

Big Data Technology notes

(For CCEE and Theory)

Prepared by : B.Dheeraj Chandan
Course : DBDA(Big Data Analytics)
Batch : March-August 2024
PRN number : 240350125017

LinkedIn : <https://www.linkedin.com/in/balivada-dheeraj-chandan-022b91223/>

Big Data Technologies

Big Data - Technology that handles large data sets, data storages.

Skills : Data analysis, Programming(Python, R ,SQL)

Tools : BigData tools(Hadoop, Spark),Data Visualization,ML and Statistics

Sources : IOT, Transactional ,Web, govt..

Big Data adoption -

Drivers (Adoption of high big data)

Barriers (high cost, lack of skilled personnel ,data privacy)

Research and changing nature of data repositories - Trend and Focus

Trends : Advancements in cloud computing , distributed data bases ,real time database systems.

Focus : High data storage, retrieval speed ensuring data integrity and security.

Hadoop versions - V1.0 / V2.0 / V3.0

Version 1.0 : Introduced HDFS and MapReduce

Limitation in terms of scalability and resource management

Version 2.0 : Introduce YARN (),Separates resource management from processing layer.

Enhancing scalability and flexibility
MapReduce + Graph processing , Interactive and Querying

Version 3.0 : Supports erasure coding

Improved storage efficiency,
Addition of docker container ,
Enhancement to HDFS, YARN, MapReduce

Hadoop & Components :-

HDFS - High throughput access to application data.

Splits files into large blocks & nodes, distribute across nodes

YARN:- Schedules jobs and manages clusters resources,
Allows multiple data processing engines to run on Single hadoop cluster

MapReduce:- 2 main tasks Map (Processes & filter data) and Reduce (Pygogates & Sammarises)

Hadoop vs Traditional DB:-

RDBMS handle

DFS (Distributed File Systems) :-

Data stored across multiple Servers & locations. High scalability.
Availability and reliability.
Distributed environment enabling efficient data processing using tools like MapReduce, YARN & Hadoop Components.

Components of HDFS: Name Node & Data Node

Namenode : Master Node manages file system namespace,
Regulates access to files
Datanode : Worker/Slave nodes that stores & retrieves block when they are told. Report back to Name node with block that they are storing. Periodically.

HDFS daemons:-

Name Node - Manages metadata and file system namespace
It keeps on tracking.
Data Node - Daemons that manage storage on individual nodes.
Handles read and write requests from a client.
Secondary Namenode - Merging namespace image & edit logs to reduce load on Name Node.

HDFS follows master-slave architecture where master is referred to Name node and Slave referred to Data node.

Scaling & Rebalancing :

Scaling : Horizontally adds data nodes to cluster.
High storage capacity.
Enhances parallel processing capabilities.

Rebalancing : Distribution of data Node becomes imbalance due to data growth or node additions/removals [HDFS balancer & redistributes]

Replications: Replicates data blocks across multiple data nodes

Rack Awareness : It uses a rack awareness algo to improve data readability and bandwidth utilization.

Data Pipeline :- Client writes to Datanode in sequence, 1st data node and lastdata node

Nade failure Management : It handles node failures automatically. If data node fails the name node automatically initiates the replication of the block.

Zookeeper- ZooKeeper is used to manage and coordinate distributed processes, maintain configuration and state, and handle tasks like leader election and failure recovery.

HDFS architecture : Manages & Process data sets

Design principles: Replications, rack awarness, High availability

Features can build big data solutions.

Hadoop Operation Mode:-

Standalone Local Mode (HDFS on a single Machine)

Pseudo distributed (HDFS runs on single machine simulated distributed environment)

Fully Distributed (Runs on multiple machines separate nodes: Namenode, DataNode, Resouce node and Secondary Name Node)HDFS and YARN daemons [Prodution environment]

Core-site.xml : Configuration Hadoop proposties

HDFS : Configuration for Data Node Nane Node

MapRed : MapReduce job ; history (location and framework)

Yarn : Configuration for YARN ; Resource Manager & Node Manager

Resource Manager : Manages resources & schedules jobs

Node manager : Manages individual nodes resources & monitors resource usage

Data Node : Stores actual data block

Name Node : Manages file system namespace and metadata

Secondary NameNode : Periodically merges mamespace image & edit logs.

Job Tracker(MRVI) : Manages job scheduling & task trauking.

Résource Manager: handles resource management

Application Master: Takes care of Job scheduling & task Monitoring

Task Tracker(MRVI)→ Executes individual tasks assigned by job tracker.

Resource Manger (YARN) : Manages resources & application

Node Manager (YARN) : Manages resource & containers on a single node

Application Manager (YARN) : Manages life cycle of a single application

DataFlow : Client requests a name node and then name node accepts and gives address of data node so that client can directly access Data node directly for read and write operation.

Scaling : HDFS is horzonatal scaling, High storage capacity.

More Data Nodes are added to cluster for high storage capacity.

Horizontal scalability: Adding date nodes to cluster & handle

Load Balancing :

NN federation :

Rebalancing in HDFS-

Ensures even distribution of data across Data Node.

HDFS balancer Tool: To redistribute accross DataNode via cmd

Threshold Based Balancing & Threshold parameter to determine When rebalancing is necessary.

Incremental Balancing: Rebdancing is performed incremenatally to minimize disruption to the system.

Optimal performance & resource utilization HDFS a robust solution

MapReduce for parallel processing

HDFS commands:-

Data split into blocks & Blokes are replicated 3 ,128 MB each. Data Nodel & ensures fault, Tolerance & high availability.

hadoop fs -chmod /path - To change permissions

hadoop dfsadmin -safemode - To enter into safe mode

hadoop dfsadmin -safenode - To leave

hadoop dfsadmin - To refresh the list of Data Nodes

NameNode : 9870 and ResourceMaanager : 8088 port numbers.

Hadoop MapReduce : Processing large datasets in a parallel across a datasets. Map - Process input data and generates key value pairs. Reduce - It aggregatesintermediate key-value pairs to produce final output.

Map Reduce Daemons :

Resource manager - manages resources

Node Manager - Monitors resources

Application Master - Manages execution of a single map reduce

Anatomy of a MapReduce Job Run

Job Submission - Client submit job to Resource Manager

Input splits - Splits into logical chunks possessed by separate Map tasks

Map phase - Proceses input & generates key-value pairs

Shuffle and sort - Key value pairs shuffled and sorted by keys.

Reduce Phase - Process grouped key-value pairsto generate final output.

Output - final output written to HDFS or another storge System

Postioners & Combiners

Partitioner - It determines which reducers a particular key-value pair will go to
(to control data distribution)

Combiner - Performs local aggregation of intermediate key value pairs that they sent to reducers.

Input Formats

Input splits & Records - Logical chunks to Map tasks. Generates key-value pairs.

Text input - Default input format . Each line treated as a record

Binary input - Sequence file input format to read binary key value pairs stored in hadoop sequence file

Multiple input - Allows input for different data sources in the same job

Output Formats :-

Text output - Key-value pairs written as text lines default output format

Binary Output - Sequence output format (Writing binaries key value pairs to hadoop sequence files).

Multiple output - Allows multiple data to multiple output from a single job

Distributed Cache -

Distributes read only data needed by MapReduce jobs to all nodes in the cluster.

add CacheFile(): Adds a file to a distributed cache

add Cache Archive() : Adds an archive to the distilate cache

Hadoop MapReduce paradigm - large scale data processing by dividing tasks into smaller subtasks .Distributed Cache performed by making read only data readily available to all nodes.

Hadoop DataTypes:

[6 Primitive] - Int, Long, Float, Double, Boolean, Text

[3 Compositée] - Array, Map ,Null

Java & MapReduce

Mapreduce written in Java.

Robust API for implementing Mapper, Reducer

Map Reduce components:

- Mapper : Define map function
- Reducer : Define reduce function,
- Drive class: configures job & specifies input output paths, Mapper & reducer etc

Map Only - No reduction or aggregate is needed. Reducer is omitted or set to null Reduce only

Reducer only - Preprocessing has already done only aggregation is needed

Combiner & Partitioner :

Combiner - Acts as mini Reducer to reduce data transfer

Partitioner - Controls the partitioning of intermediate map outputs

job.setCombinerClass(),
Job.setPartitioningClass()

Counters - Tracks the progress or count the events by .getCounter()

Schedulers :- Its task for job scheduling

FIFO scheduler - Jobs executed in order they arrived

Fair scheduler - Jobs are allocated resources that all jobs get an equal share of resource over time

Capacity scheduler - To run a Hadoop jobs in a shared ,mult tenant cluster to ensure certain capacity guarantee.

Custom Writables : To define custom data types like
Writable
Writable Comparable

Compression : Reduces amount of data transferred & stored by following methods
.setCompressOutput()
.setCompressorClass()

Complex Map Reduce Programming

Multiple MapReduce jobs can be chained together output of one job becomes input for next which is also called chaining job.

Different stages of processing having multiple mappers and reducers

Map reduce Streaming: Writing MapReduce jobs in any language that can read from standard input and standard output

Python & Map Reduce :- Python scripts for both Mapper & Reducer

MapReduce on image dataset

Image processing (resizing, format Conversion, feature Extraction) using Python Libraries : PIL or OpenCV

Implementation:

Mapper - Process each image ,performs necessary image processing tasks

Reducer - Aggregates results if needed.

Hadoop ETL

ETL (Extract Transform Load): Extracting data from various sources, Transform it into suitable format similar storage system. Loading into Hadoop's HDFS

Extraction: Extract/Load data from files

Transform: Clean up data or arrange data to standard or specific

Load : Transformed data into HDFS, HBASE or Hadoop compatible storage

ETL functions :-

 Data Extraction:

 Sqoop (data between Hadoop & RDBMS)

 Flume (collects, aggregates-Streaming data into HDFS)

 Kafka (data streams processed in real time and loaded in HDFS)

 Data Transformation:

 pig (Pig Latin to handle complex data)

Hive (SQL like query for data transform & analysis)
Spark (API for complex data transforms using RDDs and Dataframe)

Data Loading:

HDFS (Storage for large datasets)
HBase (NoSQL DB runs on top of HDFS, read and write access)
Hive (Acts as dataware house storing data in a structured format)

Need for ETL Tools :- Complex data Management, Automation, Scalability, Integration

Advantages of ETL :- Efficiency, data Quality, Compliance, Flexibility ,Performance, User Friendly

HBase-

Open Source, non-relational, distributed database.
Written in Java
Handles data across commodity servers.
Provides real time read/write access to big data
Features : column oriented storage, Scalability, Consistency, Integration with Hadoop, High availability , Horizontal scaling

Architecture :

Components -

HBase Master - Manages Cluster (Table creation, deletion & balancing regions)
Region servers - Handles read/write requests.
Manages HFiles, & WAL (Write ahead log):
Zoo Keeper - Coordinates distributed process &
Provides Configuration &
Synchronization services.

Data Model : Tables, Rows, columns, column families ,cells.

Storage : HFiles, WAL (Write-Ahead Log)

Region - Subset of table horizontally partitioned row

key range - Region Split fically splits it Automatically

Apache Hive -

Data warehousing and SQL-like query language tool built on top of Hadoop.
Allows users to query, summarize, and analyze large datasets stored in Hadoop's HDFS (Hadoop Distributed File System).
Converts SQL-like queries into MapReduce jobs that run on Hadoop, making it easier for users to work with big data without needing to write complex MapReduce code.

Key Features:

- SQL-like Interface: Uses HiveQL, a query language similar to SQL.
- Schema on Read: Data is interpreted during query execution, not when stored.
- Supports Large Datasets: Designed to handle and process petabytes of data.
- Extensibility: Supports user-defined functions (UDFs) for custom processing.
- Integration: Works with HDFS, HBase, and other Hadoop ecosystem tools.

Use Cases:

- Data Warehousing: For storing and querying large volumes of structured data.
- Data Analysis: For performing data summarization, querying, and analysis on big data.

Hive is widely used in big data analytics for its simplicity, scalability, and ability to integrate seamlessly with the Hadoop ecosystem.

