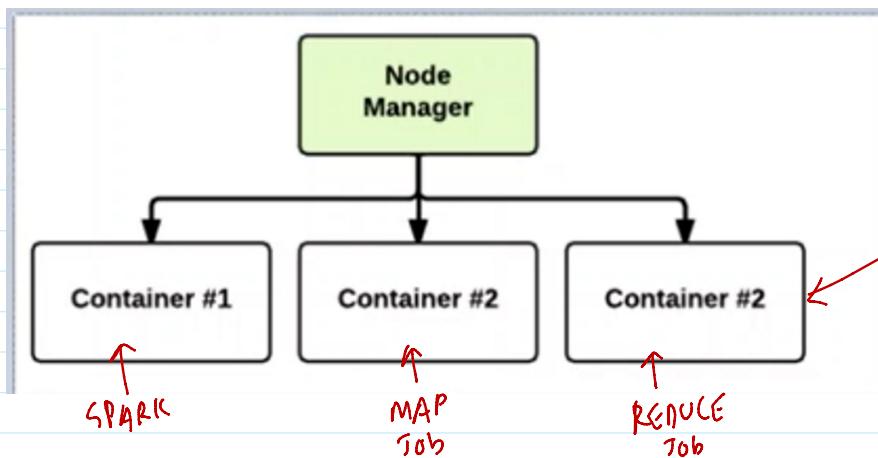
Application Master

- works @ application level
- manages lifecycle of Application & Managers Resources from RM
- load on RM is reduced.

Job Tracker → Resource Manager
Application Master → Application Manager

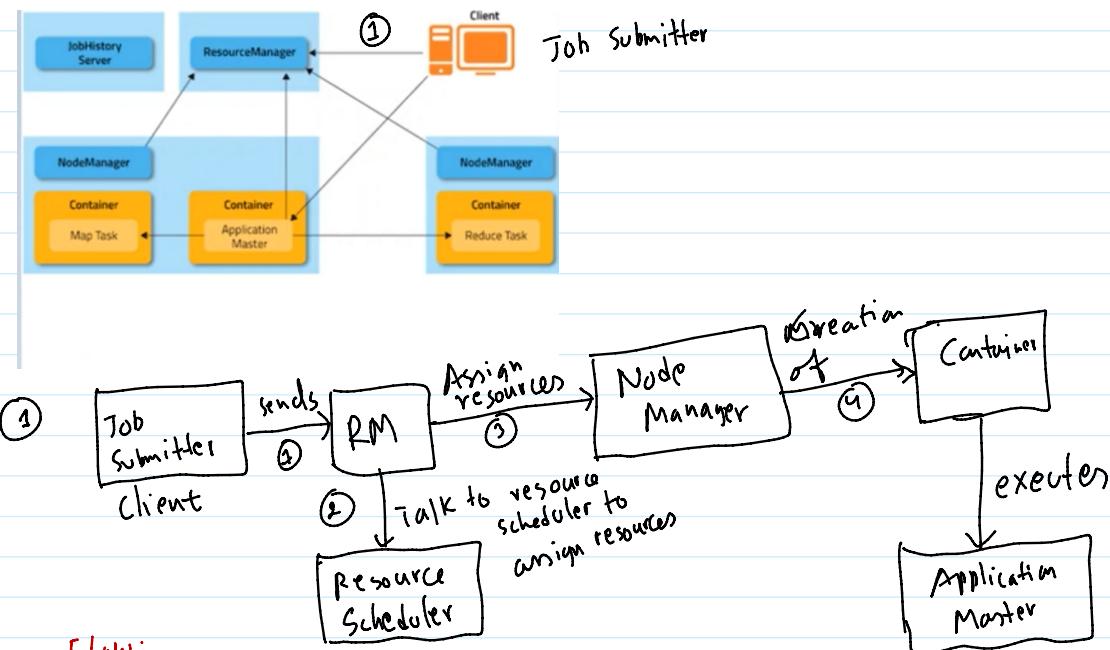
TaskTracker → Node Manager
(basic unit where program is executed)
Slots → Container
1.X → 2.X

