

Principles of Big Data Management (CSE-417)

Unit-1

Big Data characteristics:

Volume, Velocity, Variety, Veracity, Value
Size Speed of Various Quality Usefulness
MB, GB data types of data
PB emails, tweets, etc.

- Structured, unstructured, semi-structured
- RDBMS Multimedia = web pages, .CSV

* Big Data:

branch
It is a field dedicated to the analysis, processing and storage of large ~~masses~~ collections of data that frequently originate from disparate sources.

~~structured data~~
~~unstructured data~~ = ~~text~~ ~~image~~ ~~video~~
~~semi-structured data~~

Analysis: It is a process of examining data in order to find the relationships, patterns, insights and trends among the data

Analytics: It includes development of analysis methods, scientific methods and automated tools used to manage the complete data life cycle

Goal of Data Analysis is to support better decision making

Different types of Data Analytics:

- descriptive: Events that have already occurred Eg: Sales volume ^{over past 12 months}
- diagnostic: Aims to determine cause of phenomenon in past ^{behaviour}
- predictive: Predict outcome (based on patterns found in historical ^{present})
- prescriptive: Actions should be taken after results of predictive, ^{Commission}

Descriptive Ex: Monthly ^{Commission} sales earned by sales agent

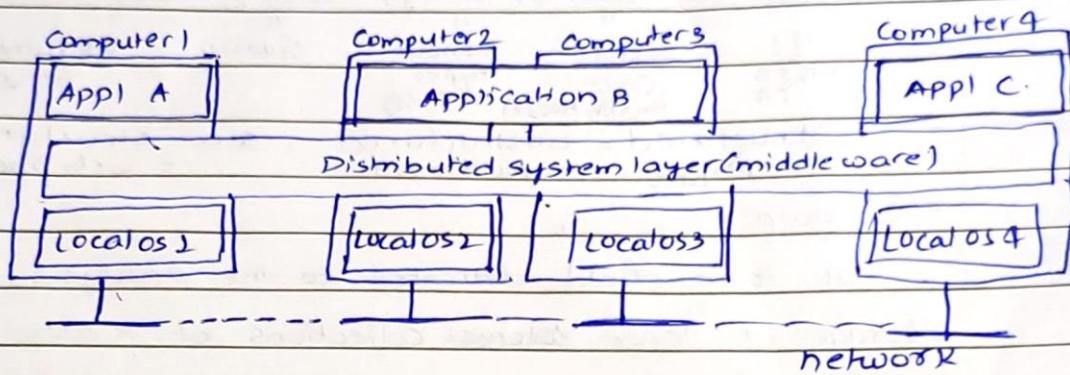
Diagnostic: Ex: Item 2 sales less than item 3 why?

Predictive: Patient survival rate if Drug B is administered instead of drug A

Prescriptive: Among 3 drugs which one provide best results

Distributed System (DS):

It is a collection of independent computers that appears to its users as a single coherent system
- computing power & storage is shared.



Advantages:

- scalability
- Reliability
- Availability
- Communication

Big Data Storage Concepts:

cluster systems:

- Collection of similar computer nodes connected by high speed LAN

- Nodes run on similar OS

These will be 2 nodes

1) Master Node

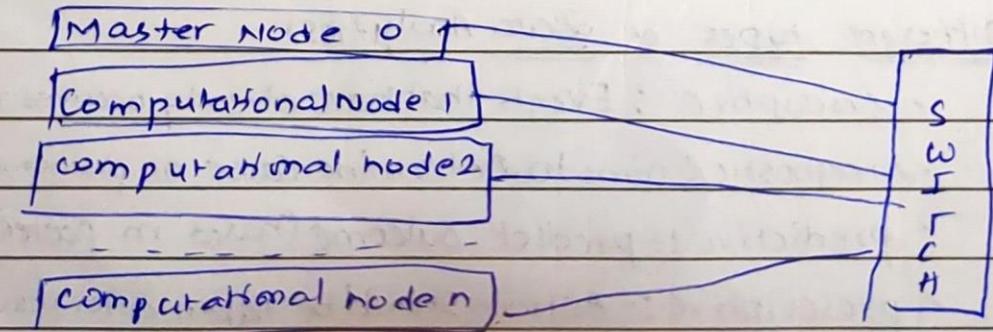


It controls storage

allocation, bandwidth

Job scheduling

- Computing power and storage are shared among users.



Filesystem:

- A File is named collection of related info that is recorded on secondary storage (Hard drives, DVD's)
- A File consists of Data or Program two parts

i) Collection of files

For storing data & programs

ii) Directory structure

Provides info abt all files in the system

- FS is responsible for organization, storage, naming, retrieval, sharing & protection of files.

Centralized File System:

FS which stores files & directories in a single computer system is known as CFS.

Disadvantages:

Distributed File System (DFS):

- FS that manages the storage across a network of machines is called DFS.
- It is a client/server based application that allows client to access and process data stored on the server.
- A DFS has to provide single system image to the clients even though data is stored in multiple computer systems.