Big Data Technology - Notes
By - B. Dheeraj Chandan
(240350125(017))

3 | 7 | 24

- Big data helps to analyze, systematically extract information deals with large or complex data by data processing applications.
- Data can be unstructured, semi-structured, Quasi-structured & Structured data.

Importance :- cheaper storage, cloud based technologies, improvements in data processing techniques & analytics.

Penetration of Technology - IoT, smart devices.

Ever increasing data, ↑ of unstructured data.

Unstructured data : Text document, PDF, Img, videos

Quasi-structured data : clickstream

Semi Structured data : Spreadsheet XML B.D.~~Big Data~~

Fully Structured data : Database

Industries :- Reuse, Aerospace, Sports, Retail --

Characteristics of big data :- Volume, Variety, Velocity, Veracity

- Volume → Quantity of generated & stored data.

- Variety → Type & nature of data

- Velocity → speed of data generation & processing time

- Veracity → Defn. of big data which refers to data quantity & the data value.

Big data ∈ [Data size, type of data, speed of data processing, data quality & data value]

Challenges :

2024/8/8 18:40

→ Scalability, timeline & cost challenge before making a move to big data.

→ Not all data is useful (data contains more noise)

→ Lack of technical talent

→ Data Security cornerstone of all infrastructures.

→ lot of organisation and human resistance when it comes to making a transition of big data.

Scalability (able to handle the big/large amount of data in the system) Data usability should be reusable.

Lack of technical talent: Automated management, Organisation is investing in ramping up technology skills of workforce.

Data Security : Security systems need to be used by Hadoop admins - stakeholders.

Platforms : Cloudera, AWS, Google, Azure, Digital Ocean, MapR, IBM, Intel.

For Hadoop installation we have to do many configurations.

What is Hadoop :-

→ Open source Software Framework for storing large amount of data & running appli. parallelly on clusters of commodity hardware.

- Ability to store large amount of data and type & process it quickly (Computing power - Capable of processing data very fast).

Fault Tolerance, Flexibility, Low cost, Scalability

FT: Data & appli. processing are protected against hardware failure.

Flexibility: Like RDBMS no need to preprocess data before saving. If 1 node fails then we can use as much of data & decide how to use it like redirected later. ex:- It can store unstructured data like text, img, videos -

- Low Cost: Free open source framework, uses commodity hardware to store large amount of data

- Scalability: Handling large amount of data by adding more nodes.

→ Created by Doug Cutting

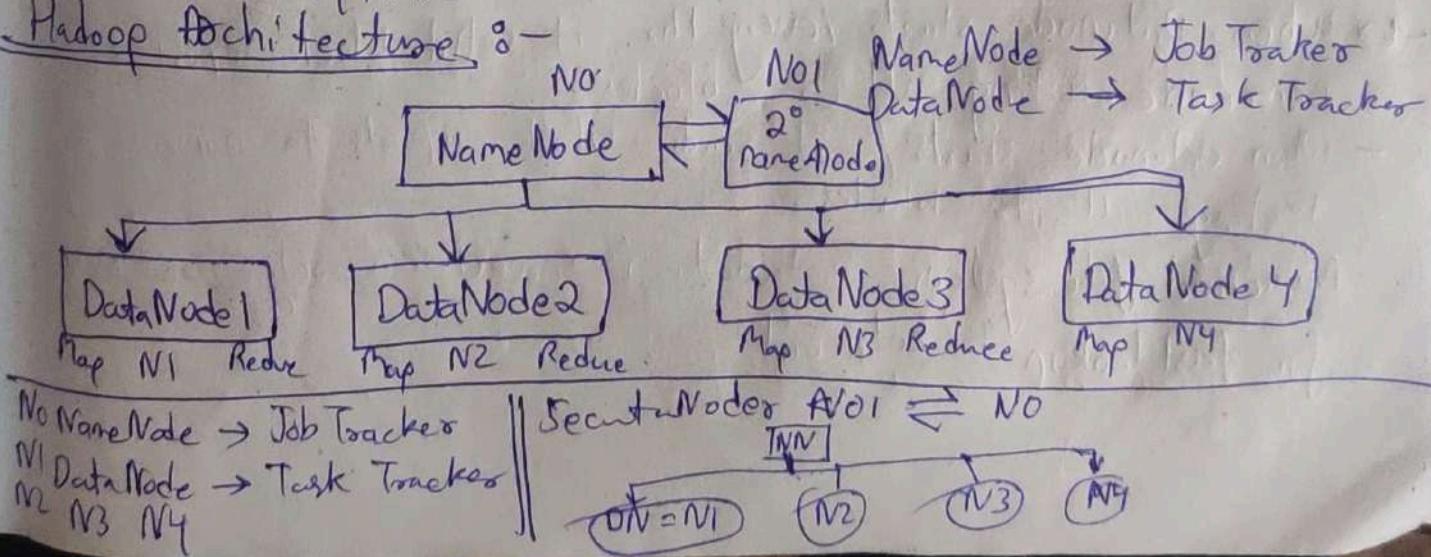
Evolution :- Hadoop evolution → Hadoop 1, Hadoop 2, Hadoop 3

Hadoop 1 - Highly distributed, fault tolerance storage, data local possible

Hadoop 2 - Allowing better fault tolerance, ↑ ecosystem now allows more specialized applications & in memory processing

Hadoop 3 - It has enhancements around fault tolerance, optimization etc..

Hadoop Architecture :-

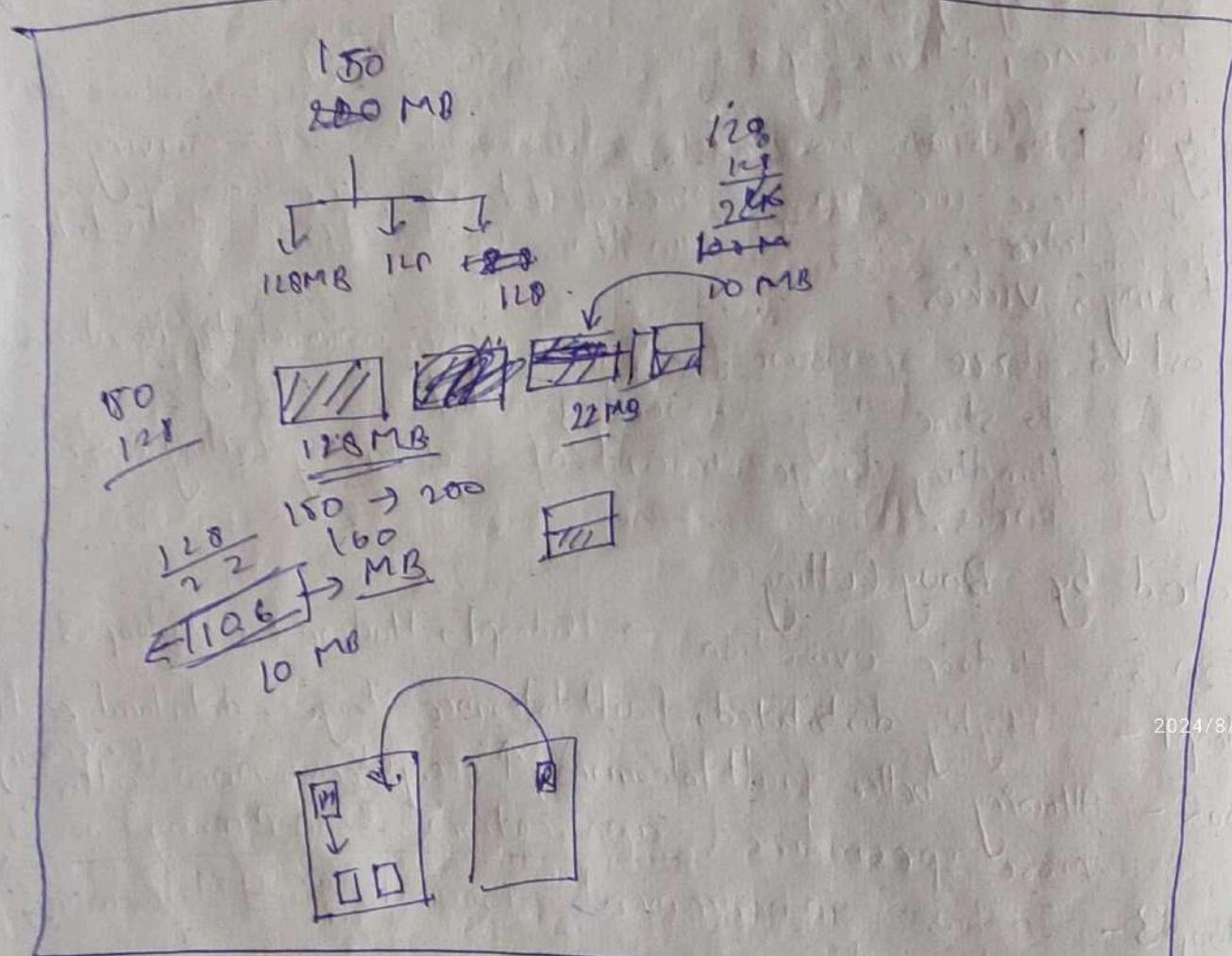


Architecture Components: Namenode, datanode, secondary namenode
Job Tracker, Task Tracker

Namenode - Filesystem + Metadata, info stored in local disk & saved in 2 forms: namespace image & edit log.

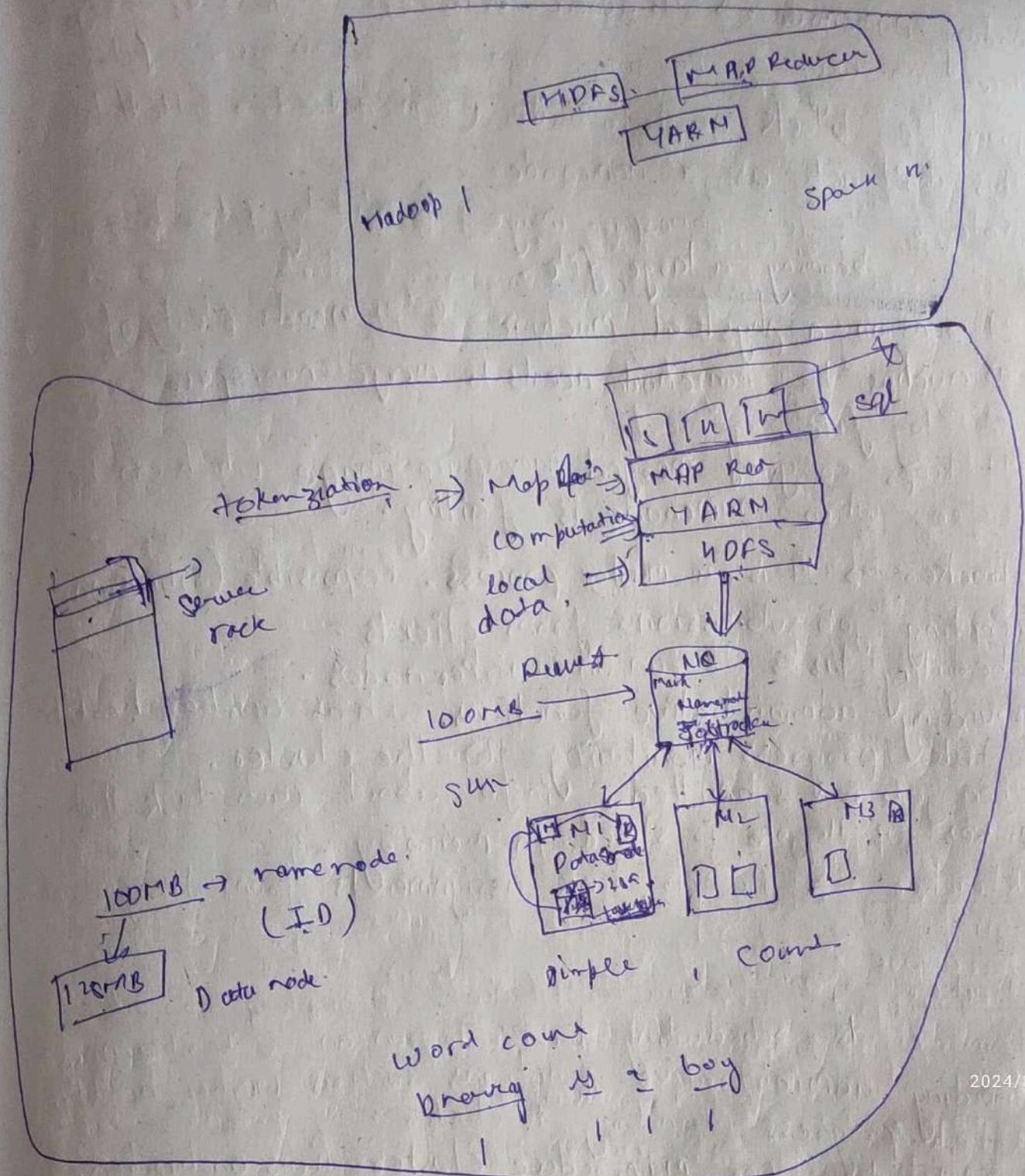
It knows all datanodes. Info is reconstructed from data node when the system starts. Cli

Client connects namenode to perform filesystem operations.
Datanodes always reports their status to the namenode so that Namenode has full view of all datanodes in a cluster.



