Unit 2

# RDBMS vs NO-SQL

| Relational Database | NO-SQL Database |
|---|---|
| (1) Relational database supports a powerful query language | ① NO SQL database supports a very simple query language |
| ② Relational Database has a fixed schema | ② NO-SQL database has no fixed schema |
| ③ Relational database follows acid properties (Atomicity, Consistancy, Isolation & Durability) | ③ NO-SQL database is only eventually consisten. |
| ④ Relational database supports transactions (also complex transaction with joins) | ④ NOSQL database do not support transactions (support only simple transactions. |
| ⑤ RDMS manages only Structured data | ⑤ NO-SQL database can manage Structured Unstructured & Semistructured data |
| ⑥ Relational database has Centralized structure | ⑥ NO-SQL databases can handle big data or data in a very high volume |
| ⑦ Relational database are used to handle moderate volume of data | ⑥ NO SQL database has decentralized structure |
| ⑧ Relational Database have a single point of failure with fail our | ⑦ NO-SQL database have no single point of failure |

# Types of NO-SQL Database (Architecture pattern)

① Key Value Stores :- The main idea is using a hash table where there is a unique key & a pointer to a particular item of data.

But it is insufficient when you are only interested in giving or updating part of a value, among other disadvantages. Key-value pair storage database stores data as a hash table where each key is unique & the value can be a Json, BLob (Binary large objects) string etc.

Eg

| Key | value |
|-----|-------|
| Name | Preeti |
| Birthday | 21/12/2001 |

(2) Column family stores :- These were created to store and process very large amounts of data distributed over many machines there are still keys but they point to multiple columns. The columns are corranged by column family.

eg Cassandra, HBase, HyperTable

these database are mainly used to manage data warehouse business intelligence, CRM, library card catalogs.

(3) Document Database :- These were inspired by lotus notes & are similar to key value stores. The model is basically versioned documents that are collection of other key-value collections. The semi-structured documents are stored in formats like Json.
Document database support querying more efficiently.

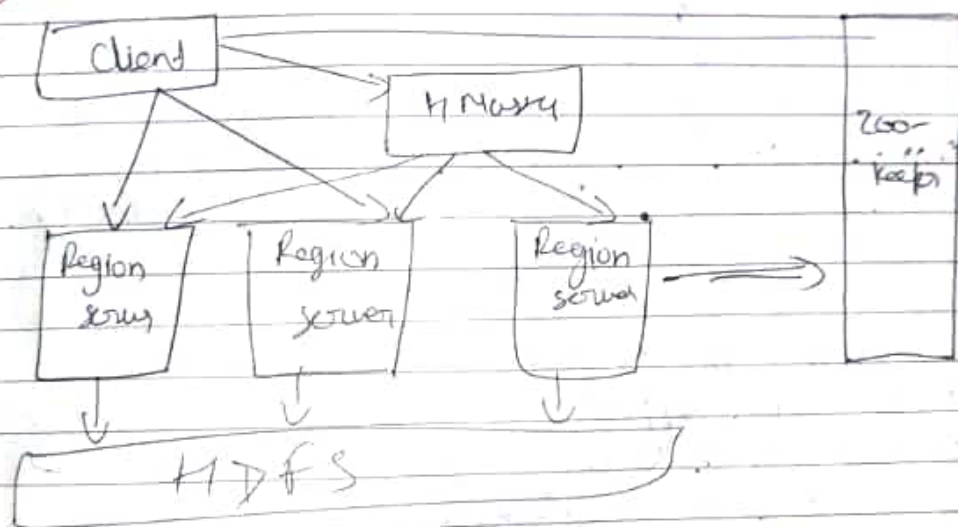Eg CouchDB, Mango DB, Amazon, Simple DB, Riack
lotus notes.

(4) Graph Databases :- Instead of tables of rows & columns & the rigid structure of SQL, a flexible graph model is used which again can scale across multiple machine. Graph base databases mostly used for social network, logistics, spatial data.

Eg Neo4J, Infolnaid, Infinite Graph etc.

