

Assignt-2

Date :

Page :

1) Let the use and Responsiblity of following of Hadoop eco-System.

→ a) HDFS →

→ It stores the data across the multiple machines in a fault tolerant manner.

→ HDFS replicates data to ensure fault tolerance & reliability.

→ It is designed to scale horizontally as more data is added to the System.

b) YARN

→ It allocates & Manages cluster Resources

→ It schedules Programs tasks & jobs submitted to clusters

c) Flume

→ Flume collects data from sources like logs, web servers, social media & transfer it to Hadoop storage.

→ It can perform lightweight data transfer during data ingestion.

d) Sqoop

→ Sqoop is used to import data from relational database into Hadoop & export data from Hadoop to Database.

→ It can perform Hadoop to Database.

e) Zookeeper.

→ Helps in coordinating Distributed Programs ensure the have consistent view of Shared Resources.

f) Hive

→ Allows users to write SQL-like queries to extract insights from data stored in Hadoop.

2) Explain Responsibility of Name & Data Nodes & Diff between Normal File System & HDFS

→ Namespace Manager Manages File System Namespace which means it keeps track of file & Directories Structure.

→ Metadata Manager - Stores meta about each file such as file name etc.

→ Block Manager - keeps track of locations & replication.

→ Heartbeat & Health Monitoring

It receives heartbeat signals from DataNodes to ensure they are alive & monitors their health.

Data Nodes

→ Data Storage:- Data Node stores the actual Data Blocks.

→ Block creation & Deletion - DataNodes create,

delete & replicate data blocks is instructed by the Name Node.

→ Heartbeat & Block Report - Periodically send heartbeat to the Name Node

HDFS

Normal File System

* Designed for Distributed Data Storage.

* Data is stored in Single Disk.

* HDFS separates Metadata from Data

* Meta Data & Data are correlated.

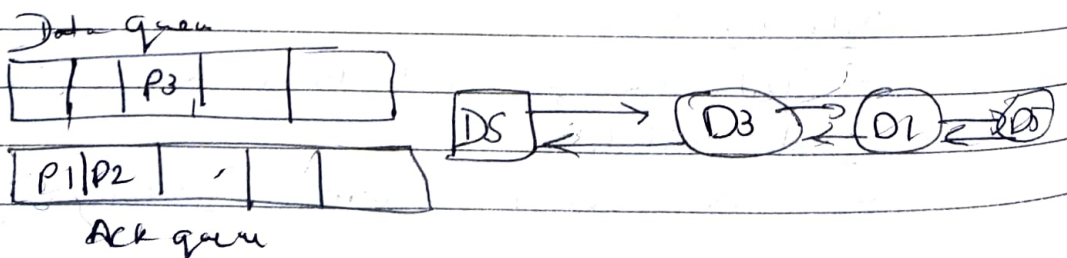
* HDFS replicates Data Blocks to ensure fault tolerance

* Replication & fault tolerance are typically not built in

* Follows a write-once & read many model

* Read-Write can be done many times

3) With Diagram show two write operation is performed in HDFS. Let the steps to recover from failure during HDFS write operation.



→ Let us assume DNL failed while writing P3

→ DS puts packet P3 back to Data Queue & communicate failure of DN1 to NameNode.

→ NameNode creates new id block -002 & add DN1.