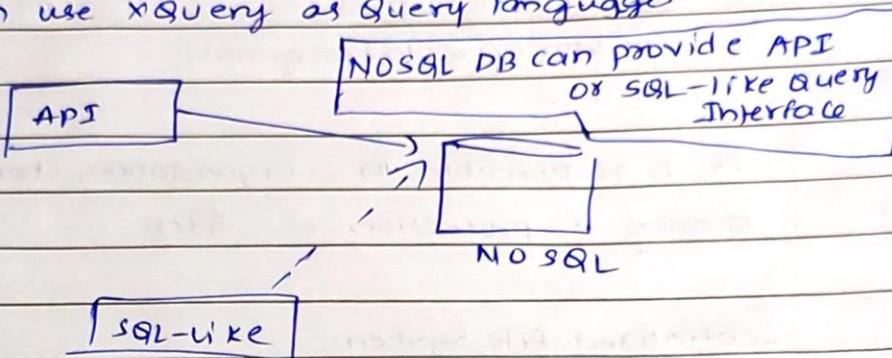
Advantages:

- Network Transparency: client doesn't know abt location of files. clients can access files from any computer connected in network.
- Reliability: keeping multiple copies for file.
- Performance: parallel access of data is possible.
- Scalable: can add multiple storage to computer systems in network

NOSQL (Not only SQL)

- It's a non-relational Database that is highly scalable, specifically designed to store structured & unstructured data.

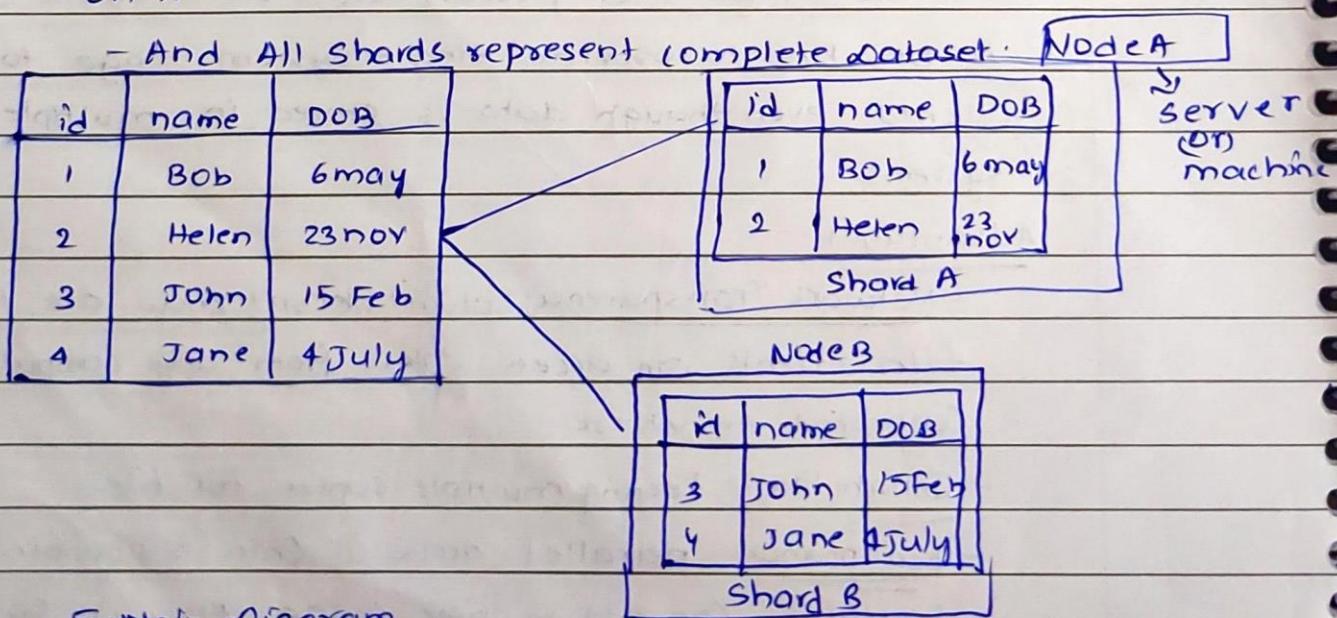
A NOSQL DB is optimized to store XML files with often use XQuery as query language.



Sharding:

It is process of horizontally partitioning a large dataset into smaller, more manageable datasets called "shards". These shards are distributed across multiple nodes where a node is a server.

- Each shard is stored on a separate node.
- And each node is responsible for only the data stored on it.
- And All shards represent complete dataset.



Explain diagram

Ex:

Advantages: In case of a node failure, only data stored on that node is affected.