⑤ Multit model. Database ? These are designed to handle multiple data models against a single integrated backend. they are a brand-new in the No-SQL world, there will be much more buzz around this type of database in the future.

Ans OrientDB is a multi-model database, combining No-SQL types, Orient DB is graph database, where each node is a document.

# HBase Architecture

H-Box Architecture has 3 Components, HMaster, Region Server, Zookeeper.

① HMaster ⇒ the implementation of master server in HBase is HMaster. It is a process in which regions are assigned to region server as well as DDL (create, delete table) operations It monitor all Region server instances present in the cluster.
It has many features like controlling, load blancing, failover etc.

(2) Region server :- HBase Tables are divide horizontally by row key range into regions. Regions are the basic building elements of HBase Cluster that consists of the distribution of tables & are Comprised of Colourn families Region server runs on HDFS data Node which is present in Hadoop cluster. The default size of Region is 256MB

(3) Zookeeper :- It is like a coordinator in HBase :-
It provides services like maintaining configuration information, naming, providing distributed synchronization, server failure notification etc.
Clients communication with region servers via Zookeeper.

# Advantage of HBase
1 Can Store large date sets
2 Database can be shared.
3 Cost effective from gigabytes to petabytes.
4 High availability through failover & replication.

# Disadvantages of HBase
1 No Support SQL structure
2 No transaction support
3 Sotored only on key.
4 memory Issues on the Cluster.

# CAP Theorem (also Called Brower's Theorem)

Network partition means there is a break in the network: or there is no connection b/w the data stores so the scenario where there is a break b/w different nodes in a distributed data stores is called partition tolerance. It is nothing but no data replication for a particular node or maybe group of nodes or we can call it isolated data so that any database or data store is isolated from other databases in the distributed system because these is partition tolerance.

- Consistency = means all Clients see the same data at the same time no matter which node they connect to for this is to happen whenever data is written to one node it must be instatly provided or replicated to all the other node in the system before the write is derived successfully.

- Availabilty = means that any Client making a sequence for data gets a response even if one or more node are down another way to state this all working nodes in the distributed system return a valid response for any request without any exception.

- Partition Tolerance = Acc. to cap theorem partition tolerance is must so when a network fails we have to choose Consistency or availability as a distributed database System is bound to have partition in a real world system due to network failure or some other reason therefore partition tolerance is a property we can not avoid while building our system. So a distributed system will either choose to give up consistency or availability but not partition tolerance

Eg In a distributed system if a partition occurs b/w 2 nodes It is impossible to provide consistent data on both the nodes & availability of complete data therefore in such

a scenario we either choose to compromise on Consistency
or on availability.
Distributed database is either Characterised as:-
                CP or AP

# HDFS (Hadoop Distributed file system)

It is used for storing massive data.
The system contains memory store of h files, memory store
is fast like the cache memory anything that is stored
onto the h base is stored here initially later the data
is transfered & saved in h files as blocks of the memstore
is flushed.

HDFS - It is the component of Hadoop which is responsible
for storing the data it has 2 components.

① Name Code : as data is segmented into various blocks
so name node stores the meta data of the various blocks
the info which this meta data holds is the name of the
blocks, size of the block, location of the block of the
original file.

② Data node : It holds the actual data i.e. block

In map Reduce layer — it has 2 components

④ Job Tracker : Divides the computer into various task. It assigns
these tasks to the various task trackers. once the computation
is performed by the various task trackers. the result is
accumulated & Output is generated by the job tracker.
It also takes the feedback from the various task trackers
as well as it monitors the status of all these task
trackers.

③ Task Tracker : It performs the actual computations.

A features of Hadoop.

① fault tolerant : as the data is stored in various nodes of the cluster & a replicated copy of the original block is also stored somewhere in the clusture so if a particular system crashes then the data is not lost.

② Scalable : It means that no of system can be increased or decreased on the basis of requirement of the application

③ low Cost :- As hadoop is open source i.e it is freely available & it is build on low cost comodity hardware which makes the entire cost very low.

④ Can store unstructured data.

⑤ faster Computation :- as hadoop distributes processing among the various node in the clusters i.e why it performing faster computation.