

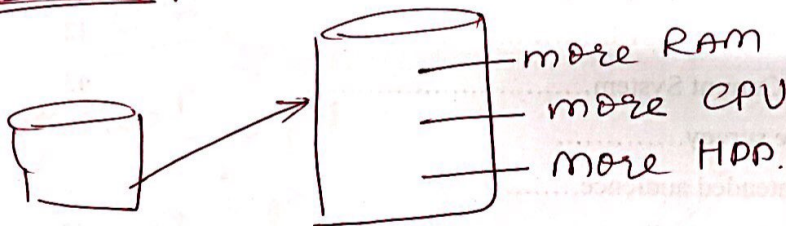
* NoSQL Database

- Non-relational DBMS, that does not require a fixed schema.
- avoids joins + easy to scale.
- NoSQL is used for Big Data + real time web Applications.
- "Not only SQL" or "Not SQL"

* Why NoSQL?

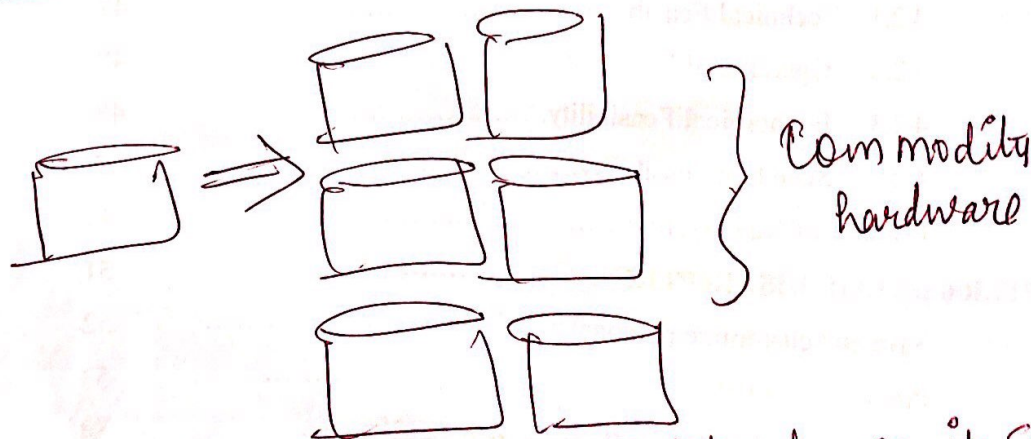
- The system response time becomes slow when you use RDBMS for massive volumes of data.

- Scale-up - with RDBMS (Vertical Scaling)



- expensive process.

- Scale-out - NoSQL - (horizontal scaling)



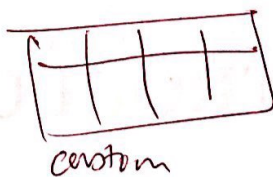
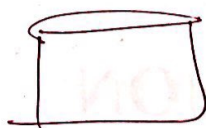
- NoSQL database is non-relational, so it scales out better than relational Databases.

* Features of NoSQL:

- ① Non-relational.
- ② Schema-free - or relaxed schema
Implicit Schema

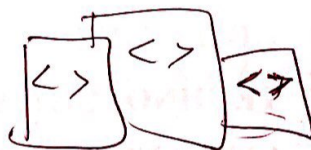
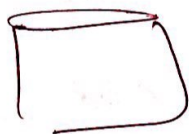
RDBMS:

select Name, age from customers



NoSQL DB:

Item [price]
Item [Discount]