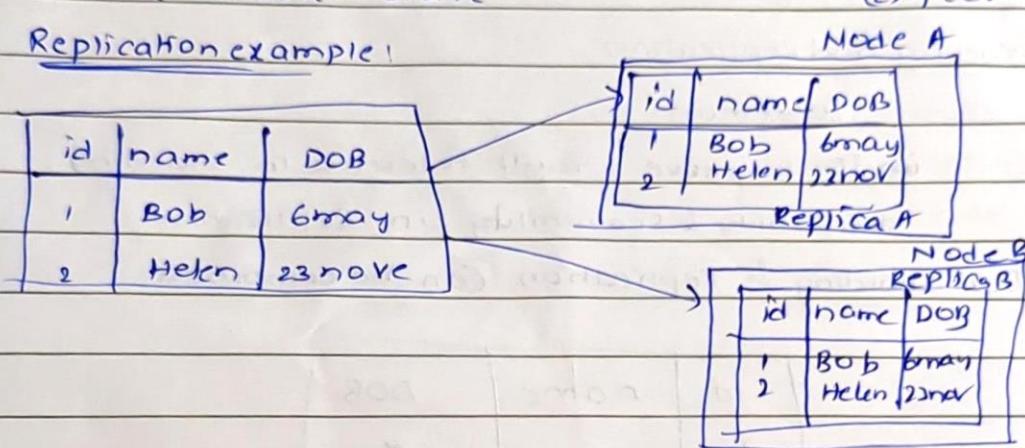
Replication: It stores multiple copies of a dataset, known as replication on multiple nodes.

- It provides scalability and availability, due to the fact that the same data is replicated on various nodes.
- Fault tolerance is also achieved since data redundancy ensures that data is not lost when an individual node fails.

There are two methods that are used for replication:

(1) Master-slave

Replication example:

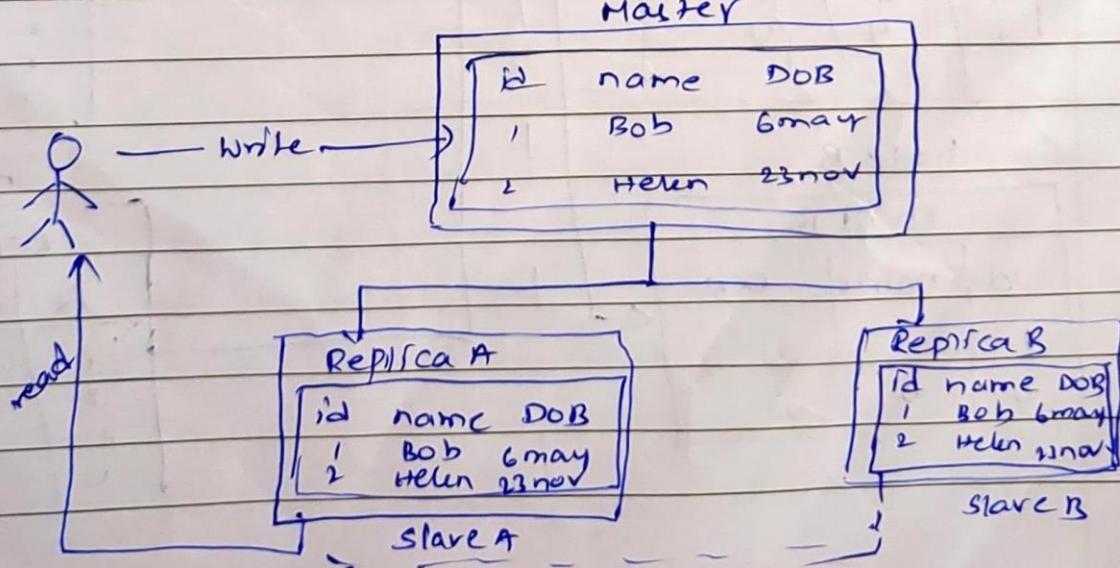


(2) peer-peer

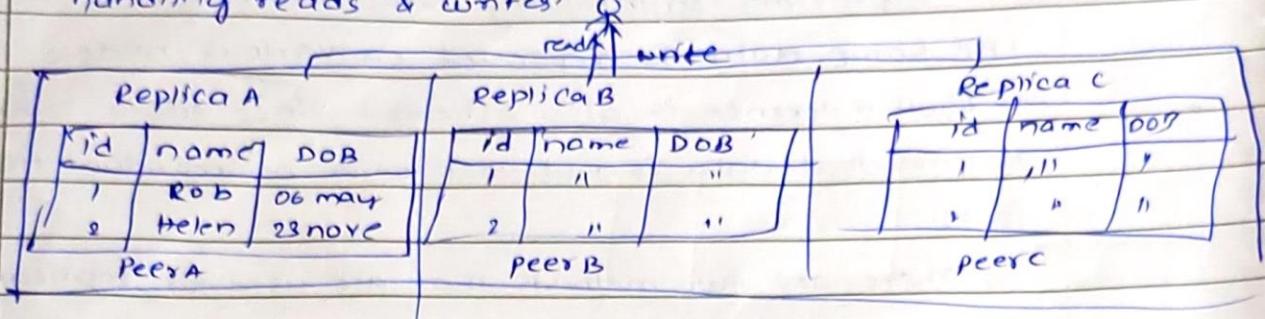
(1) Master-slave

Here nodes are arranged in master slave configuration and all data is written in master node

- Then data is replicated to multiple slave nodes
- All write requests ^{Including} ~~are~~ insert, update and delete occur in Master node
- All read requests are fulfilled in Slave node.



Peer-peer: All nodes operate at same level (no master & slave)
 Each node is a peer and equally capable of handling reads & writes.