2024/8/8 18:41

- When Datanode starts / every 1hr it sends block report to namenode, listing all its blocks.
- When datanode starts / it sends block reports to namenode, listing all its blocks.
- Namenode file system metadata is stored entirely from RAM for fast lookup & retrieval



Datanode — Workhorses of file system, store & retrieve data from block when they are told by namenode. They report back to namenode periodically with lists of blocks that they are storing.

Secondary Namenode → To avoid losing of all files of the system since no way of knowing how to reconstruct files from the block on data nodes (namenode to resile to failure). This don't acts as a namenode.

*Role : Merge namespace image & edit log file to avoid becoming a large file.

It runs on a physical machine & it needs plenty of CPU as much as the namenode needs to merge namespace image & edit log file.

It keeps copy of merged namespace image which can be used when namenode is failing.

JobTracker → It is a master process. Responsible (submission).

- accepting job submissions from clients.
- scheduling tasks to run on worker nodes.
- providing administrative functions such as worker health & task progress monitoring to the cluster.
- When job is submitted, info about each task that makes up the job is stored in memory.

Tasktracker → Accepts task assignments from the jobtracker

4/7/24

Job Tracker → It is a master process, accepting job submissions from clients.

Scheduling tasks to run on worker nodes. Providing administrative functions such as worker health & task progress monitoring to cluster.

- 1 jobtracker per map reduce cluster. Runs on hardware. Since failure of master results in failure of all running jobs.

- Clients & task trackers communicate with jobtracker by RPC (Remote Procedure Calls)

- TaskTracker informs Job Trackers as to their current health & status by way of regular heartbeats. Each heartbeat contains no. of map & reduce tasks available, no. occupied, detailed info. about currently executing tasks.

- A tasktracker is dead after period of no heartbeat.
 - Job tracker uses a thread pool to process heartbeats & client requests in parallel.
 - Job submitted & info. about task that makes up job is stored in memory.
- Task Tracker :-
- Accepts tasks from Jobtracker, executes locally, reports progress back to job tracker periodically.
 - Single tasktracker on each worker node.
 - Tasktracker & datanodes run on same machine, makes each node both a compute node and storage node.
 - Configured with specific no. of map & reduce tasks slots which runs in parallel.
 - Tasktracker uses a list of user-specified directories holds intermediate map o/p & reducer i/p. Required bcz of large data to fit in memory for large jobs. Many jobs run in parallel.

Hadoop Distributed File System :-

- Manages the storage across network of machines.
- They are network based, distributed filesystems more complex than regular disk filesystems.
- Hadoop comes with distributed filesystem called HDFS.
- HDFS stores large files with streaming data platforms running on clusters of commodity hardware.

- HDFS features → Can process 100s of TB of data. 2024/8/8 18:41
- It built around idea of data processing pattern Write Once and read many times.
 - Don't require expensively, highly reliable hardware to run on.
 - Designed to run on clusters of commodity hardware.

- HDFS limitations → Appln. require low-latency access to data. HDFS optimized to delivering high throughput data which may be high expensive. HBase is best for low latency access.

- It is not suitable for small sized large no. of files.
- File in HDFS written to by a single writer. Writers are always made at the end of file.

Components of HDFS

- NameNode (Job Tracker) NO ; DataNode (Task Trackers) N1, 2, 3, 4
- File can be larger than any single disk in the network.
- It simplifies storage of subsystem Storage subsystem deals with blocks, simplifying storage management.
- Blocks fit well with replications for providing fault tolerance & availability.
- HDFS blocks are larger than normal disk blocks.
- By making the block size large make the transfer of data in less time.
- Files in HDFS broken into chunks and stored as independent units.
- Map tasks work at a block level.

HDFS Architecture :-

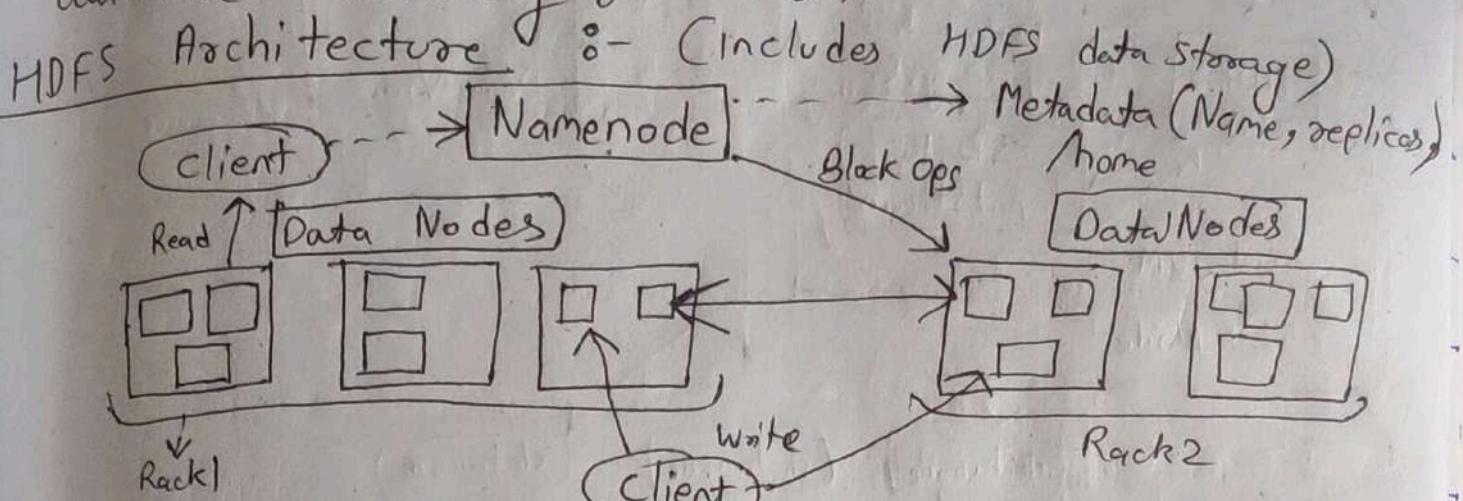
↳ Distributed file system, Designed to run on commodity hardware, Highly fault tolerance & high throughput access to application data & suitable for applications that have large datasets

[DFS, Designed to run on commodity hardware, high fault tolerance, high throughput access to application data & suitable for applications that have large datasets.]

HDFS - Assumption & Goals :-

Hardware failure, streaming data access, Large datasets simple consistency Model, Moving computation is cheaper than moving data, Portability across Heterogeneous Hardware & Software Platforms.

→ Architectural goal of HDFS : Detection of faults & quick automatic recovery from them.



→ Namenodes & DataNodes : HDFS build using Java, Java can run Namenode or Datanode software -

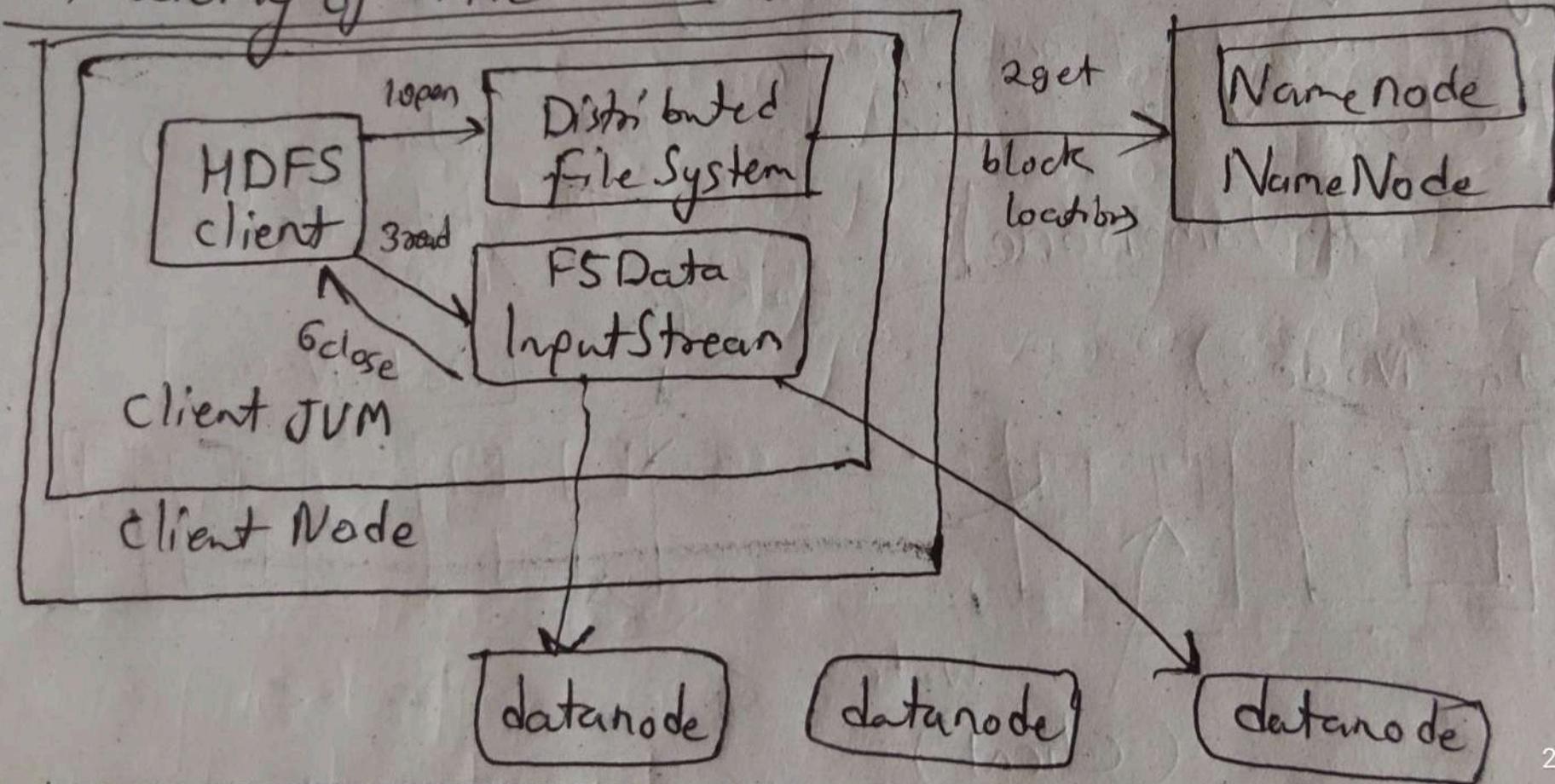
- HDFS has master/slave architecture.
- HDFS cluster consists of single namenode.
- No. of datanode of one namenode.
- HDFS save/stores user data in files. A file split into one or more blocks stored in a set of data nodes.
- = Name node executes file system namespace operations. [Opening, closing, renaming files & directories]
- Data nodes responsible for read & write requests from file system.
- Mapping of blocks to dataNode.
- DataNode performs block creation, deletion & replication upon instruction from Name Node.