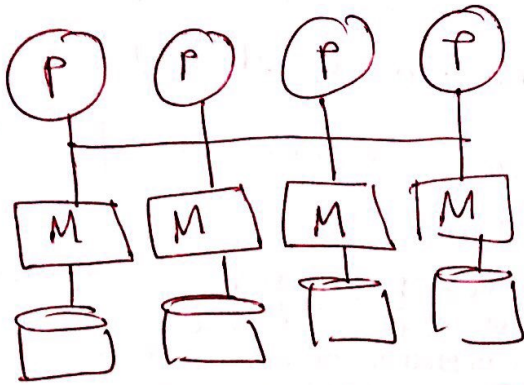


③ Simple API

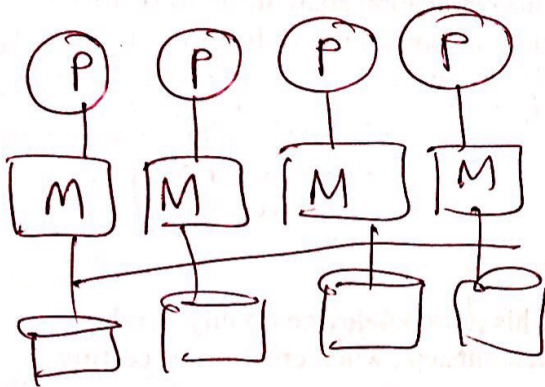
④ Distributed:

- auto scaling & fail-over capabilities
- often ACID concept can be sacrificed for scalability & throughput.
- Mostly no synchronous replication between distributed nodes
- only provide eventual consistency.
- Shared nothing Architecture.
enables less coordination & higher distribution.

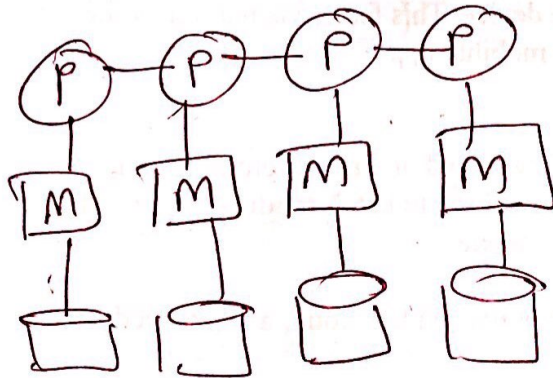
Oracle 11g - shared memory



Oracle RAC - shared ~~RAE~~ Disk



Nosql - shared nothing



* Types of NOSQL Database:

- ① Key-value
- ② Column-oriented
- ③ Graph based
- ④ Document orient

① Key-value:

- Data is stored in key-value pair.
- It is designed in a such a way to handle lots of data & heavy load.
- Store data as hash table where each key is unique, and the value can be JSON, BLOB, string etc.

→ Website → Domain name

Key	Value
Name	
Age	
Occupation	

key → value
Cust1237 → xyz
Cust1237 → { name: "xyz", street: " ", city: " ", state: " " }

→ used as collection, dictionaries, arrays etc.