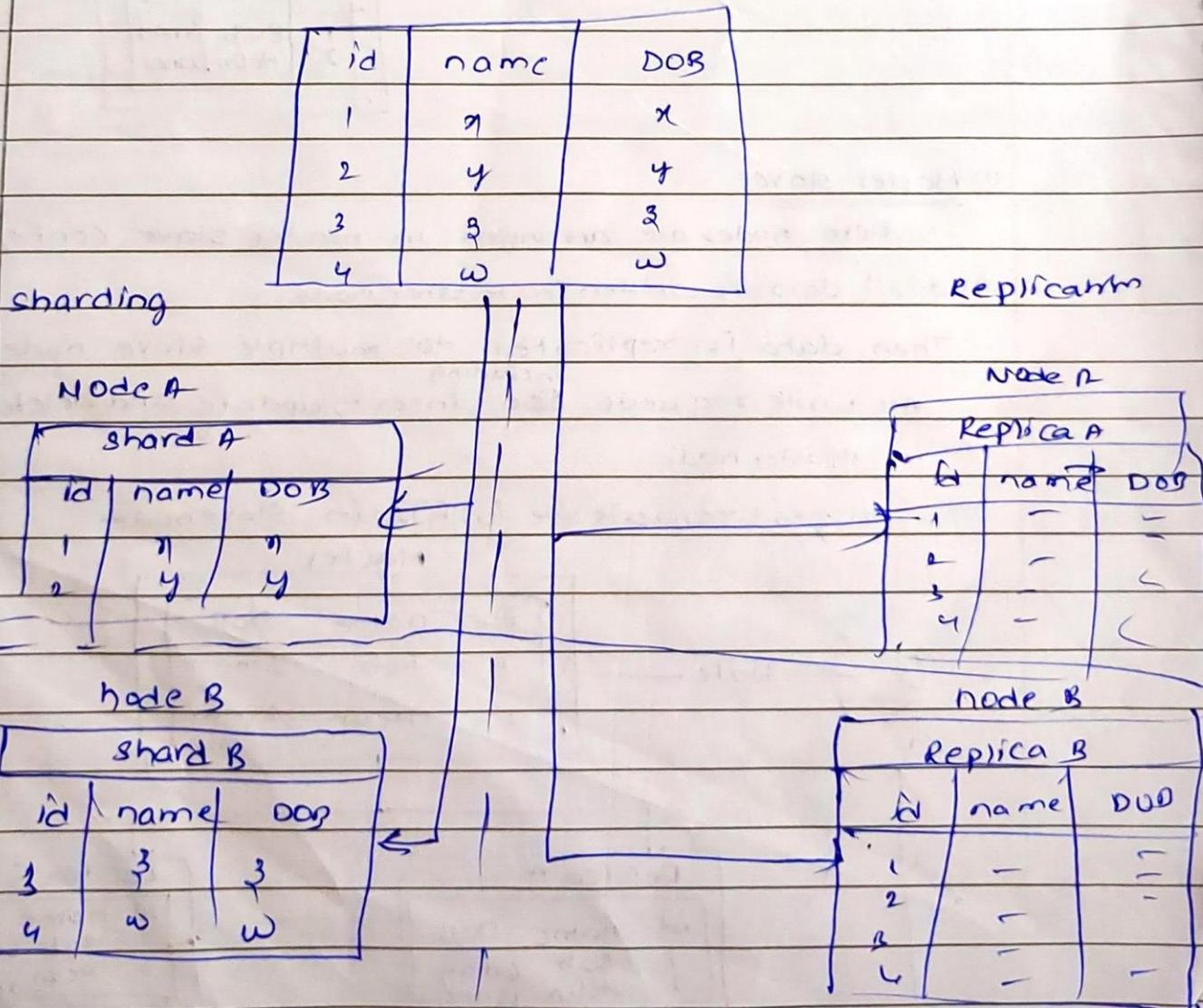


Sharding and Replication

To remove

TO improve fault tolerance in sharding,
 availability & scalability in replication.

Both sharding & replication can be combined.



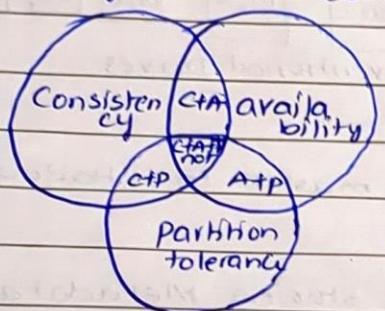
CAP Theorem / Brewer's Theorem:-

Consistency, Availability, and Partition Tolerance

Consistency: A read ~~reading~~ from any node results in same data across multiple nodes diagram

Availability: A r/w request will always be "acknowledged" in the form of a success or failure.

Partition tolerance: The database system can tolerate communication outages that split cluster into multiple silos and can still serve read/write requests.
In any communication outages and still persist r/w



cap theorem states that a distributed database system can only provide two of three properties.