2024/8/8 18:42

HDFS Architecture :- - NameNode & DataNode, Data replication, Replica placement, The persistency of file system metadata, Robustness [Data disk failure, heartbeats & re-replication, Data Integrity], Data Organization [Data blocks, Replication Pipelining]

HDFS Architecture :- - NameNodes & DataNodes, File System Namespace, Data Replication [Replica placement, The persistency of file system Metadata, Robustness (Data disk failure, heartbeats & Re-replication), Data Integrity, Meta data disk failure, Snapshots, Data Organization (Data blocks, Replication Pipelining)]

Anatomy of File read :-



2024/8/8 18:42

spth

```
import os
from pyspark import SparkConf, SparkContext
from pyspark.sql import SparkSession
```

import os

```
from pyspark import SparkConf, SparkContext
```

from [Core-site.xml, hdfs-site.xml, mapped-site.xml, yarn-site.xml]
core, hdfs, mapped & yarn

Core-site.xml → This file contains configuration settings of Hadoop

- hadoop.tmp.dir
- fs.defaultFS

hdfs-site.xml → This file contains configuration setting of HDFS

- dfs.replication
- dfs.permission

mapped-site.xml → The file is for configuration related to MapReduce framework

- mapreduce.framework.name

yarn-site.xml → This file configuration for the YARN framework which manages various applications

- yarn.nodemanager.aux-services
- yarn.resourcemanager.hostname
- yarn.nodemanager.aux-services
 - aux-services
 - aux-services.mapreduce.shuffle.class → resource.memory-mb
- yarn.resourcemanager.hostname
 - hostname
 - webapp.address

2024/8/8 18:43

core-site.xml → Configuration settings of Hadoop
mapred-site.xml → Configuration of mapreduce framework
hdfs-site.xml → Configuration of HDFS settings
yarn-site.xml → Configuration of YARN framework
manages various applications.

Core-site.xml :

↳ hadoop.tmp.dir
io.file.buffer.size
fs.toash.interval
fs.default.name

mapred-site.xml :

↳ mapreduce.map.memory.mb
mapred.reduce.tasks
mapred.child.java.opts
mapred.job.tracker
mapred.local.dir
mapred.child.ulimit
mapred.tasktracker.map.tasks.maximum
mapred.tasktracker.reduce.tasks.maximum
io.sort.md
io.sort.factor
mapred.compress.map.output
mapred.map.output.compression.codec
mapred.output.compression.type
mapred.job.tracker.handlers.count
mapred.jobtracker.taskScheduler
mapred.reduce.parallel.copies
tasktracker.http.threads
mapred.reduce.slowstart.completed.maps

hdfs-site.xml :

- ↳ dfs.balance.bandwidthPerSec
- dfs.block.size
- dfs.datanode.data.dir
- dfs.namenode.name.dir
- dfs.namenode.checkpoint.dir
- dfs.permissions.supergroup
- dfs.datanode.du.reserved
- dfs.namenode.handler.count
- dfs.datanode.failed.volumes.tolerated

dfs.hosts

dfs.hosts-exclude

yarn-site.xml :

- ↳ yarn.nodemanager.aux-services
- yarn.nodemanager.auxservices.mapreduce.shuffle-class
- yarn.resourcemanager.hostname
- yarn.resourcemanager.webapp.address

→ hdfs dfsadmin -repo

2024/8/8 18:43

```

<c>
  <p>
    <v></v>
  </p>
</c>
  
```

`return(date, open, close)`

`sdt = spark.SparkContext.textFile("...")`

`sdt = sdt.map(parseLine)`

`sdt.take(2)`

`o_sdt = spark.read.option("InferSchema", "True").`

`.load("...")`

`o_sdt = o_sdt.toDF([" ", " "]).addMap(lambda row: [row[0], row[1], row[2] * 1000], 2)`

`o_sdt = o_sdt.filter(lambda x: x[2] > 2000)`

core-site.xml :

↳ `hadoop.tmp.dir`

mapred-site.xml :

↳ `mapreduce.map.memory.mb`

`mapred.reduce.tasks`

`mapred.child.java.opts`

`mapred.job.tracker`

`mapred.local.dir`

`mapred.child.ulimit`

`mapred.tasktracker.map.tasks.maximum`

`mapred.tasktracker.reduce.tasks.maximum`

`mapred.job.tracker.handler.count`

`mapred.jobtracker.taskscheduler`

HDFS-site.xml :

↳ `dfs.block.size`

`dfs.datanode.data.dir`

`dfs.namenode.data.dir`

2024/8/8 18:43

yarn-site.xml :

- yarn.nodemanager.aux-services
- yarn.nodemanager.auxservices.mapreduce.shuffle.class
- yarn.resourcemanager.hostname
- yarn.resourcemanager.webapp.address
- mapreduce.map.memory.mb, mapred.reduce.tasks, mapred.child.java.opts

Hadoop → open source framework designed to handle large amount of data which ranges between peta to exa bytes.

Data handles : structured, semi-structured & unstructured data

Components → HDFS : Provides throughput output access to system

YARN : Manages & schedules resources in cluster

MapReduce : Process the large datasets with parallel, distributed algorithm.

- HDFS architecture → Master/Slave architecture

Namenode (Job Tracker)	- Master	[Manages file system namespace, controls access to nodes]
DataNode (Task Tracker)	- Slave	[Manage storage attached to nodes, Data is splitted and distributed across each data node.]

NameNode → HDFS manages the file system & (metadata & namespace) of file system

Datanode → Data blocks responsible for read & write requests from file system.

Secondary NameNode → Not a backup of Namenode. Periodically it merges namespace image with edit log to prevent the edit log from becoming too large.

- YARN → Resource management layer of Hadoop

Allows multiple data processing engines to handle data stored in a single platform. Providing resource management & job scheduling.

Components → Resource Management : Allocates resources

Node Manager : Monitors resource usage on each node and reports to resource manager.

Application Master : Manages the lifecycle of applications running on YARN cluster

MapReduce programming model → It involves 2 steps
2 steps : Map & Reduce
Map - Process i/p data & generates key Value pairs.

Reduce - Takes % from Map tasks & agg segregates it to produce final result.

Combiner in MapReduce → Optional function that can perform a local reduce task on the output of the map tasks. It helps in minimizing amount of data transferred b/w Map & reduce tasks.

HDFS → Ensures data reliability through replication - data block is replicated across multiple datanodes. Maintains 3 copies of each block.
Block is amount of data to read or write

Hadoop performance → Tuning no. of Mappers or reducers
Combiners to reduce data transfer
configuring hadoop parameters

Speculative execution → Process where duplicate copies of slow tasks are executed on other nodes.

1st task to complete is used & other tasks are killed
It helps in reducing overall job execution time.

Setup Hadoop cluster → - Install hadoop on all nodes 2024/8/8 18:46
- configure core-site.xml, hdfs-site.xml,
mapred-site.xml, yarn-site
- Formatting namenode
- Start the HDFS & YARN daemons

Hadoop-env.sh → To set environment variables that effect Hadoop daemons - JAVA_HADOOP.

Hadoop-1.x. → Uses JobTracker / TaskTracker for MapReduce

Hadoop-2.x. → Uses YARN resource management & job scheduling. It supports multiple data processing frameworks other than MapReduce

Issues & troubleshoot Hadoop → Namenode failures, network connectivity problems & data corruptions.

Troubleshooting involves checking log files, using monitoring tools & verifying configuration settings.

↳ Install

update

group add

system configuration

disk format

authorized key

ssh root@localhost

download hadoop

change owner

configuration in .bashrc

check java & hadoop version

xml properties

jps

hdfs input & output

jar file

pwd, cp, Variable data files, cat
locate, Paas, chown, FTP

df, HFS+, Useradd, -bashrc, ipconfig,
scp, NTFS, kill, has variable directory,
Sas, apt-get update

/root/.ssh

cloud computing

↳ security, to access or transfer data
via internet.

example: - Laptop, Network, Tablet Computer
Database, mobile, server, smartphone

It is easy to access it present

a speed storage to
public private & hybrid cloud

↳ yum package

chmod is to change file permission

default interactive shell - Bash

Nano - text editor for modal editing capabilities

Systemctl to start-stop & manage services

ps to list currently running process

/bin contains configuration files in Linux

Ubuntu - Often used in enterprise environment

SSH for secure remote login

VMware is native to Linux & provides hardware assisted virtualization - on

11 | 7 | 24

2024/8/8 18:46

19/7/24

Big Data

Big Data : Handles large data sets, data storage.

- Skills → Data analysis & programming [Python, R, SQL]

- Tools → Big data tools [Hadoop, spark], DV, ML & statistics

- Sources → IoT, transactional, web, govt..

Big data Adoption : Drivers & Barriers

(Adoption of Big data ↑) (↓ Cost, lack of skilled person)
- Rel, data privacy

Research & changing nature of data repositories :

↳ Trends & Focus

Trends - Advancements in CC, distributed databases, distributed

Focus - ↑ data storage, retrieval speed, ensuring data integrity & security

Data sharing & reuse practices & their implications for Repository data curation :

practices -

Implications -

Hadoop versions :

V1.0 :- HDFS & Map Reduce

- Limitations in terms of scalability & resource management

V2.0 :- YARN, (Separates resource management from processing layer).

enhancing scalability & flexibility

MapReduce + [Graph processing, Interactive]

V3.0 :- Supports erasure coding, Querying, improved storage efficiency, addition of docker containers.

enhancements to HDFS, YARN, MapReduce,

2024/8/8 18:47

Hadoop & Components :-

- ↳ HDFS :- ↑ throughput access to application data.
 - splits files into large blocks & distribute across nodes.
- YARN :- schedules jobs and manages clusters resources,
 - Allows multiple data processing engines to run on single Hadoop cluster
- MapReduce :- 2 main tasks
 - Map (Processes & filters data)
 - Reduce (Aggregates & Summarizes)

Hadoop vs Traditional DB :-

- ↳ RDBMS handle

DFS (Distributed File System) :-

- ↳ Data stored across multiple servers & locations, for availability, reliability & scalability.
- ↳ Distributed environment, enabling efficient data processing using tools like MapReduce, YARN & Hadoop Components.

Components of HDFS :- NameNode & DataNode

- ↳ NameNode - Master Node manages file system namespace
 - Regulates access to files
- ↳ DataNode - Worker/slave nodes that stores & retrieves block when they are told.
 - Report back to NameNode with block that they are storing. Periodically.

HDFS daemons :-

- ↳ NameNode → Manages metadata & filesystem namespace. It keeps on tracking.
- DataNode → Daemons that manage storage on individual nodes - r and w request from a client
- 2nd NM → Merging namespace image & edit log to reduce load on NN.

↳ HDFS follows master-slave architecture

Name Node

master-slave architecture

Data Node

Scaling & Rebalancing :-

Scaling → Horizontally adds datanodes to cluster. ↗
↑ storage capacity

Enhances parallel processing capabilities.

Rebalancing → distribution of dataNode become in balance
(Data Node)

due to data growth or node additions/removals
[HDFS balances & redistributes]

Replications :- Replicates data blocks across multiple data
-Nodes.

Rack Awareness :- It uses a rack awareness algo to improve
data reliability & bond with utilization.

Data pipeline :- Client writes to Data Node in sequence
1st DataNode then last next 2nd DataNode

Node failure Management :- Handles Node failures
automatically. If DN fails the NN automatically
allocates initiates the replication of the block

ZooKeeper

HDFS architecture :- Manages & Processes data

2024/8/8 18:48

sets

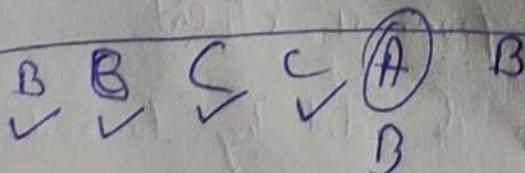
Design principles → Replications, rack awareness,
high availability

Features can build big data solutions.