4) Explain how HDFS ensure high Availability & Service of Data

→ Data Replication

→ Data Distribution

→ Data Block Placement

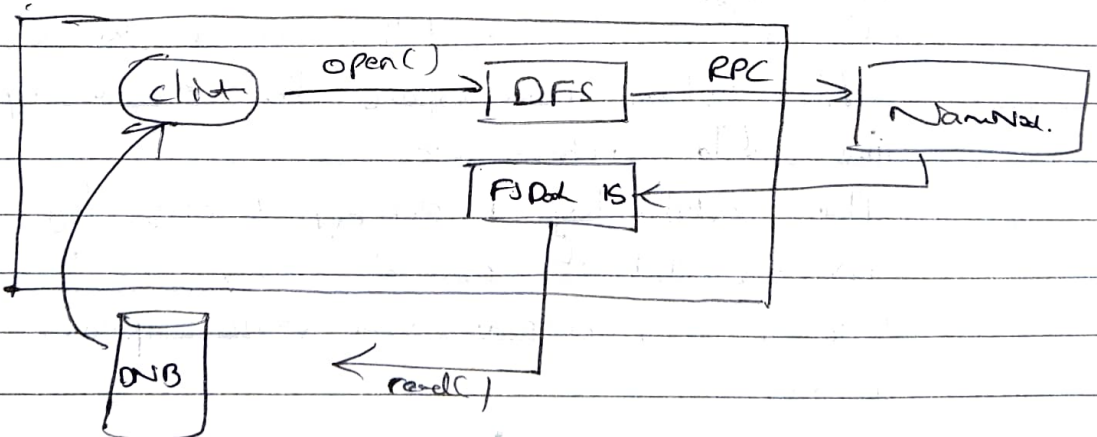
→ Heartbeats & Health Monitoring

→ Secondary Name Node

→ Load Balancing

→ Manual Failover

5) With Diagram Explain, how read operation is performed in HDFS.



* CLNT not a open()

* Through DFS object RPC is made.

- * NameNode search in "metadata" about required file.
- * It find out the Datanode which store the reqd. of data.
- * Use FSDF goes to datanode.
- * If one of the Data node is down, it goes to next nearest Datanode.
- * Reading operations happen sequentially.

Q) Explain the terms heartbeat, checkpoint & edit logs.

→ Heartbeat → It is a periodic signal sent by a datanode to the master server in HDFS cluster. known as the namenode. The signal serves as a way for datanodes to communicate their health.

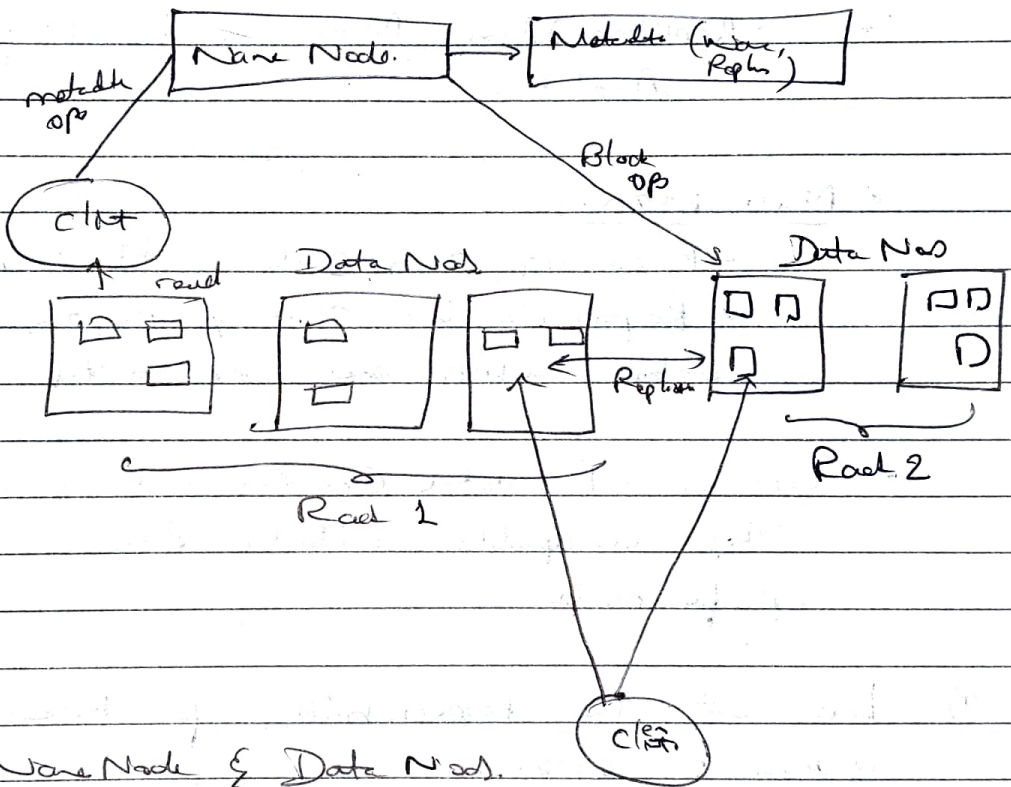
Edit log → It is a log file that records all changes or modifications made to the file system namespace. This includes operations like file creation etc.

The edit log is an append only log means that it records changes as they happen and does not overwrite previous entries.

Checkpoint → This is a mechanism to periodically save a snapshot of the

File System metadata & Data structure

7) With a neat Diagram, Explain architecture of HDFS.



Name Node & Data Node.

→ HDFS is a master-slave architecture.