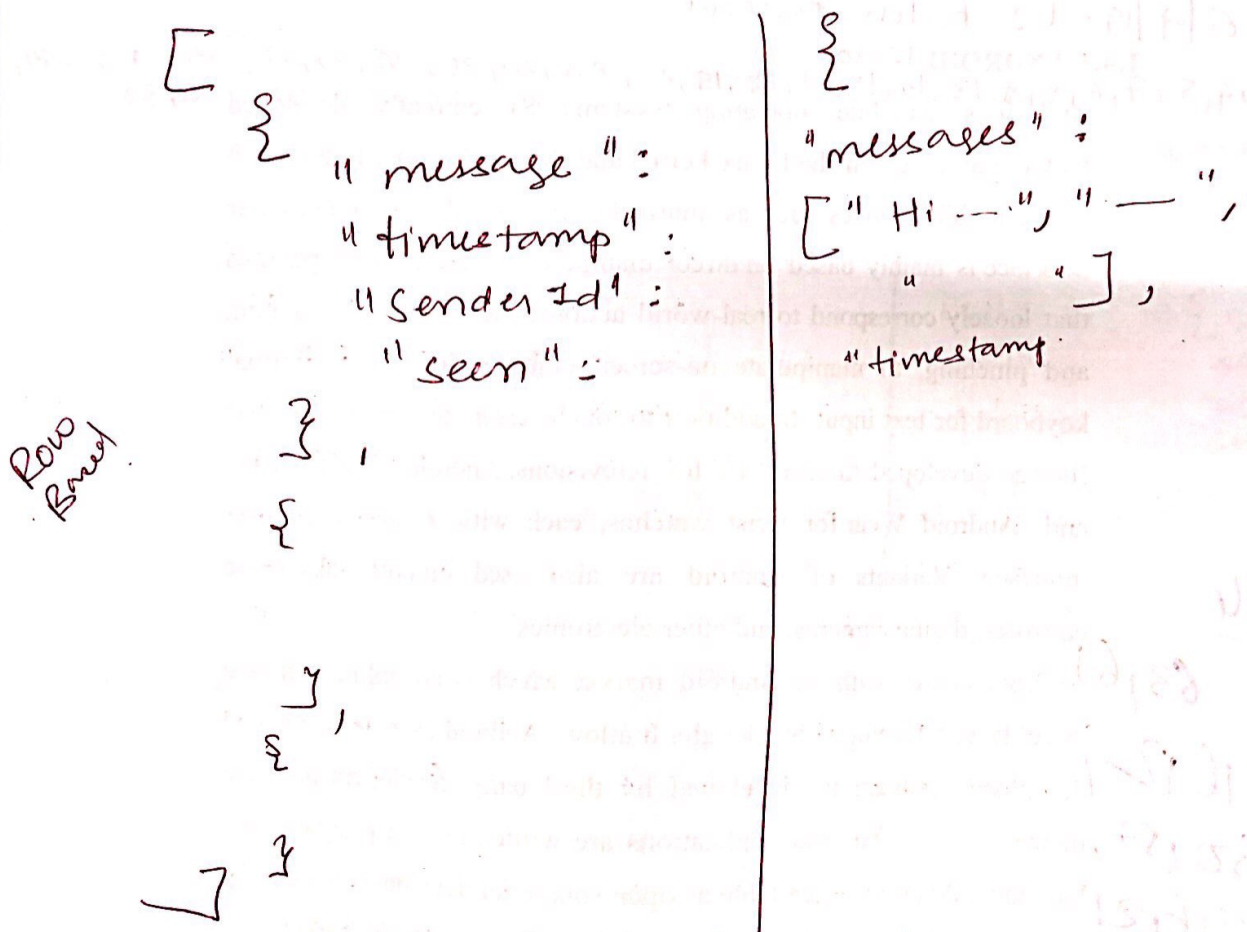
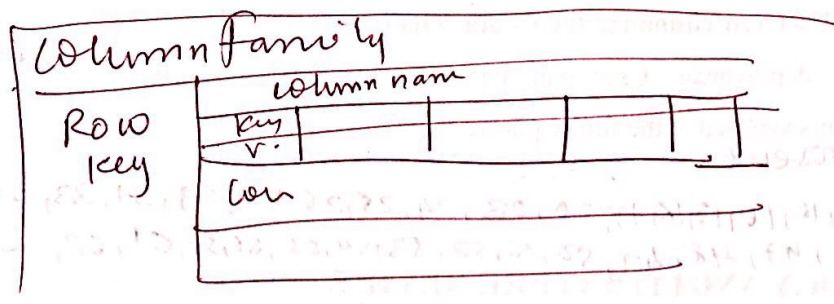
→ Best for → Shopping Cart

→ ex: Redis
Dynamo
Riak.

→

② Column-Based

- Based on BigTable paper by Google.
- every column is treated separately
- Values of single column databases are stored contiguously.



- usually deliver high performance on aggregation queries like sum, count etc. as data is readily available in column
- App Best for - BI, CRM, etc.
- ex. HBase, Cassandra,

③ Document Oriented:

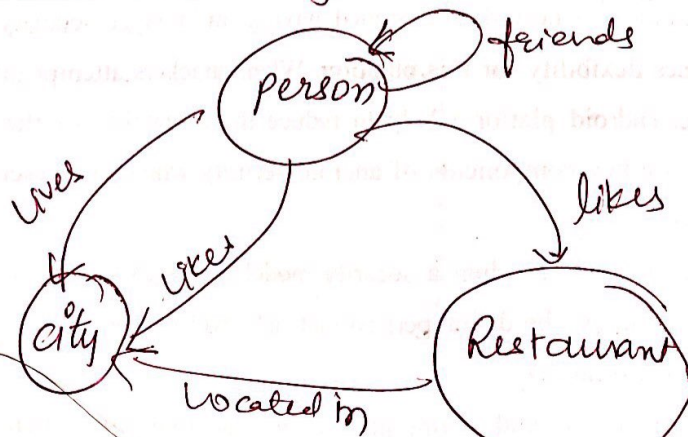
- stores & retrieve data as a key-value pair but value part is stored as document.
- document is stored in JSON or XML format