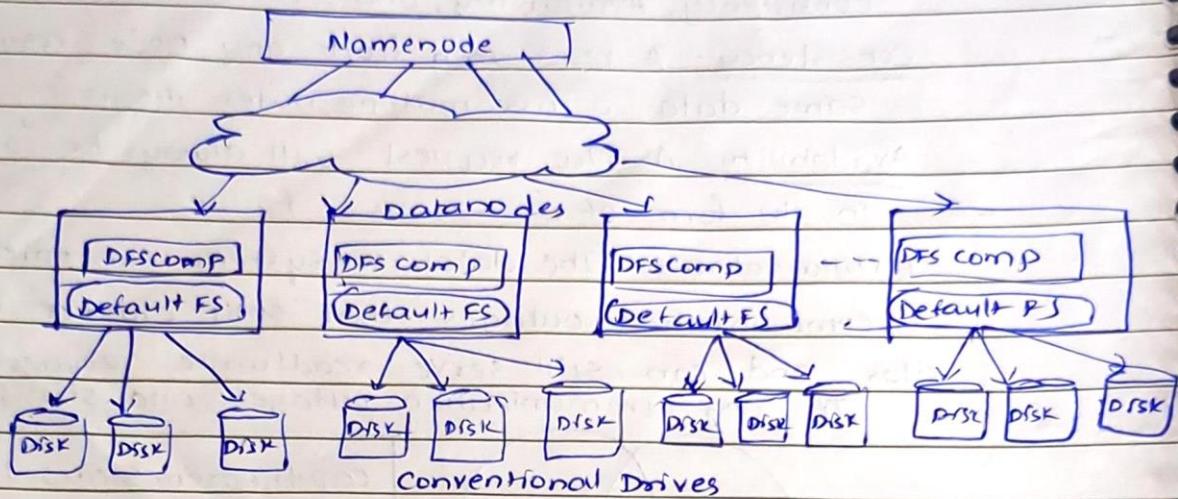
BASE:- It is a database design principle based on the CAP theorem

BASE = basically available, soft state, eventual consistency
Basically Available: This property refers to fact the database will acknowledge client's request in the form of success or failure notification.

soft state: This property refers to the fact that the state of database can change overtime, even without any explicit user intervention. This happens due to the effects of background processes and updates to data.

Eventual consistency:- This property refers to the fact that it is the state in which reads by different clients, immediately following a write to the database, may not return consistent results

Hadoop Distributed filesystem:

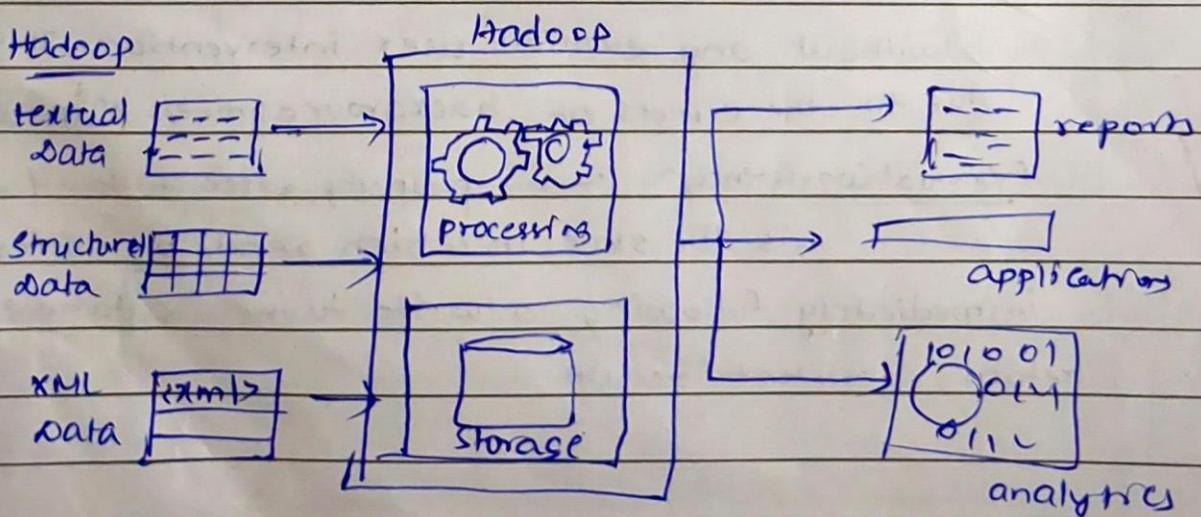


Name Node (Master)

- It works as a master in a Hadoop cluster that guides datanodes.
- Mainly used for storing Metadata (location of file, size of file etc.)

Datanode (Slave)

- works as a slave mainly used for storing Data in a hadoop cluster
- runs HDFS Client program
- Executes Map & Reduce tasks