↳ Drop → Deletes entire table

Delete → Removes row from table

Truncate → Removes all row data



20/7/24

Hadoop Operation Mode :-

- Standalone Local Mode (HDFS on a single Machine) [dough & turn]
- Pseudo distributed (HDFS runs on single machine simulated distributed environment)
- Fully distributed (Runs on multiple machines ; HDFS and YARN separate environment)
 - Nodes : NN, DN, RM, NM
 - [Production Environment]

- ↳ core-site.xml :
- * HDFS : Configures Hadoop properties
 - mapped : " DataNode & NameNode
 - Yarn : " MapReduce ; job history & Job Framework

Resource Manager → Manages resources & schedules jobs
 Node Manager → Manages individual nodes resources & monitors resource usage

Data Node → Stores actual data block
 Name Node → Manages file system namespace & metadata
 Secondary NN → Periodically merges namespace image & edit log.

Job Tracker (MRVI) → Manages job scheduling & task tracking.

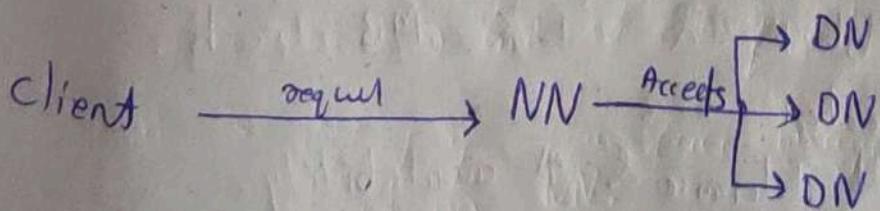
- Resource Manager : Handles resource management
- Application Master : Takes care of Job scheduling & task monitoring.

Task Tracker (MRVI) → Executes individual tasks assigned by job tracker.

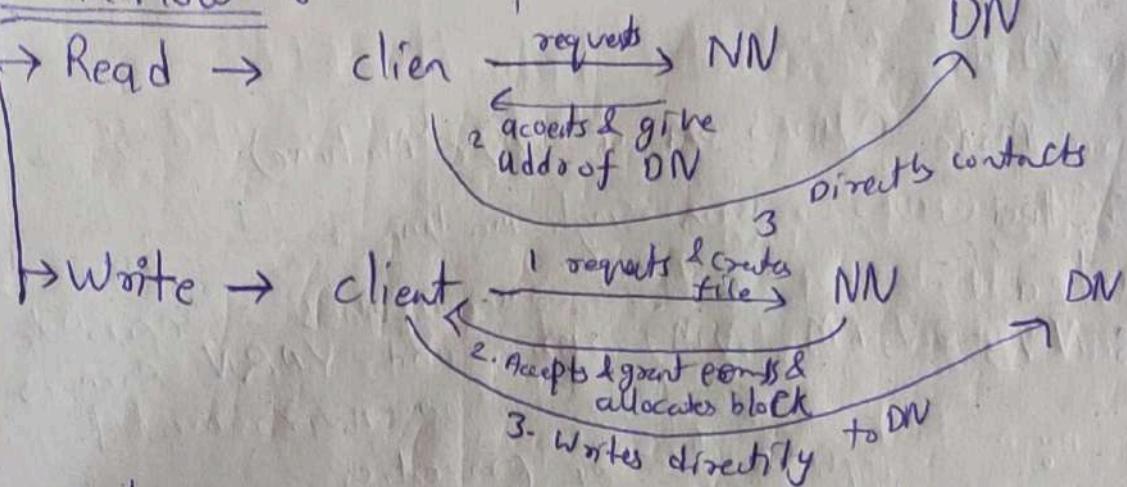
Resource Manager (YARN) → Manages resources & applications

Node Manager (YARN) → Manages resource & containers on a single node

Application Master (YARN) → Manages life cycle of a single application



Data Flow :-



Scaling → HDFS is horizontal scaling. More DN are added to cluster to ↑ storage capacity

- Horizontal scalability : adding data nodes to cluster & handle
- Load Balancing :
- NN federation :

Rebalancing in HDFS :- Ensuring even distribution of data across DN

- HDFS balancer Tool : To redistribute across DN via cmd
- Threshold based Balancing : Threshold parameter to determine when rebalancing is necessary.
- Incremental Balancing : Rebalancing is performed incrementally to minimize disruption to the system.

↳ Optimal performance & resource utilization
HDFS a robust solution

Map Reduce for parallel processing

HDFS commands :-

- Data split into blocks & 128 MB each.
Blocks are replicated → 3
DN & ensures fault tolerance & ↑ availability.

↳ hadoop fs - chmod /path : A permissions
dfsadmin -safe mode enters
dfsadmin -safe mode leave
dfsadmin dfsadmin Refresh the list of DataNodes.

↳ NameNode : 9870

Resource Managers : 8088

↳ Hadoop Map Reduce : processing large datasets in a parallel across a dataset

- Map Reduce tasks → Map [processes I/p data & generates key-value pair]

- Reduce to Reduce [Aggregates intermediate key value pair to produce final o/p]

Map Reduce Daemons :-

→ Resource manager [manages resources]
→ Node Manager [monitors resources]
→ Application Master [manages execution of a file Map Reduce]

Anatomy of a MapReduce Job Run :-

↳ Job Submission - Client submit job to RM
I/p splits - splits into logical chunks processed by separate Map phase - Processes I/p & generates key-value pairs.
Shuffle & Sort - key-value pair shuffled & sorted by keys
Reduce Phase - Processes grouped key-value pair to final o/p.
o/p - final o/p written to HDFS or another storage system.

Partitions & Combiners :-

↳ Partitioner determines which reducer a particular key-values pair will go to. (To control data distribution)

Combiner performs local aggregation of intermediate key value pairs that they sent to reducers.

↳ I/p formats

Input Formats :-

- i/p splits & Records [Logical chunks to Map task generates key-value pairs]
- Text i/p [Default input format each line treated as a record]
- Binary i/p [Sequence file format to read binary key-value pairs stored in Hadoop sequence files]
- Multiple i/p [Allows i/p for different data sources in the same job]

Output Formats :-

- Text Output [key-value pairs written as text lines default o/p format]
- Binary Output [Sequence o/p format (writing binary key-value pairs to Hadoop seq.. files)]
- Multiple Output [Allows multiple data to multiple o/p from a single job]

Distributed Cache :-

↳ Distributes read only data needed by MapReduce jobs to all nodes in the cluster

- add CacheFile() : Adds a file to a distributed cache

- add CacheArchive() : Adds an archive to the distributed cache

↳ Hadoop MapReduce paradigm { Large scale data processing by dividing tasks into smaller subtasks
Distributed Cache performed by making read only data readily available to all nodes.

Hadoop DataTypes ⇒

6 (int, long, float, double, bool, Text) [Primitive]

3 (Array, Map, Null) [Composite]

Java & MapReduce :-

- ↳ Mapreduce written in Java.
- ↳ Robust API for implementing Mapper, Reducer
- ↳ Map Reduce components :-
 - Mapper : define Map function
 - Reducers : define reduce function
 - Driver class : configures job & specifies I/O paths, Mapper & Reducer etc. —

Map Only :-

- ↳ No reduction or aggregate is needed. Reducer is omitted or set to null

Reduce only :-

- ↳ Preprocessing has already done only aggregation is needed.

Combiner & Partitioner :-

- ↳ Combiner - Acts as mini Reducer to reduce data transfer
- ↳ Partitioner - controls the partitioning of intermediate map outputs `job.setCombinerClass()`, `job.setPartitionerClass()`

Counters : Tracks the progress or count the events

• `getCounter()`

Schedulers :- Job scheduling

- ↳ FIFO scheduler - Jobs executed in order they arrived
- ↳ Fair scheduler - Jobs are allocated resources that all jobs get an = share of resource over time.
- ↳ Capacity scheduler - To run Hadoop jobs in a shared, multi-tenant cluster to ensure certain capacity guarantee.

Custom Writables : To define custom data types

- Writable
- Writable Comparable

Compression : Reduces amount of data transferred & stored

- `setCompressOutput()`
- `setCompressorClass()`

Complex Map Reduce Programming

↳ multiple MR jobs can be chained together.] chaining
% of one job being % for next job

↳ different stages of processing having multiple mapped reducers.

Map Reduce Streaming :- writing MR jobs in any language
that can read from standard I/O & write to standard O.

Python & Map Reduce :- Python scripts for both Mappers &
Reducers

Map Reduce on image dataset :-

↳ Image processing (resizing, format conversion, feature extraction)
- Using Py Lib : PIL or OpenCV

↳ Implementation :-

Mapper - Process each image, perform necessary
image processing task

Reducers - Aggregate results if needed.

Hadoop ETL :-

↳ ETL (Extract Transform Load) :- Extracting data from
various sources, Transform it into suitable format
loading into Hadoop's HDFS or similar storage system.

- Extraction : Extract/Load data from files

- Transform : Clean up or convert data to standard or
specific

- Load : Transformed data into HDFS, HBase or
Hadoop compatible storage.

ETL functions :-

- Data Extraction : Sqoop (data b/w Hadoop & RDBMS)

Flume (collects, aggregates - streaming data into HDFS)

Kafka (datastreams processed in real time or loaded) in HDFS

- Data Transformation : pig (pig latin to handle complex data)

Hive (SQL - query for data transform & analysis)

Spark (API for complex data transform)
using RDDs & DP

2024/8/8 18:49

- Data Loading :- HDFS (1^o storage for large datasets)
- HBase (NoSQL DB runs on top of HDFS, r/w access)
- Hive, (Acts as data warehouse storing data in a structured format)

EVent for ETL Tools :- Complex data Management, Automation, Scalability, Integration

Advantages of ETL :- Efficiency, data Quality, Compliance, Flexibility, Performance, User Friendly

HBase :-

→ Open Source, non-relational, distributed database.

→ Written in Java

→ Handles data across commodity servers.

→ Provides real time read/write access to big data.

Features :- column oriented storage, Scalability, Consistency, Integration with Hadoop, High availability.

→ Horizontal scaling.

HBase Architecture :-

→ Components →

- HBase Master : Manages Cluster (Table creation, deletion & balancing) 2024/8/8 18:49

- Region Server : Handles r/w requests. Manages HFiles & WAL (Write-Ahead Log)

- ZooKeeper : Coordinates distributed process & provides configuration & synchronization services.

→ Data Model →

Tables, Rows, columns, col families, cells.

→ Storage →

HFiles, WAL (Write-Ahead Log)

→ Region →

Region - Subset of table horizontally partitioned by row key range

Region Split - Automatically splits it

Installation :-

- Prerequisites : Hadoop cluster installed & running
Zookeeper ensemble configured & running
- Steps :-
 - Download & Extract HBase
 - Configure HBase
 - Start HBase

Operations :-

- Creating a Table - Create
 - Listing Tables - list
 - Inserting data - put
 - Retrieving data - get
 - Scanning data - scan
 - Deleting data - delete
 - Dropping table - disable
drop
- ↳ counting rows
- count

HBase commands :-

- hbase shell - starting the shell
- listing commands - help

Java Client API for HBase :-

- Establishing connection - HBaseConfiguration.createC()
• createConnection()
- CRUD operations -
 - Create : .getTable() , .addColumn()
 - Read : .get() , .getValues()
 - Update : HBase overwrite the value
 - Delete : .addColumn() , .delete()

2024/8/8 18:49

Admin API :-

- Creating :- .getAdmin(), .addFamily(), .createTable()
- Listing Tables :- .listTables()
- Disabling & Dropping Tables :
 - . disableTable
 - . deleteTable

↳ HBase - Scan, Count, Truncate
(scan , count , truncate)

HBase Security

Authentication : kerberos

Authorization : Access Control (Defining permissions)

Encryption : Data encryption (Encrypt using HDFS zones)
Network encryption (Encrypt data in transit)

↳ Robust support for storing & processing large datasets with read-time &/w access.

HIVE

Features : SQL like Query Language, schema on Read, Integration with Hadoop, support for diff.. storage types -

Components : UI, MetaStore, HiveQL process Engine, Exchange Engine, HDFS

Installation :

- Prerequisite : Java JDK, Hadoop, SSH

Datatypes : string, int, bigint, Boolean, float, double, timestamp, Date

Complex types : Array, Map,
Struct, UnionType

String functions : CONCAT(), LENGTH(), UPPER(), LOWER()

Date functions : CURRENT_DATE(), DATE_ADD(), DATE_SUB()

Hive tables : Managed tables , external tables
(manages data) (manages metadata only)

Partitions & Buckets :-

Partitions → Divides ~~tables~~ into parts based on values of particular column.