JSON object → document

- Best for - cms system, blogging platform, e-commerce.
- ex. CouchDB, MongoDB, Rak, Lotus notes, Amazon Simple DB

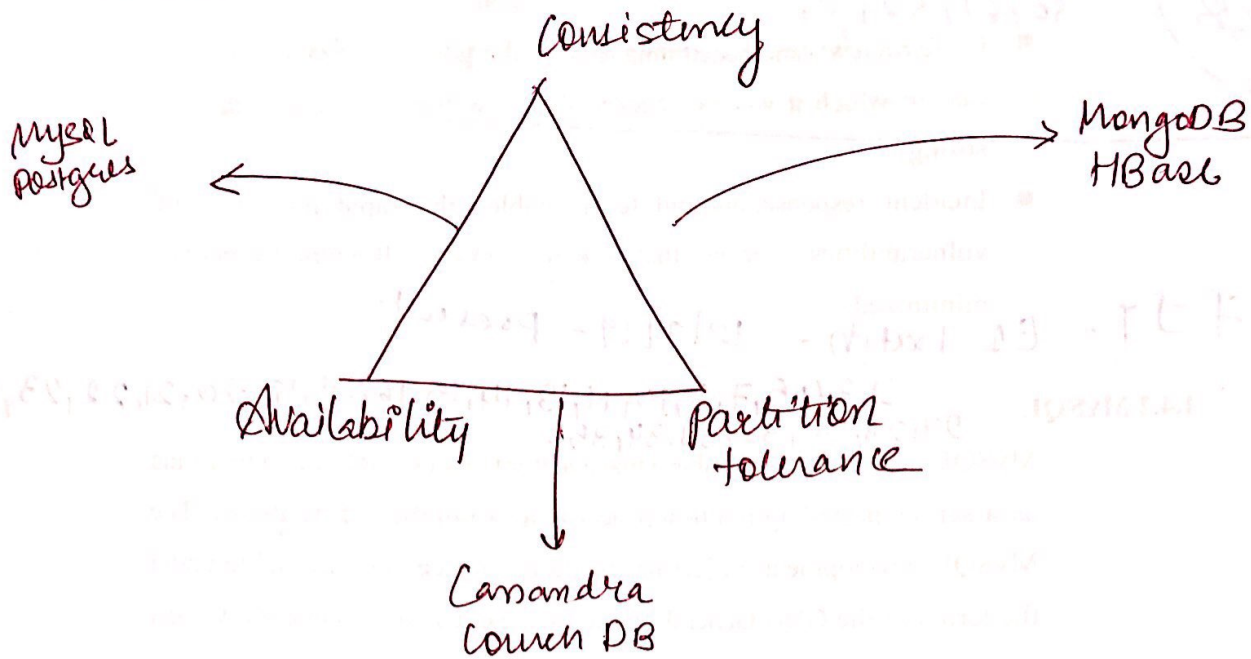
④ Graph Based

- A graph type database stores entities as well the relations amongst those entities.
- entity → node
relation → edge
- every node and edge has a unique identifier.



- Best for Social networks, logistics, spatial data
- ex. Neo4j, or graph, OrientDB etc.

* CAP Theorem!



Consistency: data is the same across the cluster, so you can read or write from/to any node and get the same data

Availability: Ability to access the cluster even if a node in the cluster goes down

Partition tolerance: cluster continues to function even if there is a "partition" between 2 nodes
communication break
i.e both nodes are up, but can't communicate.

→ In order to get both Availability & Partition tolerance, you have to give up consistency.

MOBILITY

➤ Supports stunnel 4 connections for access to MySQL

→ Consider if you have 2 nodes.