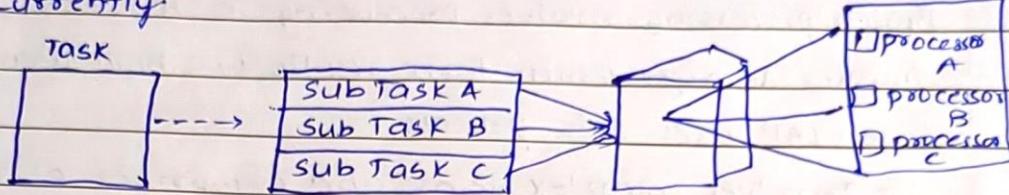


Unit - II

Big Data Processing Concepts

Parallel Data Processing

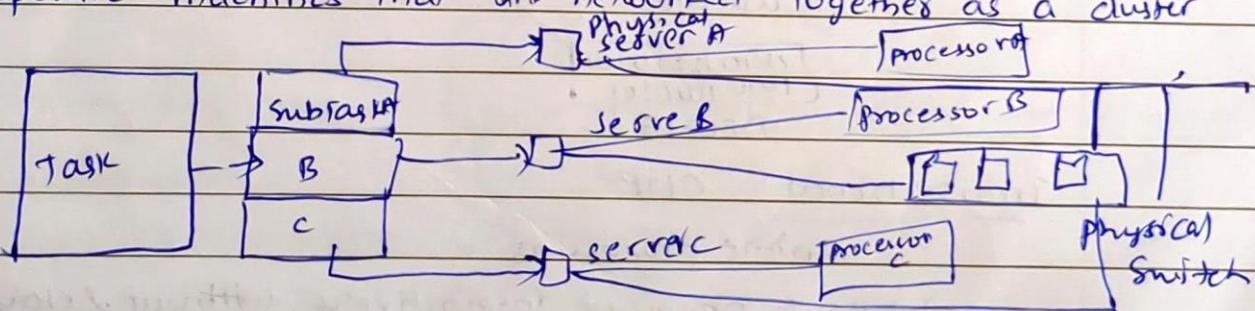
- It involves simultaneous execution of multiple tasks that collectively comprise a larger task.
- The goal is to reduce the execution time by dividing single larger task into multiple smaller tasks that run concurrently.



A task can be divided into 3 subtasks that are executed in parallel on 3 diff processors within the same machine

Distributed data processing

- It is also closely related to parallel data processing in that the same principle of "Divide & Conquer" is applied
- However, distributed is always achieved through physically separate machines that are networked together as a cluster



Hadoop

It's a open source framework for large scale data storage.

- It is versatile framework that provides both processing and storage capabilities.

Processing Workloads:

It is the amount & nature of data that is processed within a certain amount of time.

→ Two types

(1) Batch / offline processing

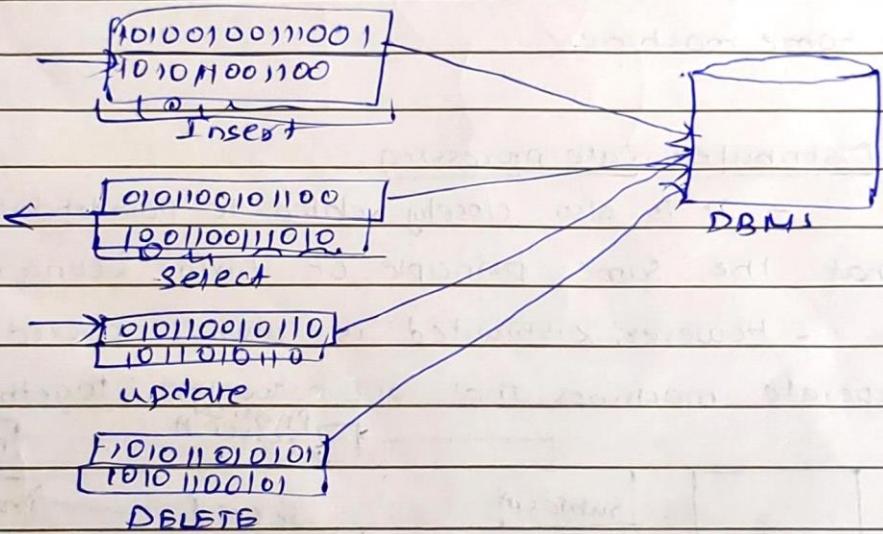


(2) Transactional

Batch processing involves processing in batches and usually imposes delays which turn results in high latency responses.

→ OLAP uses this process.

→ Involves complex joins and comprises of group of two queries.



Transactional: OLTP

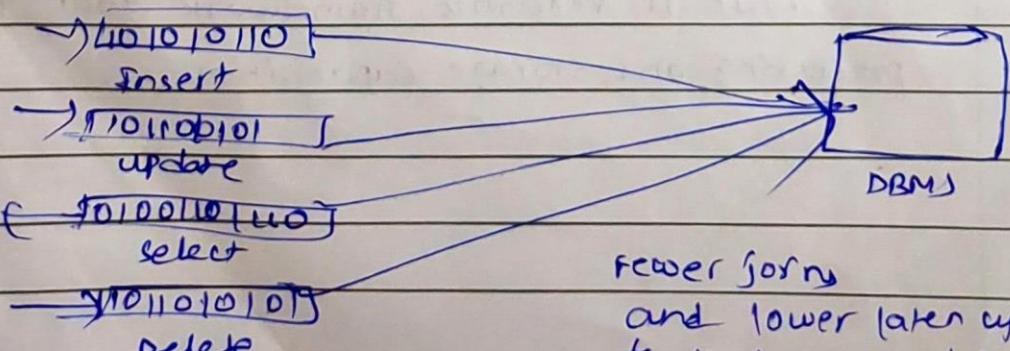
online processing

→ Data is processed interactively without delay.

→ low latency responses

→ small amt of data with random rand w

→ fewer joins



fewer joins

and lower latency responses than batch workloads

- Majority of Big Data processing occurs in Batch Node.
- Predictive, prescriptive in Batch processing.