Buckets → Divided into more manageable parts

Used for more efficient sampling & joining.

Storage formats :- Text file, Sequence file,

- Text files : Default storage format (text-table)

- Sequence files : seq-table

- ORC : (Optimized Row Columnar)

- Parquet : parquet-table

Importing data :- LOAD DATA command by creating external tables.

Querying Data :-

→ Sorting data : ORDER BY, SORT BY, DISTRIBUTE BY

→ Aggregating data : COUNT, SUM, AVG, MIN, MAX

Map Reduce scripts :- Hive can translate HiveQL into MapReduce jobs.

Joins & Subqueries :- SELECT, FROM, WHERE,
INNER JOIN, SUBQUERY

Views are virtual tables created using the result of a query.

Map and Reduce Side joins to Optimize Query

Map side and Reduce side joins to optimize query.

map side join : smaller dataset fit into memory.

Reduce side join : Default join in HIVE, suitable for large datasets.

Data manipulation with HIVE :-

Inserting data, Updating data, deleting data

UDFs (User Defined Functions) :-

↳ Creates custom functions to extend HiveQL capabilities -
can be written in Java & other languages.

- Creates Java class extending UDF.

- Implement the evaluate method

- Compile & add the JAR to HIVE.

Custom mapReduce in HIVE

↳ Scripts can be integrated into HIVE queries
(TRANSFORM).

Writing HQL scripts :-

↳ HQL allows for the batch processing of HIVE queries.

AQL script is a text file containing HIVEQL commands.

For executing Hive : hive -f script.hql

2024/8/8 18:50

Hadoop :- Open Source \$, to handle large datasets - Allows
for storage & processing of big data.

Components : HDFS, YARN, MapReduce

HDFS :- Provides scalable & reliable data storage designed
to span large clusters of commodity servers
Architecture : NameNode, DataNode, 2^o Node, HDFS daemons (NameNode & DataNode)

HBASE :- Distributed, scalable,
stores and manages large amount of sparse data
Architecture : HBase Master, Region Servers, Zookeeper
HFile

Hive :- Data warehouse built on top of Hadoop.

- Data summarises, querying & analysis.

Architecture -

Metastore

Hive Server

Hive Driver

Execution engine

↳ Hadoop : Provides distributed storage (HDFS) & processing data (Map Reduce)

HDFS : Handles large datasets with high reliability & scalability

HBase : NoSQL DB built on top of HDFS, real-time read/write access to large dataset

Hive : Data warehousing that provides SQL like interface for querying & analyzing data stored in Hadoop.

Apache Spark :-

↳ open source, distributed computing system. Provides interface for programming entire cluster with implicit data parallelism & fault tolerance.

→ cluster with implicit data parallelism & fault tolerance

→ To overcome limitations of Hadoop MapReduce framework and provide faster, more efficient data processing.

Linking with Spark →

↳ Hadoop

Apache Hive

Apache HBase

Apache Kafka

Initializing spark :-

↳ `SparkSession.builder()`

RDD (Resilient Distributed Datasets) :-

↳ Properties : (Immutable, Fault tolerance, Parallel processing)

• RDDs fundamental data structures of Spark.

Creating RDDs :- • `parallelize(data)`

↳ RDDs : Low Level APIs provide transformation actions.

↳ `pyspark.sql import SparkSession`

`sparkSession = builder().appName("None").getOrCreate()`

`sc = spark.sparkContext`

`data = [, , ,]`

`rdd = sc.parallelize(data)`

~~# rdd = sc.textfile(" --- ")~~

Shuffle Operations :-

↳ `reduceByKey` & `groupByKey`

RDD persistence :-

↳ It can be persisted in memory or disk to avoid recomputation.

Removing data :- To unpersist data [`rdd.unpersist()`] 2024/8/8 18:50

Shared Variables :-

- Broadcast Variable : send only value to all nodes

- Accumulator : variables added to from workers & read only by drivers.

Deploying to a cluster :-

- standalone mode : spark manages own cluster

- YARN : spark appl. on Hadoop cluster

- Mesos : cluster manager used to run spark.

Map Reduce with Spark :-

- Spark [with hadoop] → Hadoop HDFS for storage
YARN for resource management
- [without Hadoop] → Use local file systems or other storage solutions like S3.
- Spark's standalone cluster manager.
- ↳ Storage : with Hadoop - HDFS
without Hadoop - Local file system or cloud storage
- Reason Management : with Hadoop - YARN
without Hadoop - Standalone or Mesos

Data Preprocessing :- Step in data analysis pipeline.

Transforms raw data into a clean & usable format.

Steps → Data cleaning

Data Transformation

Data Reduction [(PCA), Sampling]

Feature Engineering

EDA (Exploratory Data Analysis) :-

2024/8/8 18:50

↳ Analyzing data sets to summarize their main character
istics, often using visual methods..

Understanding the data & uncovering patterns, anomalies,
relationships.

- Descriptive Statistics [mean, median, mode, var, SD]
- Data visualization [Plots]
- Correlation Analysis [Relationship b/w variables using correlation Coefficients]
- Outlier detection [Identifying the outliers in data]

Apache Kafka :-

- Designed for high throughput, low latency in data streaming.
For build a real time data pipelines & streamin applications.

spark streaming Architecture →

- Discretized Streams (DStreams) [Spark streams represent streams of data → Series of RDDs]
 - Receivers [To ingest data from sources like kafka, Flume & others]
 - Processing [Each batch data processed by Spark jobs, can perform transformations & actions on data]
 - Output Operation [Processed data can be written to various output sinks (HDFS, DB, ...)]
- ↳ Apache spark & kafka for building scalable, real time data processing pipelines - combining capabilities of Spark for processing

kafka connect API :-

- Source Connectors (Import data $\xrightarrow{\text{from}}$ Ext.. syst into kafka)
 - Sink Connectors (kafka $\xrightarrow{\text{export}}$ Ext.. systems)
- ↳ Apache Kafka :- written in java & scala.
- Messaging, Storage, Processing
 - Topics, Partitions,
 - Producers & Consumers
 - Brokers
 - Zoo Keepers

Apache Spark

2024/8/8 18:51

21/7/24

Apache kafka :-

→ Open source stream processing platform. Written in Scala & Java.

Functionalities : Messaging, storage, processing
 Categories where records are stored

Concepts : Topics - Segments of topics that allow parallel processing

Partitions - Data written to Kafka topics

Producers - Data Read from Kafka topics

Consumers - Kafka nodes handle storage & retrieval of records.

Brokers - Manages & coordinates Kafka brokers.

Zookeepers -

Apache Spark :- Open source unified analytics engine for large scale data processing.

Provides high level API in java, scala, Python, R

Components : Spark Core, Spark SQL, spark streaming, MLlib, GraphX

Spark SQL :- Spark module for structural data processing.

It provides DataFrame API & SQL queries.

Features : DF & Datasets → DF (Data organized into named columns)

Dataset (Provides benefits of RDD optimized advantage of Spark SQL)

SQL Integration → Allows executing SQL queries along with complex analytic algorithms.

Catalyst Optimizer → Query Optimizer for executing queries efficiently.

Data Source API → Allows integration with variety of data

Spark MLlib :- ML lib built on top of Spark.

Features : Algo - Classification, Regression, Clustering, Collaborative filtering etc.

Factorization - Factor extraction, LDA, dimensionality reduction.

Pipelines - Tools for constructing, evaluating & tuning ML

Utilities - Statistics, data handling & other utilities.

Predictive analysis with Spark :-

↳ Using historical data to predict future outcome

steps - Data Preprocessing (cleaning & transforming)

Feature Engineering (Creating new features from existing)

Model building (LR, DT)^{data}

Model Evaluation (Metrics like accuracy, precision,)

Model deployment (Integrating model into production)
Systems to make predictions

Workflow in Spark :-

↳ Initialize Spark Context → Load & preprocess data → Feature
Engineering
 ↓
 evaluate Model ← Train Model ← Split data
 ↓
Deploy model

Linking kafka with Spark :-

↳ Kafka used as source and sink for Spark streaming application.

Steps : - Setting up kafka

- Creating Spark Streaming Context

- Consuming data from kafka

- Processing data

- Producing data to kafka

2024/8/8 18:51

→ Apache kafka & Apache spark powerful tools for real time, scalable, data processing pipelines.

→ Spark SQL and MLlib provides robust frameworks for data processing & ML, enabling predictive analysis.

→ Integrating kafka and spark, creates efficient streaming applications that handles large volumes of data with low latency.

- Benefit of dashboard actions : To link visualization and enable user interactivity.
- Automatic layout helps to design dashboards that adapt to different screen sizes.

Hadoop

- ↳ purpose of data replication : Ensure fault tolerance & data availability.
- ↳ HDFS_High Availability (HA) : provides high availability by allowing a standby Namenode.
- ↳ ACLS (Access Control Lists) in HDFS allows admin to do to manage access to files & directories. (User permissions)
- ↳ Kerberos authentication to authenticate users & services in a secure manner.
- ↳ YARN → MapReduce jobs
- ↳ Pig is a Hadoop component which provides a high-level abstraction over MapReduce.
- ↳ MapReduce Framework is to process the data in parallel across a distributed cluster.
- ↳ MapReduce framework provides fault tolerance through data replication.
- ↳ Hive is a tool in hadoop ecosystem which is used for data warehousing & SQL like queries.
- ↳ Apache Hbase is to provide real time read/write access to large datasets.
- ↳ Sqoop tool is used to import & export data b/w Hadoop & RDBMS.
- ↳ Apache Flume is used to data ingestion from various sources.
- ↳ Job Tracker : Coordinates MapReduce jobs
- ↳ Task Tracker : Execute tasks assigned by Job Tracker
- ↳ Secondary NameNode : Merges namespace image & edit log ; to help in recovering the file system metadata in case of a failure.
- ↳ Feature of Hadoop data storage model : Data stored in large, fixed size blocks.

- ↳ Rack Awareness: distribution across the different racks to improve fault tolerance & reduce bandwidth usage.
- ↳ Inter-rack Awareness is to authenticate users & services in a secured manner.
- ↳ Hadoop Key Management Server (KMS) is to manage & store cryptographic keys for data encryption.
- ↳ Namenode: Manages distributed File System (DFS)
Maintains metadata of HDFS & manages file system namespace.
- ↳ Job Tracker (+1v)
Resource Manager (+2v) :
(YARN)
 - ↳ Coordinates Map Reduce jobs
 - ↳ Manages scheduling & resource managing, execution of MapReduce jobs
- ↳ DataNodes: It stores the output data.
- ↳ Monitoring (Cluster monitoring Hardware) is handled by cluster monitoring tools [like, Ambari, Ganglia or Nagios] (used in hadoop ecosystem for monitoring purpose)
- Resource Manager is role of managing resources in the cluster in YARN
- Task Tracker in (Ver1) or Node Manager in (Ver2+) responsible for executing tasks assigned by Job Tracker or Resource Manager
- Job Tracker in (ver1) or Resource Manager in (ver2+) is responsible for scheduling & managing execution of jobs in the cluster.
- Hadoop (Key Management Services) KMS is to manage & store cryptographic keys for data encryption.
It handles generation, management & storage of encryption keys for securing data in HDFS.
- Monitoring the health of the Hadoop cluster is done by external tools like Ambari, Ganglia or Nagios used for system & application monitoring & including Hadoop clusters.