→ There is break betn n/w communication between X & Y.

→ so they can't sync updates.

→ At this point you can either:

① Allow the nodes to get out of sync.
(giving up consistency)

② Consider the cluster to be down
(giving up availability).

→ All the combinations available are:

CA: data is consistent between all nodes as long as all nodes are online.

CP: data is consistent between all nodes, & maintain partition tolerance by becoming unavailable when node goes down

AP: nodes remain online even if they can't communicate.

* Strict Consistency:

ACID compliance:

- ① Bank balance is \$50
- ② deposit \$100
- ③ queried from any ATM anywhere, is \$150
- ④ withdraw \$40
- ⑤ queried from any ATM anywhere, is \$110

* Eventual Consistency: wether report

- To have copies of data on multiple machines to get high availability & scalability.
- Thus, changes made to any data item on one machine has to be propagated to other replicas.
- Data replication may not be instantaneous as some copies will be updated immediately while others in due course of time.
- These copies may be mutually, but in due course of time they become consistent.

* BASE → Basically Available, Soft State, Eventual consistency

→ DB is available all the time as per CAP theorem

Soft state

→ even without an input, the system may change

Eventual consistency

→ The system will become consistent over time.

* HBASE

- designed to provide quick random access to huge amounts of structured data.
- column-oriented database & the tables in it are sorted by row.
- The table schema defines only column families, which are key-value pairs.
- A table have multiple column families & each column family can have any no. of columns
- Subsequent column values are stored contiguously on the disk.

→ In HBase:

- - Table - collection of rows.
- ROW - collection of column families
- Column family - " " columns.
- columns - " " key value pairs.

Row id	Column family				Column family			Column family		
	col1	col2	col3	col4	col1	col2	col3	col1	col2	col3

- Rows are sorted alphabetically by the row key, as they are stored.
- For this reason, the design of the row key is very important.
- The reason is to store data in a such a way that related rows are near each other.

→ A common low key pattern in website domain.

→ store them in reverse.

e.g. }
 org.apache.www,
 org.apache.mail,
 org.apache.jira

This way all of the apache domains are near each other in the table, rather than being spread out.

→ A table in HBase would look like.

Rowkey	column family - Personal		column family - office	
	Name	Residence phone	phone	Address.
00001	John Paul	415-111-1234	415-212-55	1021 market st
00002				

```

{
  "00001":
  {
    "personal":
    {
      "Name":
      {
        "Timestamp1": "John"
      }
      "Residence phone":
      {
        "ts1": "___",
        "ts2": "___"
      }
    }
  }
  "office":
  {
    ...
  }
}
  
```

* Row-Oriented Vs Column Oriented

RDBMS - Row-oriented
HBase - Column-oriented.

- Row-oriented databases store table records in a sequence of rows.
- Column-oriented databases store table records in a sequence of columns.
i.e. the entries in a column are stored in contiguous locations on disk.

e.x

Cust Id	Name	Addr	Product Id
1	John	US	231
2	Paul	UK	520

→ Row oriented → 1, John, US, 231
2, Paul, UK, 520.

→ Column oriented → 1, 2, John, Paul, US, UK, 231, 520.

→ When amount of data is very huge, like in terms of petabytes or exabytes, we use column-oriented approach.

→ So, the data of a single column is stored together and can be accessed faster.

→ HBase tables has following components,

Rowkey	column family			
	customers.		products	
cust#0	Name	city & country	prod. name	price

Annotations in the diagram:
 - A bracket above the first column is labeled "Row key".
 - A bracket above the last two columns is labeled "column family".
 - A bracket on the right side of the last two columns is labeled "column qualifiers".
 - A bracket on the right side of the last two columns and the row below is labeled "cell".

① Tables: Data is stored in table format. In HBase, here, tables are in column-oriented format.

② Row key: used to search records which makes searches fast.

③ column families: Various columns are combined in a column family.

- These column families are stored together which makes the searching process faster as data belonging to same column family can be accessed together in a single seek.

④ column qualifiers: column's name

⑤ cell: data is stored in cells.
 - The data is dumped into cells which are specifically identified by rowkey & column qualifiers.