Batch processing with MapReduce:

- Map Reduce is a widely used Batch processing framework.
- It is highly scalable, Reliable and it is based on divide & conquer rule principle; which provides built-in fault tolerance and redundancy.
- Map Reduce can process schema less datasets.

Map and Reduce Tasks:

MapReduce:

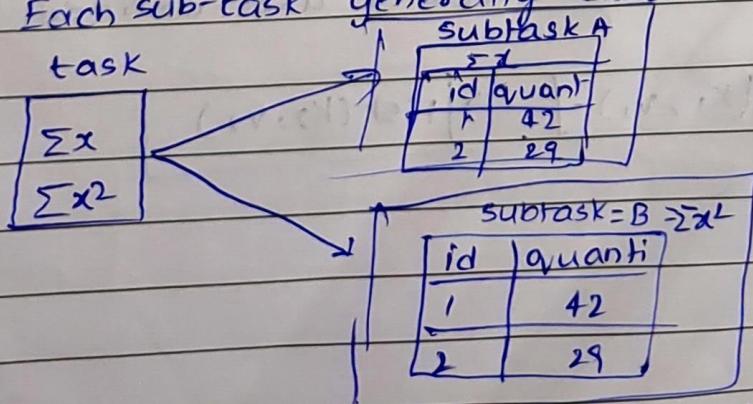
- programming model for data processing
- it is based on divide and conquer principle
- it divides big problem into smaller problem
- A dataset is broken into multiple part and operations are performed on each part independently in parallel.

The divide & conquer is generally achieved using one of following approaches.

(i) Task parallelism: (Dividing tasks)

- Dividing a task into subtasks and running each sub-task on a separate node.

- Each sub-task generally executes a different algorithm.

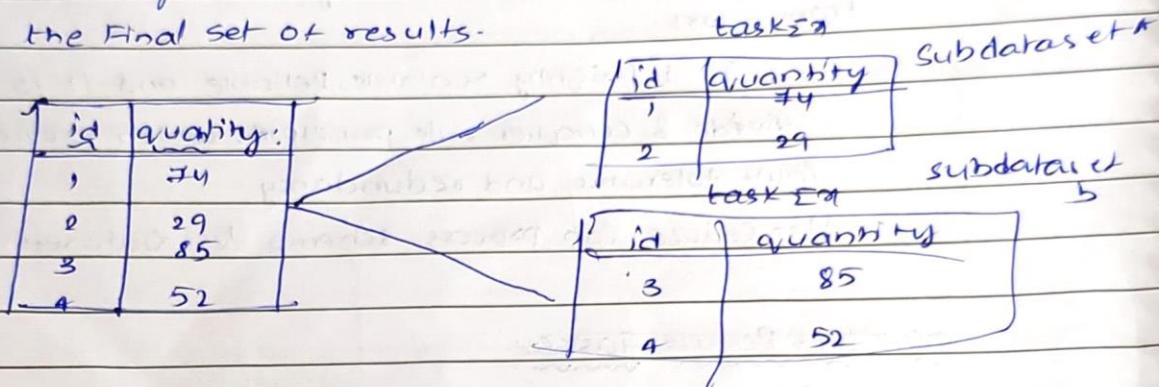


task div in to sub tasks and run on 2 different nodes on same data.

- Finally multiple subtasks is joined together to obtain final result.

Data parallelism: (Dividing dataset)

- Dividing dataset into multiple datasets and processing each subdataset in parallel.
- All are processed using same algo.
- Finally subdatasets are joined together to obtain the final set of results.



Map and Reduce Tasks:

Map tasks:

- map
- combine (optional)
- partition

Reduce tasks

- shuffle & sort
- Reduce

Map:

→ The first stage is Map, here dataset file is divided into multiple smaller splits

→ Converted into key-value pair

Key - ordinal position of record ; Value - Actual record

→ Then these key-value pairs for each split are then sent to a map funcn.

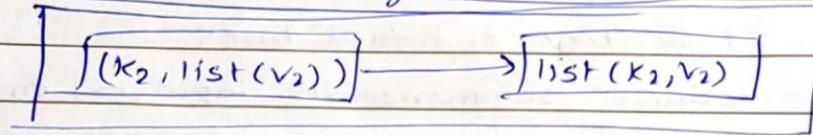
→ map function executes user-defined logic.

$$(K_1, V_1) \rightarrow \text{list}(K_2, V_2)$$

Combine:

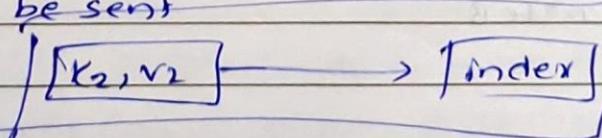
Generally, output of the map funcn is handled directly by reduce funcn

- However, Map and Reduce run on different nodes, This requires moving data btw Mappers and Reducers. This requires lot of bandwidth and processing latency
- In large datasets, the time taken to move the data btw Mapper and Reduce will be excep^{tive}
- Due to this reason MapReduce engine provides an optional combine funcn. that summarises mapped output before it gets processed by the reducer.