- Hadoop : Open source framework, process large datasets in a distributed computing environment & performs parallel computing across the cluster of machines.
- Components : HDFS, MapReduce, YARN
- Zookeeper : is a centralized service for maintaining configuration info., naming, providing distributed coordination for distributed application. Often used to manage coordination & configuration of Hadoop clusters and services like HBase & kafka. It manages failures.
- HBase : Distributed, scalable, NoSQL database built on top of HDFS. Provides real time &/or access to large datasets. Supports structured data as a key value pairs & designed to handle large tables;
- Application : real time analytics & random real time &/or access to big data.
- Hive : Data warehouse infrastructure built on top of Hadoop. Provides tools ~~like~~ for data summary, query and analysis & SQL like language called HiveQL. Abstracts complexity of Hadoop's MapReduce programming. Used for batch processing querying large data sets in HDFS.
- Apache Spark : Open Source, distributed computing system. In memory data processing engine, faster than hadoop Map Reduce. It supports batch processing, real time stream processing & ML, graph processing. Offers API in Java, Scala, Python & R. 2024/8/8 18:52
- Apache Flume : Centralise data store, Used to ingest data into Hadoop for further processing and analysis. Collects data from various web servers, appln. servers, & other sources & store it in HDFS, HBase.
- HDFS (Hadoop Distributed File System) : Stores large amount of datasets reliably & stream data sets at high bandwidth. Splits large files into blocks (128 MB). Providing fault tolerance, ensures data redundancy by replicating each block multiple times across different nodes.

- ↳ DFS (Distributed File System) : It refers to HDFS. Manage & store data across a distributed network of machines. Provides ↑ throughput access to application data & capable of storing v large files.
- ↳ YARN (Yet Another Resource Manager) : Management layer of hadoop. Separates resource management & job scheduling functions from the data processing components. Allows multiple data processing like MapReduce spark. to run & process data stored in HDFS. It contains Resource Manager (Arbitrates resources among all applications) Node Manager (Monitors resource usage on each node)

↳ Map Reduce : Programming model & implementation for processing large datasets in a distributed fashion.

Map Phase (processes i/p & produces key value pairs)

Reduce Phase (Intermediate results to produce the final o/p - Written in Java).

↳ PIG : High level platform for creating MapReduce Scripting language (Pig Latin) Abstracts complex MapReduce programming model. Used for data transformation, data aggregation & analysis tasks.

↳ Oozie : Workflow scheduler system to manage Hadoop jobs. Define sequence of actions to be executed. in a workflow (MapReduce jobs, Pig scripts, Hive queries) Supports scheduling & chaining of workflows enabling complex data processing pipelines in Hadoop.

↳ Sqoop : Tool for transforming bulk data between Apache Hadoop & structured data stores such as RDBMS. Imports data from external database into HDFS or HIVE. exporting data from hadoop back to database. minimizes amount of manual coding required.

↳ Flume : For efficient collecting, aggregating & moving large amount of log data or other data streams from various sources to a centralized data like HDFS or HBase.

↳ HCatalog : Table & storage management layer for Hadoop provides schema & storage layer for Hadoop tools. Allows to share data between different tools like Pig, MapReduce & Hive. Provides unified view of data stored in HDFS.

Hadoop Daemons :-

↳ NameNode : Manages metadata & namespace of HDFS, keeps track of data where data is stored across the cluster. Operations (open, close, renaming files/directories).

↳ Secondary NameNode : Periodically (checkpoints of HDFS metadata to prevent data loss). Merges edit logs & namespace images (fsimage). Keep NameNode's memory usage efficient.

↳ Resource Manager : Part of YARN, Manages resources across the cluster. Schedules and allocates resources to applications based on their needs.

↳ Node Manager : Manages resources of a single node in the YARN cluster. Monitors resource usage (CPU, memory, disk) & reports to resource manager.

↳ Job History Server : Tracks completed MapReduce jobs & stores information about history for debugging & auditing purposes.

↳ Task Tracker (Deprecated in v2) :

In Hadoop 1, Task Tracker managed task running on individual nodes.

In Hadoop 2, Replaced by Node Manager Under the YARN framework.

↳ (A, B, B, C, B, B, C) , $(A, B, A, B, B, A, D, A, C, C, B, A, B, B)$

→ Histogram : frequency of data points within specified range within a specified range. [distribution, shape & spread of data]

→ Year() to extract year from a date value in many

→ EOMONTH() last day of the month for a given date.

LOD (Level of Detail) :- Data analytics & visualization

To define how data should be aggregated in calculations

Concepts of LOD :-

Granularity : Level of detail or specificity in data.

Aggregation : Process of summarizing data

Types of expressions : Fixed, Include, Exclude

- Fixed :-

- Include :-

- Exclude :-

Use cases of LOD :-

- Consistent Benchmarks

- Dynamic Scoping

- Removing unwanted details

Gantt chart : For project management & scheduling. It displays timeline for tasks or activities, showing their start and end dates for how they overlap.

↳ Plots :-

Histogram → Distribution of data points across a continuous variable. Frequency distribution of a continuous variable. Distribution of data points in a dataset. To display distribution of data. Frequency of events.

- ↳ Scatter plot : Relationship b/w numerical variable . Correlation between 2 variables
- ↳ Bar chart : Comparison of sales / Quantity across category . Comparison across categories , compare categorical data ,
- ↳ Gantt chart : The progress of tasks over time in a project
- ↳ Line chart : Trend over time , changes in a variable over time , trends over time , Time series data
- ↳ Tree map : Treemap displays (Proportional data in a hierarchical data).
- ↳ Stacked Map chart Mosaic Plot : Relationship b/w 2 categorical variables
- ↳ Pie chart : Proportional relationships b/w parts of whole , To illustrate parts of a whole.
- ↳ Area chart : To show cumulative total over time

Apache Spark :- Open source , distributed computing system provides a fast & general purpose cluster computing framework . Large scale data processing and handles batch processing , interactive queries , Stream processing , ML & graph processing . Speed , ease of use , wide range of workloads .

- Speed \rightarrow 100x times faster than Hadoop MapReduce

- Ease of use \rightarrow Java , Scala , Python , R .

- Unified engine \rightarrow Cover a wide range of data analytics :

- Batch processing , stream processing , ML , Graph processing .
- Distributed \rightarrow Operates distributed clusters , scale horizontally across multiple machines .
- Resilient Distributed Dataset (RDD) \rightarrow Data structure of Spark . Distributed collection of object can be processed in parallel
- ~ Data Frame & dataset API , In memory Computing , Integration

↳ Use cases of Apache spark :-

Data processing & ETL - Extract Transform Load
Real time data processing - Apache, Flume, Kafka, Kinesis. real time analytics and monitoring.

ML -

Interactive data Analysis -

Graph Processing -

Architecture of Apache Spark :-

Driver - Controls the flow of data and execution of tasks.

cluster Manager - Allocates resources across the cluster. builtin cluster manager, Apache Mesos, Hadoop YARN

Workers - Responsible for executing tasks assigned by the driver. Each worker node hosts one or more executors.

Executors - Executes individual tasks & keep data in memory. They communicate with the driver & perform the actual computation.

HBase :- Distributed, scalable, high performance 25/7/24
NoSQL database built on (HDFS). Real time r/w access to large datasets, designed to handle large amount of sparse data.

key features → Realtime r/w access, Horizontal scalability, Column oriented storage, Integration with hadoop ecosystem, fault tolerance.

Architecture → Region servers → Region , HMaster, Zookeeper

Use cases → Real time data processing , Large scale data storage

Hive :- Data warehousing & SQL like query language built on top of Hadoop . Data analysis & querying using HiveQL or HQL similar to SQL like language. Handles large datasets & integrates seamlessly with Hadoop.

Feature → SQL like interface , Data warehouse infrastructure , Integration with hadoop , supports various data formats , extensible

Architecture → Metastore , Driver , Execution engine Client Interface

Use cases → Data warehousing , Batch processing

HBase

- NoSQL DB for real time access to large datasets.
 - fast read/write capabilities
 - Built ontop of HDFS
- Hive
- SQL like language HiveQL(HQL) for batch processing & querying large datasets.
 - Built on top of Hadoop

Airflow :- Open source platform to programmatically author , schedule , monitor workflows . Used to define DAGs (Directed Acyclic Graphs) of tasks . Enabling the orchestration & management of complex data pipelines

Features :- DAGs (Directed Acyclic Graphs) , Tasks , Operators (Python Operator , Bash Operator , SQL operator , Sensor Operator)

Scheduling , Task dependencies , Extensibility , Monitoring , Logging , Retry & alert Machines , Scalability , Integration with cloud and external services

Uses : ETL pipeline , Data analytics & ML , Data integration , Data quality & monitoring , Batch processing

```
from pyspark.sql import  
from pyspark.sql import SparkSession  
spark = SparkSession.builder.appName().getOrCreate()  
df = spark.read.csv(" ", header=True)  
df.show()
```

```
from pyspark import SparkContext  
sc = SparkContext(" ", " ")  
rdd = sc.textfile(" ")  
for line in rdd.take(5)  
    print(line)
```

```
sc.stop()
```

1 → Datasource API read
2 → SparkContext read
3 → Pandas to RDD

① from pyspark.sql import SparkSession, Row
spark = SparkSession.builder.appName(" ").getOrCreate()
data1 = spark.read.csv(" ", header=True)
data2 = spark.createDataFrame([Row(,),
 Row(,)])

from pyspark.sql.functions import col
unique_ele = spark.groupby(" ").count().filter(col() >)

```
from pyspark import SparkContext  
sc = SparkContext()  
rdd = sc.textfile()
```

```
from pyspark.sql import SparkSession  
sc = SparkSession()
```

↳ # DataFrame API

```
from pyspark.sql import SparkSession  
spark = SparkSession.builder.appName().getOrCreate()  
from pyspark.sql import Row  
data = [Row(id=1, name="A"),  
        Row(id=2, name="B"),  
        Row(id=3, name="C")]  
df1 = spark.createDataFrame(data)  
df1.show()  
  
df2 = spark.read.csv("", header=True)  
df2.show()  
unique_words = df2.groupby("").count().filter(col("") > 1)  
    .dropDuplicates()
```

↳ ~~from pyspark.sql import SparkSession~~ # API
spark = SparkSession.builder.AppName().getOrCreate()

```
from pyspark.sql import Row  
data = [Row(ids=1, name="a"), Row( ), Row( )]  
df = spark.createDataFrame(data)
```

```
df.show()  
from pyspark.sql.functions import col  
dup_df = df.groupby("name").filter(count() > 1)
```

drop_df = df.dropDuplicates(["IDs?", "name"])

2024/8/8 18:54

↳ ~~from pyspark import sparkContext~~ # Text file

sc = SparkContext("local")

rdd = sc.parallelize([" ", " ", " "]).rdd

for i in rdd.collect():

print(i)

words = rdd.flatMap(lambda line: line.split(" "))

unique = words.distinct()

words_rdd = words.collect()

for i in words_rdd:

print(i)

```
from pyspark.sql import SparkSession  
spark = SparkSession.builder.appName().getOrCreate()  
  
#data = spark.read.csv("")  
#data.show()  
data = pd.DataFrame([{"": [1, "a", 1]}])  
df = spark.createDataFrame(data)  
d)
```

↳ from pyspark.sql import SparkSession
spark = SparkSession.builder.appName().getOrCreate()

~~#~~ from pyspark.sql import Row
data = [Row(), Row(), Row()] ✓
df_1 = spark.createDataFrame(data)
dup_data = df_1.groupby("").count().filter(col() > 1)
drop_data = df_1.dropDuplicates(["", ""])

↳ from pyspark.sql import SparkSession
import pandas as pd
spark = SparkSession.builder.appName().getOrCreate()
data = pd.DataFrame({})
rdd = spark.createDataFrame(data).rdd
~~drop~~_rdd = rdd.distinct()
2024/8/8 18:55
df = spark.createDataFrame(data)
dup_df = df.groupby("").count().filter(col() > 1)
drop_dup = df.dropDuplicates(["", ""])

↳ from pyspark import SparkContext
sc = SparkContext()
data = ["", "", ""]
words_rdd = sc.parallelize(data)
words = words_rdd.flatMap(lambda l:l.split("")) ✓
dup_words = words.distinct()
words_list = words.collect()

Hadoop & Components :-

- HDFS :- → throughput access to application data.
 - splits files into large blocks & distribute across nodes.
- YARN :- schedules jobs and manages clusters resources,
 - allows multiple data processing engines to run on single hadoop cluster

MapReduce :- 2 main tasks Map (Processes & filters data)
Reduce (Aggregates & Summarises)

Hadoop vs Traditional DB :-

- RDBMS handle

DFS (Distributed File System) :-

- Data stored across multiple servers & locations, providing availability, reliability & scalability.
- Distributed environment, enabling efficient data processing using tools like MapReduce, YARN & Hadoop Components.