⑥ Timestamp: Combination of date & time.
 Whenever data is stored, it is stored with its timestamp.
 - This makes it easy to search for a particular version of data.

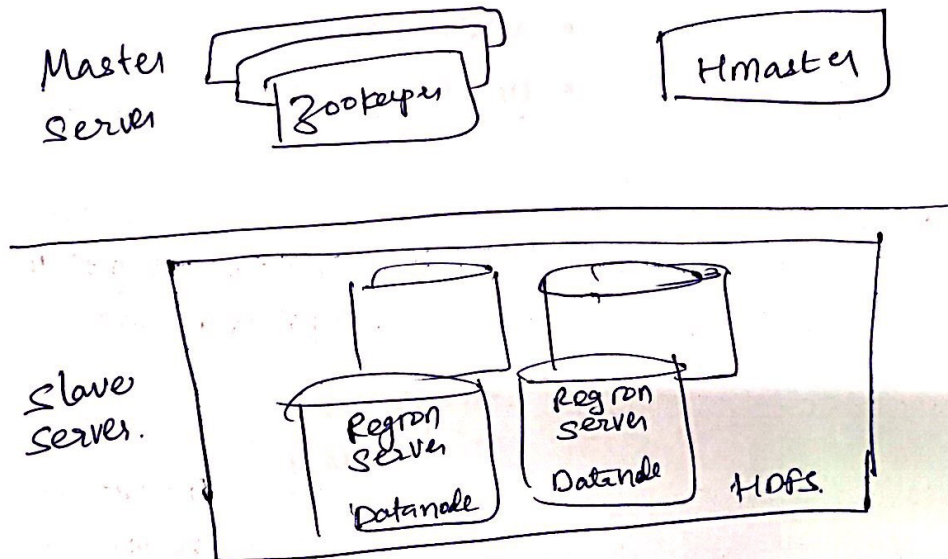
* HBASE ARCHITECTURE :

⇒ Components of HBase Architecture:

③ main components.

- ① HMaster server
- ② HBase Region server
- ③ Regions.

and 4th component is Zookeeper.



→ All these servers (HMaster, Region server, Zookeeper) are placed to coordinate & manage Regions & perform various operations inside Regions.

* Ro

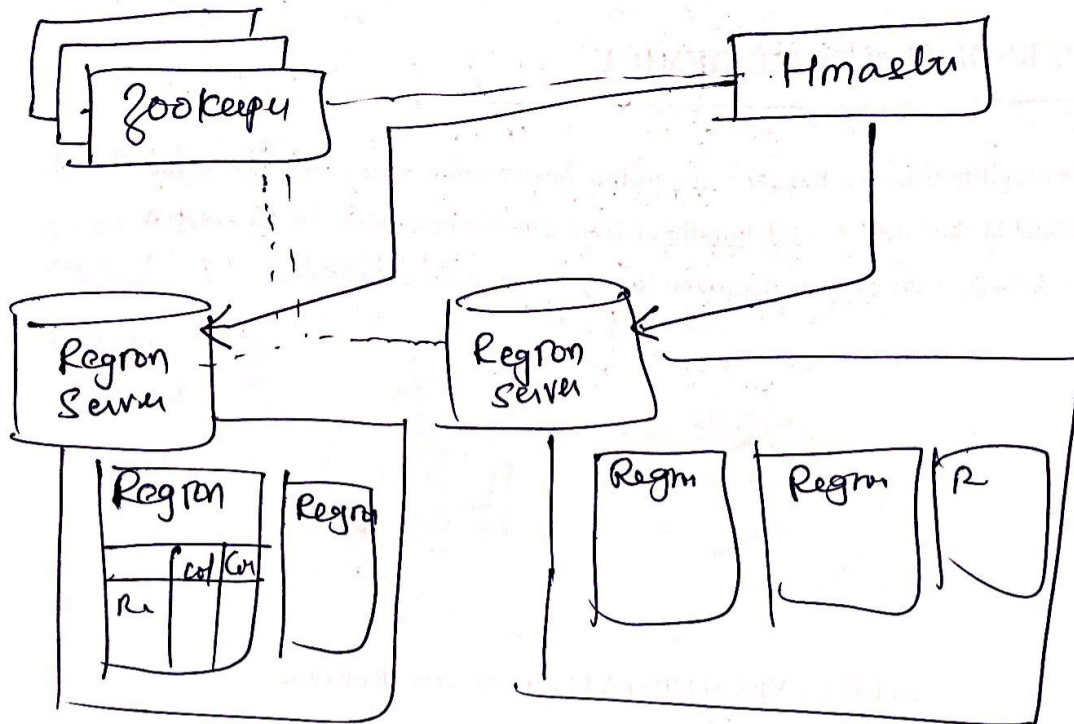
★ Regions:

- A region contains all the rows between the start key and the end key assigned to that region.
- HBase tables divided into such a way that all columns of a column family is stored one region.
- Each region contains the rows in a sorted order.
- Many regions are assigned to a Region Server, which is responsible for handling, managing, executing reads & writes operations on that set of regions.
- Region has a default size of 256 MB which can be configured according to the need.
- A group of regions is served to the clients by a Region Server.
- A region server can serve approximately 1000 regions to the clients.

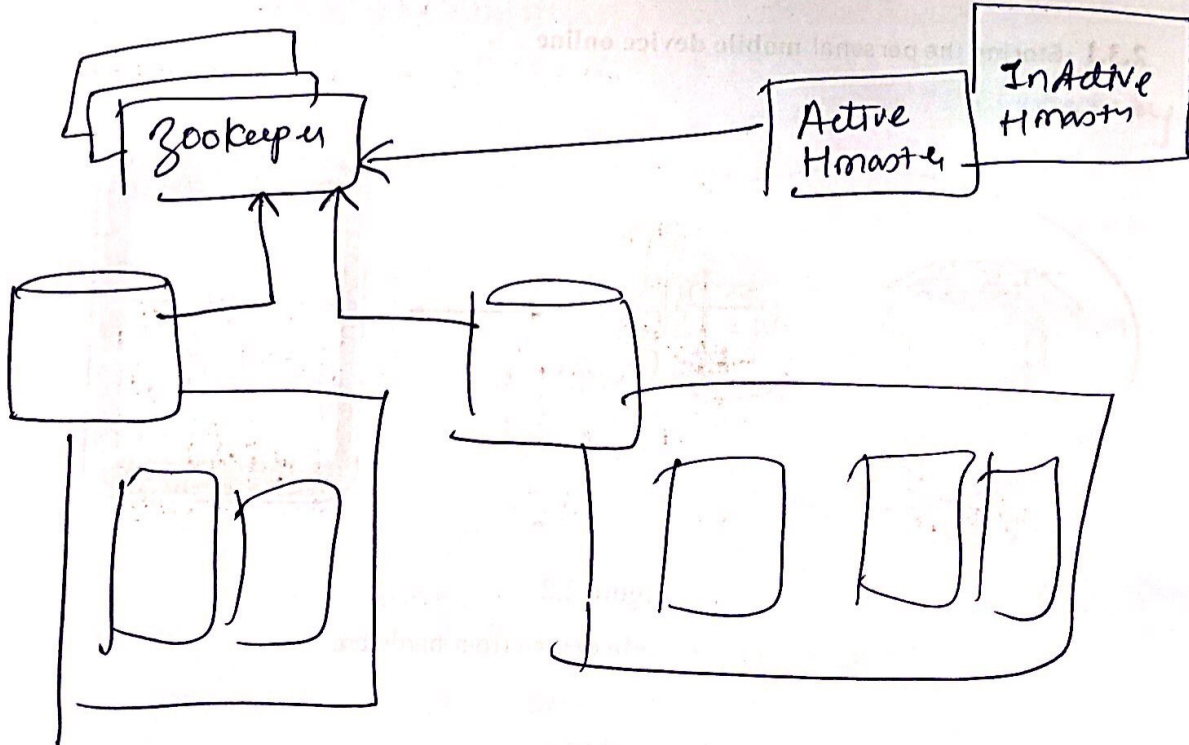
★ HMaster:

- HMaster handles a collection of Region Servers which resides on a Datanode.
- HBase HMaster performs DDL operations & assigns regions to the Region Server.
- It coordinates and manages the Region Server.
- It monitors all the Region Server instances in the cluster (with the help of ZooKeeper) & performs recovery activities whenever any region server is down.