Node Manager:



Slot (in Hadoop 1.x)

- 1) Present in slaves.
- 2) Many containers per cluster
- 3) Manage lifecycle of container
- 4) Send heart beat to Resource Manager



Flow:

- ① Client submits application to Resource Manager
- ② Resource Manager (RM) allocates container.
- ③ Resource Manager (RM) contacts Node Manager (NM)

- (1) NodeManager (NM) launches container
- (2) Container executes Application Master

HADOOP 1.X

- 1) MR does processing and management
- 2) Job tracker, Task Tracker
- 3) MapReduce Supported (PIG, HIVE)
- 4) 1000 (upto)
- 5) Namenode (single point of failure)
- 6) No windows support (X)

HADOOP 2.X

- 2) YARN → Resource Management
MR, TEZ, SPARK → Processing
- 2) Resource Manager & ApplicationMaster, Node Manager
- 3) MR, TEZ, SPARK, GIRAPH
- 4) 10,000 (upto)
- 5) Stand by Name Node.
- 6) Windows support also present

COMPONENTS in 1.X

Name Node

Data Node

Secondary Name Node

Job Tracker

Task Tracker

COMPONENTS in 2.X

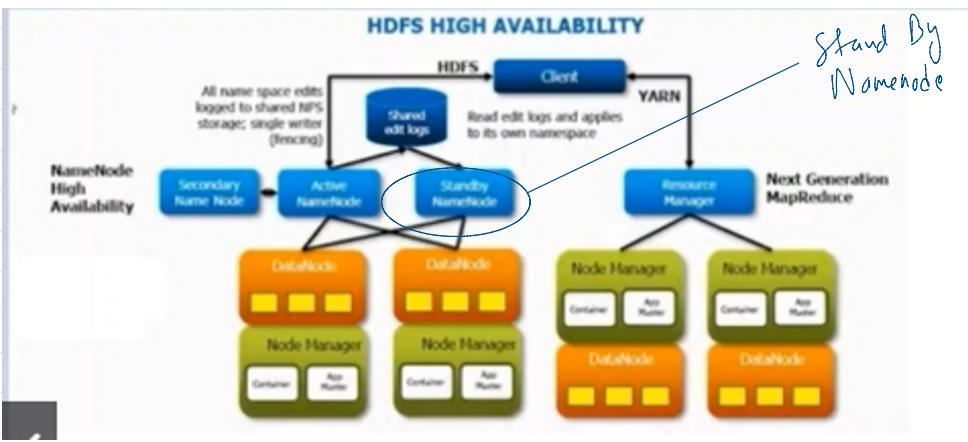
NameNode + Standby NameNode

DataNode

Secondary Name Node

Resource Manager, Application Master

NodeManager



Hadoop Commands

- ① hdfs dfs -mkdir /test1019
- ② hdfs dfs -ls /
- ③ hdfs dfs -chgrp cloudera /test1019
- ④ hdfs dfs -chmod 777 /test1019
- ⑤ hdfs dfs -copyFromLocal /home/cloudera/Desktop/file1 /test1019/file1
- ⑥ hdfs dfs -copyToLocal /test1019/file1 /home/cloudera/Desktop/file2
- ⑦ ls -l /home/cloudera/Desktop
- ⑧ hdfs dfs -count /
- ⑨ hdfs dfs -cp /test1019/file1 /test1019/file2
- ⑩ hdfs dfs -ls /test1019/
- ⑪ hdfs dfs -du /
- ⑫ hdfs dfs -du -s /
- ⑬ hdfs dfs -du -h /
- ⑭ hdfs dfs -mv /test1019/file1 /file2

(16) hdfs dfs -rm /file2

(17) hdfs dfs -tail /test1010 /file2

(18) hdfs dfs

(19) hdfs dfs -help

* hdfs dfs or hdfs fs