Components of HDFS :- NameNode & DataNode

- NameNode - Master Node manages file system namespace
 - Regulates access to files
- DataNode - Worker/slave nodes that stores & retrieves block when they are told.
 - Report back to NameNode with block that they are storing. Periodically.

HDFS daemons :-

- NameNode → Manages metadata & file system namespace. It keeps on tracking.
- DataNode → Daemons that manage storage on individual nodes.
 - and w request from a client

2nd NM → Merging namespace image & edit log to reduce load on NN.

Installation :-

- Prerequisites : Hadoop cluster installed & running
Zookeepers ensemble configured & running
- Steps :-
 - Download & Extract HBase
 - Configure HBase
 - Start HBase

Operations :-

- Creating a Table - Create
 - Listing Tables - list
 - Inserting data - put
 - Retrieving data - get
 - Scanning data - scan
 - Deleting data - delete
 - Dropping table - disable
drop
- ↳ counting rows
- Count

HBase commands :-

- hbase shell - Starting the shell
- listing commands - help

Java Client API for HBase :-

- Establishing connection - HBaseConfiguration.createC()
• createConnection()

CRUD operations -

- Create : • getTable() , • addColumn()
- Read : • get() , • getValues()
- Update : HBase Overwrite the value
- Delete : • add column() , • delete()

- Admin API :-
- Creating : - getAdmin(), - addFamily(), - createTable()
 - Listing tables : - listTables()
 - Disabling & Dropping Tables : - disableTable
- deleteTable
- ↳ HBase - Scan, Count, Truncate
(scan , count , truncate)

HBase Security

Authentication : kerberos

Authorization : Access Control (Defines permissions)

Encryption : Data encryption (Encrypt using HDFS zones)
Network encryption (Encrypt data in transit)

↳ Robust support for storing & processing large datasets with real time o/w access.

HIVE

Features : SQL like Query Language, schema on Read, Integration with Hadoop, support for diff.. storage types -

Components : UI, Metastore, HiveQL process Engine, Execution Engine, HDFS

Installation :

- Prerequisite : Java JDK, Hadoop, SSH

Datatypes : string, int, bigint, Boolean, float, double, timestamp, Date

Complex types : Array, Map,
Struct, Union Type

String functions : CONCAT(), LENGTH(), UPPER(), LOWER()
Date functions : CURRENT_DATE(), DATE_ADD(), DATE_SUB()
Hive tables : Managed tables , external tables
(manages data) (manages meta data)
(metadata) (only)

Partitions & Buckets :-

Partitions → Divides ~~table~~ tables into parts based on values of particular column.

Buckets → Divided into more manageable parts
Used for more efficient sampling & joining.

Storage formats :- Text file, Sequence file,

- Text files : Default storage format (text-table)
- Sequence files : seq-table
- ORC : Optimized Row Columnar
- Parquet : parquet-table

Importing data :- LOAD DATA command by creating external tables.

Querying Data :-

→ Sorting data : ORDER BY, SORT BY, DISTRIBUTE BY

→ Aggregating data : COUNT, SUM, AVG, MIN, MAX

Map Reduce scripts :- Hive can translate HiveQL into MapReduce jobs.

Joins & Subqueries :- SELECT, FROM, WHERE, INNER JOIN, SUBQUERY

Views are virtual tables creating using the result of a query.

Map and Reduce Side joins to Optimize Query

Map side and Reduce side joins to optimize Query.

map side join : smaller dataset fit into memory.

Reduce side join : Default join in HIVE, suitable for large datasets.

Data manipulation with HIVE :-

Inserting data, Updating data, deleting data

UDFs (User Defined Functions) :-

↳ Creates custom functions to extend HiveQL capabilities -
Can be written in Java & other languages.

- Creates Java class extending UDF.

- Implement the evaluate method

- Compile & add the JAR to HIVE.

Custom mapReduce in HIVE

↳ Scripts can be integrated into HIVE queries
(TRANSFORM).

Writing HQL scripts :-

↳ HQL allows for the batch processing of HIVE queries.

HQL script is a text file containing HIVEQL commands.

For executing Hive : `hive -f script.hql`

Hadoop :- Open source, to handle large datasets - Allows

for storage & processing of big data.

components : HDFS, YARN, MapReduce

scalable & reliable data storage designed

HDFS :- Provides cluster of commodity servers

to span large clusters of commodity servers (HDFS daemons)

Architecture : NNode, DNode, 2^o Node, NameNode & DataNode

HBASE :- Distributed, scalable, large amount of space

Stores and manage HBase master, Region Servers, Zookeeper

Architecture : HFile

Hive :- Data ware house built on top of Hadoop.
— Data summarises, querying & analysis.

Architecture -

Metastore

Hive server

Hive driver

Execution engine

↳ Hadoop : Provides distributed storage (HDFS) &
processing data (Map Reduce)

HDFS : Handles large datasets with high scalability

HBase : NoSQL DB built on top of HDFS, real time
read/write access to large dataset

Hive : Data warehousing that provides SQL like interface for
querying & analyzing data stored in Hadoop.

Apache Spark :-

↳ Open source, distributed computing system. Provides interface
for programming entire cluster with implicit data
parallelism & fault tolerance.

→ Cluster with implicit data parallelism & fault
tolerance

→ To overcome limitations of Hadoop MapReduce framework
and provide faster, more efficient data processing.

Linking with Spark →

↳ Hadoop

Apache Hive

Apache HBase

Apache Kafka

Initializing spark :-

↳ `SparkSession.builder()`

RDD (Resilient Distributed Datasets) :-

↳ Properties : (Immutable, Fault tolerance, Parallel processing) -

RDDs fundamental data structures of Spark.

Creating RDDs :- • `parallelize(data)`

↳ RDDs : Low Level APIs provide transformation actions.

↳ `PySparkSQL` import `SparkSession`

`sparkSession.builder().appName("None").getOrCreate()`

`sc = spark.sparkContext`

`data = [,,,]`

`rdd = sc.parallelize(data)`

`rdd = sc.textfile(" --- ")`

Shuffle Operations :-

↳ `reduceByKey` & `groupByKey`

RDD persistence :-

↳ It can be persisted in memory or disk to avoid recomputation.

Removing data :- To unpersist data [`rdd.unpersist()`]

Shared Variables :-

- Broadcast Variable : variable read only value to all nodes workers & drivers.

- Accumulator : variables added to from read only by

Deploying to a cluster :-

- standalone mode : spark manages own cluster
Spark appl. on Hadoop cluster

- YARN : cluster manager used to run spark.

- Mesos :

Map Reduce with Spark :-

Spark [with hadoop] → Hadoop HDFS for storage
YARN for resource management
[without Hadoop] → Use local file systems or other storage solutions like S3.
Spark's standalone cluster manager.

↳ Storage :- with Hadoop - HDFS

without Hadoop - Local file system or cloud storage

Reason :- with Hadoop - YARN

Management :- without Hadoop - Standalone or Mesos

Data Preprocessing :- Step in data analysis pipeline.

Transforms raw data into a clean & usable format.

Steps → Data cleaning

Data Transformation

Data Reduction [(PCA), Sampling]

Feature Engineering

EDA (Exploratory Data Analysis) :-

↳ Analyzing data sets to summarize their main characteristics, often using visual methods.

Understanding the data & uncovering patterns, anomalies, relationships.

- Descriptive Statistics [mean, median, mode, var, SD]

- Data visualization [Plots]

- Correlation Analysis [Relationship b/w variables using correlation Coefficients]

- Outlier detection [Identifying the outliers in data]

Apache Kafka :-

- Designed for high throughput, low latency in data streaming.
For build a real time data pipelines & streaming applications.

Spark Streaming Architecture :-

- Discretized Streams (DStreams) [Spark streams represent streams of data → Series of RDDs]
 - Receivers [To ingest data from sources like kafka, Flume & others]
 - Processing [Each batch data processed by Spark jobs, can perform transformations & actions on data]
 - Output Operation [Processed data can be written to various output sinks (HDFS, DB, ...)]
- ↳ Apache spark & kafka for building scalable, real time data processing pipelines - combining capabilities of Spark for processing

Kafka Connect API :-

- Source Connectors (Import data $\xrightarrow{\text{from}}$ Ext.. syst into kafka)
 - Sink Connectors (kafka $\xrightarrow{\text{export}}$ Ext.. systems)
- ↳ Apache Kafka :- written in java & scala.
- Messaging, Storage, Processing
 - Topics, Partitions,
 - ↳ Producers & Consumers
 - ↳ Brokers
 - ↳ Zoo Keepers

Apache Spark

21/7/24

Apache kafka :-

→ Open source stream processing platform . Written in Scala

Java .

Functionalities : Messaging , storage , processing
Categories where records are stored

Concepts : Topics -

Segments of topics that allow parallel processing

Partitions -

Data written to Kafka topics

Producers -

Data Read from Kafka topics

Consumers -

Kafka msg of brokers handles storage &

Brokers -

retrieval of records .

Zookeepers -

Manages & coordinates Kafka brokers .

Apache Spark → Open source unified analytics engine

for large scale data processing .

Provides high level API in java , scala , Python , R

Components : Spark Core , Spark SQL , spark streaming , MLlib ,

GraphX

Spark SQL :- Spark module for structural data processing .

It provides Dataframe API & SQL queries .

Features : DF & Datasets → DF (Data organized into named columns)

Dataset (Provides benefits of RDD with optimized adv of Spark SQL)

SQL Integration → Allows executing SQL queries along - side complex analytic algorithms .

Catalyst Optimizer → Query Optimizer for executing queries efficiently .

DataSource API → Allows integration with variety of data source (JSON , Parquet , Avro)

Spark MLlib :- ML lib built on top of Spark .

Features : Algo - classification , Regression , Clustering , collaborative filtering etc . -

Factorization - Factor extraction , E.T , dimension Reduction

Pipelines - Tools for constructing, evaluating & tuning ML

Utilities - Statistics, data handling & other utilities.

Predictive analysis with Spark :-

↳ Using historical data to predict future outcome

steps - Data Preprocessing (cleaning & transforming)

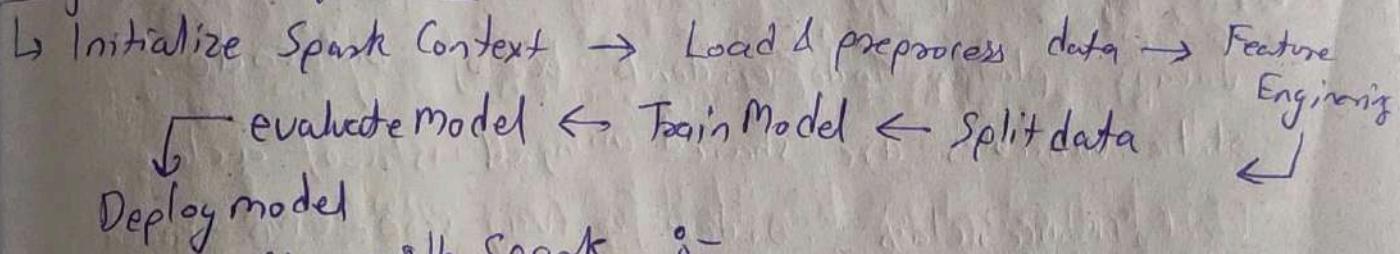
Feature Engineering (creating new feature from existing)

Model building (LR, DT) ^{data}

Model Evaluation (Metrics like accuracy, precision, recall)

Model deployment (Integrating model into production) Systems to make predictions

Workflow in Spark :-



Linking kafka with spark :-

kafka used as source and sink for spark streaming application.

Steps : - Setting up kafka

- Creating spark streaming Context

- Consuming data from kafka

- Processing data

- producing data to kafka

→ Apache Kafka & Apache Spark powerful tools for real time,

scalable, data processing pipelines.

→ Spark SQL and MLlib provides robust frameworks for data processing & ML, enabling predictive analysis.

→ Integrating kafka and spark, creates efficient streaming applications that handles large volumes of data with low latency.