→ An HDFS cluster consists of a single name node, a master server that manages the file system namespace & regulates access to files by clients.

→ There are a no. of Data nodes which are per rack in the cluster.

→ HDFS exposes a file system namespace allows users data to be stored in files.

→ Data Nodes perform block creation, deletion & replication upon Distribution.

8) List the Major Components of YARN, & Responsibility.

→ * Resource Manager.

It keeps track of available resource.

* Scheduler → allocates Resources.

* Application Master Manages the Lifecycle of Application.

* Node Manager

* Runs on every node in Hadoop.

* Responsible for starting & stopping

Containers, which are the smallest execution

* Container - Fundamental Execution unit in YARN.

* They encapsulate allocated Resources on a specific node.

9) List out the Responsibilities of Resource Manager in YARN.

→ * Resource Allocation

* Application life Cycle Manager

* Scheduler

* NodeManager Communication

* Fair and Priority

* High Availability

* Handles master Application.

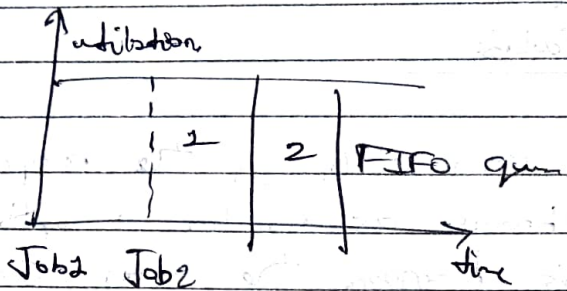
10) Explain Different Schedulers in YARN.

→ FIFO:-

→ Simple Scheduler.

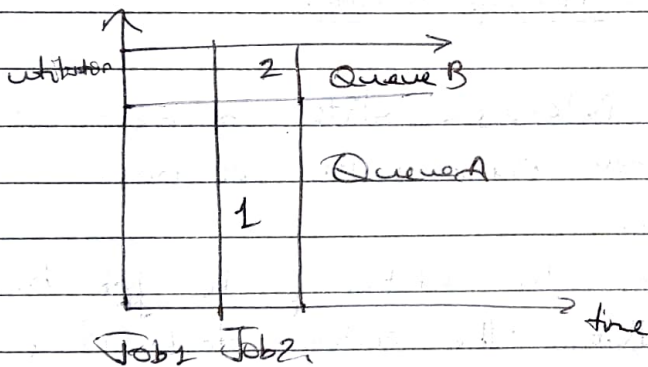
→ Request for the first Application node.

queue as allocated first.



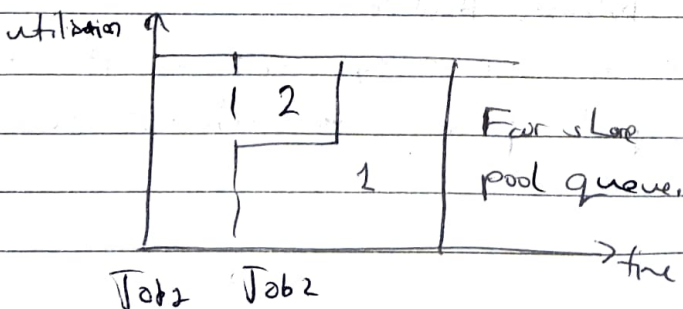
* Capacity Scheduler:-

- * A separate dedicated queue allows the small job to start as soon as it is submitted.
- * It caps overall cluster utilisation since the queue capacity is served for jobs.



Fair Scheduler:-

- No need of reserving a set amount of capacity
- It will dynamically balance resource



11) Explain Different types of failure in Map reduce Job.

→ Task Failure

- When user code in the map or reduce task throws an runtime exception
- If this happens the task JVM reports the error back to its parent application master before it exits.
- This error gets written into the user logs.
- The application master marks the task attempt as failed.

Application master failure

- * Resource Manager starts a new instance of master running in a new container.
- * Recovers the state of the job history so that the entire task need not be run again.
- * Map Reduce client polls the application master for progress. If it doesn't get response the client requests Resource Manager for the new instance address.