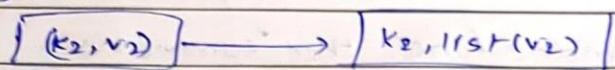
Partition:

- In this stage, if more than one reducer is involved a partitioner divides the output from the mapper or combiner into partition between reducer instances
- NO. OF partitions will be equal to number of reducers
- It helps for distribution of key-value pairs.
- It is last stage of map task
- It returns index to which particular partition should be sent



Shuffle & sort:

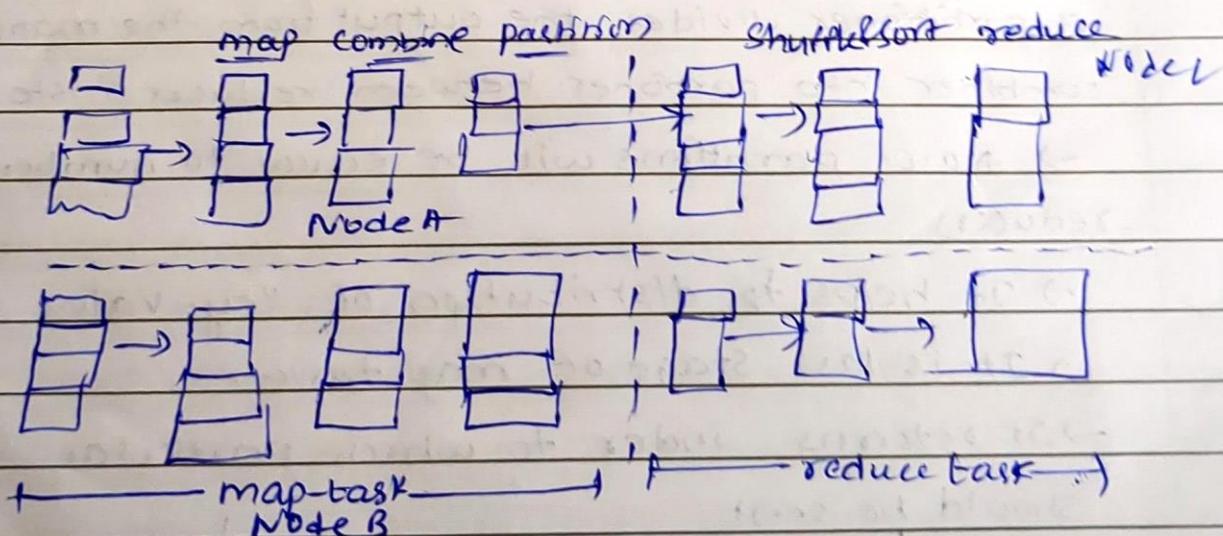
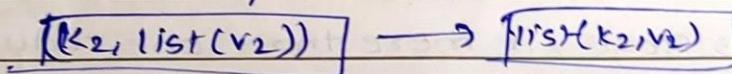
- First stage of reducer task.
- During this stage output from all partitioners is copied across network to reducer nodes, this is known as shuffling
- MapReduce engine automatically groups and sorts the K-V pairs according to keys, so that the output contains a sorted list of all input keys and their values.



Reduce:

Final stage of Reducer task.

- reducer - summarize its input (or) emit the output.
- The output of reducer i.e; key-value pairs is returned as a separate file.



Unit-4

A Customer has placed many orders

Customer					order			
<u>a</u>					<u>b</u>			
Cid	Cname	Cage	Caddr	Csal	Oid	Odate	Oamt	Cid
101	xyz	30	HYD	25K	201	10JUN	5300	101
102	PGR	45	BAN	30K	202	21 APR	2400	102
103	ABC	32	DEL	48K	203	3 Aug	3841	103
104	DEF	27	MUM	52K	204	5 Sep	2728	104

→ hdfs dfs -get /hive/joins/customer.csv

→ hdfs dfs -get /hive/joins/order.csv

hive

hive > create table customer_join (Cid INT, Cname STRING,
Cage INT, Caddr STRING, Csal Float)

ROW FORMAT DELIMITED FIELDS TERMINATED BY ',';

hive > LOAD DATA LOCAL INPATH 'customer.csv' INTO TABLE
customer_join;

hive > create table order_join (Oid INT, Odate STRING, Oamt Float
(Cid INT))

ROW FORMAT DELIMITED FIELDS TERMINATED ',';

hive > LOAD DATA LOCAL INPATH 'order.csv' INTO TABLE
order_join;

hive > SELECT a.Cid, a.Cname, a.Cage, a.Caddr, a.Csal, b.Odate
From customer_join a LEFTOUTERJOIN order_join b
ON (a.Cid = b.Cid)

hive > SELECT a.Cid, a.Cname, a.Cage, b.Odate FROM
customer_join a JOIN order_join b ON (a.Cid = b.Cid);

sum of squares

Mapper :

function mapper(input-list): ✓

output-list = []

 for index, number in enumerate(input-list)

Square = number * number

output-list.append((index+1, square))

 return output-list

reducer

function reducer(input-list): ✓

sum-of-squares = 0

 for key, value in input-list

Sum-of-square + = value

 return sum-of-squares