Node Manager Failure:

- * Resource Manager recovers it from its pool of nodes to schedule containers

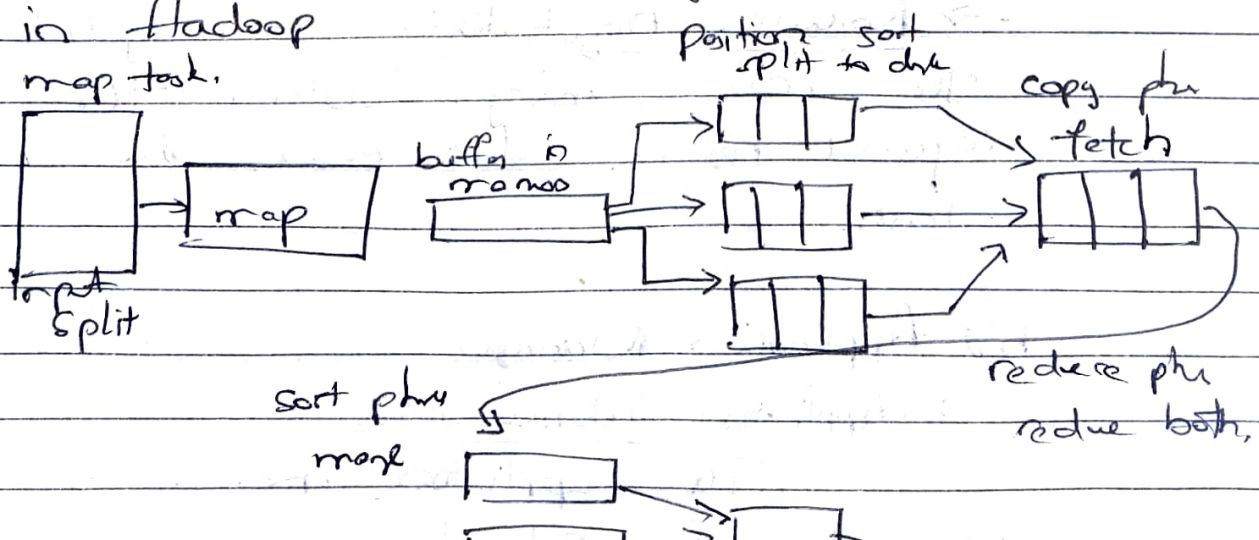
12) Give Different Reason for task failure,
Give steps for recover from application
master failure

- User code in the map or reduce task
Throws a runtime exception
- In this case task JVM reports the error
back to its parent application master before
it ends
- The application master note it is failed
and frees up the container.

Application master failure Recovery:

An application master send periodic
heartbeat to the resource manager. In
the event of application master failure
the resource manager will detect the
failure, the resource manager will detect
the failure and start new instance of
the master running in a new cluster.

13) Show how map reduce jobs are executed
in Hadoop
map task.



Reducer

The map output file is sitting on the local disk of the machine that ran the map task but now it is needed by the machine that is about to run the reduce task for the partition. The reduce task needs the map output for its particular partition from several map tasks across the cluster.

14) Give the Responsibilities of

a) Application Master

→ Application master is responsible for managing the lifecycle of a specific application or job running on the cluster. It is created when an application is submitted and monitors its progress until completion.

→ Resource Negotiation: It negotiates the resource manager to obtain the required resources for executing the tasks.

b) Application Manager→ Application Submission

The application manager is responsible for submitting new applications to the YARN cluster.

c) Node Manager:

→ Node Manager:

Node Manager is responsible for monitor the utilisation of Resources.

d) Resource Manager:

Hadoop map-reduce is a program framework for handle large-scale data process task in a distributed and parallel manner.

It has map phase & Reduce phase.

Example:- Word count.

15) With example, explain Hadoop map reduce process.

→ Hadoop map-reduce is a program framework for handle large-scale data process task in a distributed and parallel manner.

It has map phase & Reduce Phase.

Ex:- Word count

Map-Reduce:-

In the Map phase the input data is divided into chunks and each chunk is processed.

Input data:-

→ Hello world, how are you? Hello world?

→ Mapper task tokenize the input into words & emit pairs for each word.

(Hello, 1)

(world, 1)

(how, 1)

(are, 1)

(you, 1)

(Hello, 1)

(World, 1)

Each word is a key and the value associated with it is 1.

Shuffle Sort:-

Hadoop framework sorts & shuffles intermediate key value pairs group them by key ensure that all key value pairs with same keys are grouped.

Reduce Phase:-

Reducers take groups of sorted & shuffled key value pairs.

output of the Reducers will be the sum of each words.

→ (Hello, 2)

(world, 2)

(how, 1)

(are, 1) (you, 1)