→ It provides an Interface for creating, deleting & updating tables.



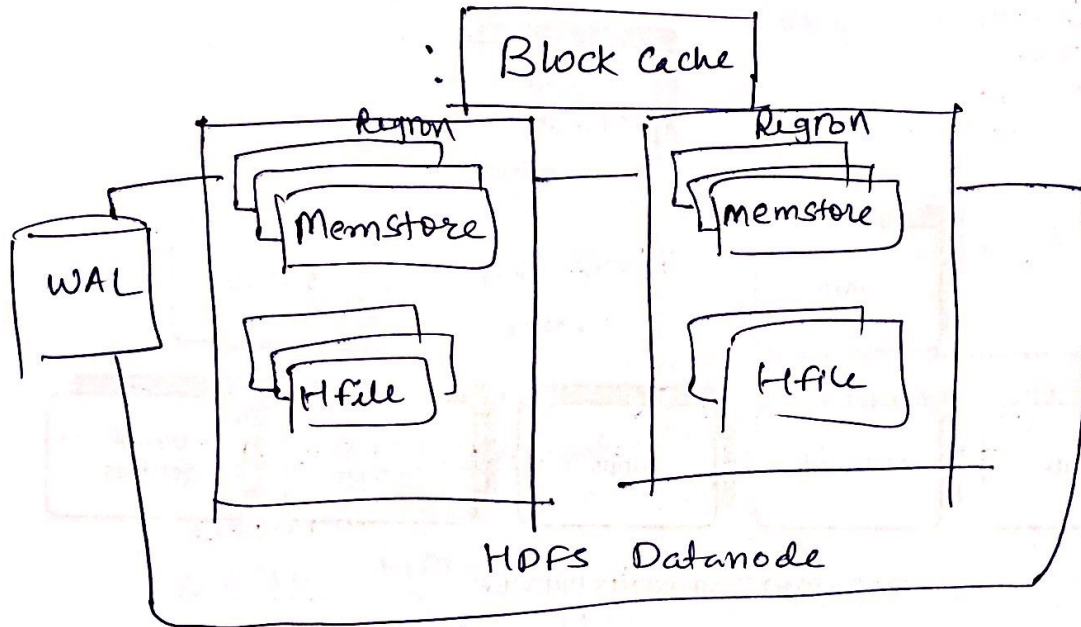
★ Zookeeper - The Coordinator



- Zookeeper acts like a coordinator inside HBase distributed environment.
- It helps in maintaining server state inside the cluster by communicating through sessions.
- Every Region server along with HMaster sends continuous heartbeat at regular interval to Zookeeper and it checks which server is alive and available.
- Here, Inactive server act as a backup for active server.
- Active master sends the heartbeat to the Zookeeper while Inactive HMaster listens for the notification send by active HMaster.
- If the active HMaster fails to send a heartbeat, the session is deleted and the Inactive HMaster becomes active.
- While if Region server fails to send a heartbeat, the session is expired & all the listeners are notified about it.
- Then HMaster performs suitable recovery actions.
- Zookeeper also maintains the .META server's path, which helps any client in searching for any region.

* Region Server:

→ Region Server maintains various regions running on top of HDFS.



① WAL:

→ Components of Region server:

① WAL:

- It is a file attached to every region server inside the distributed environment.
- The WAL stores the new data that hasn't been persisted or committed to permanent storage.
- It is used in case of failure to recover datasets.

② Block Cache:

It stores the frequently read data in the memory.

- If the data in Block cache is least recently used, then that data is removed from Block cache.

③ Memstore: It is write cache.

- It stores all the incoming data before committing it to the disk or permanent memory.
- Memstore for each column family in a region.
- The data is sorted in lexicographical order before committing it to the disk.

④ HFile: It stores the actual cells on the disk.

* Search mechanism

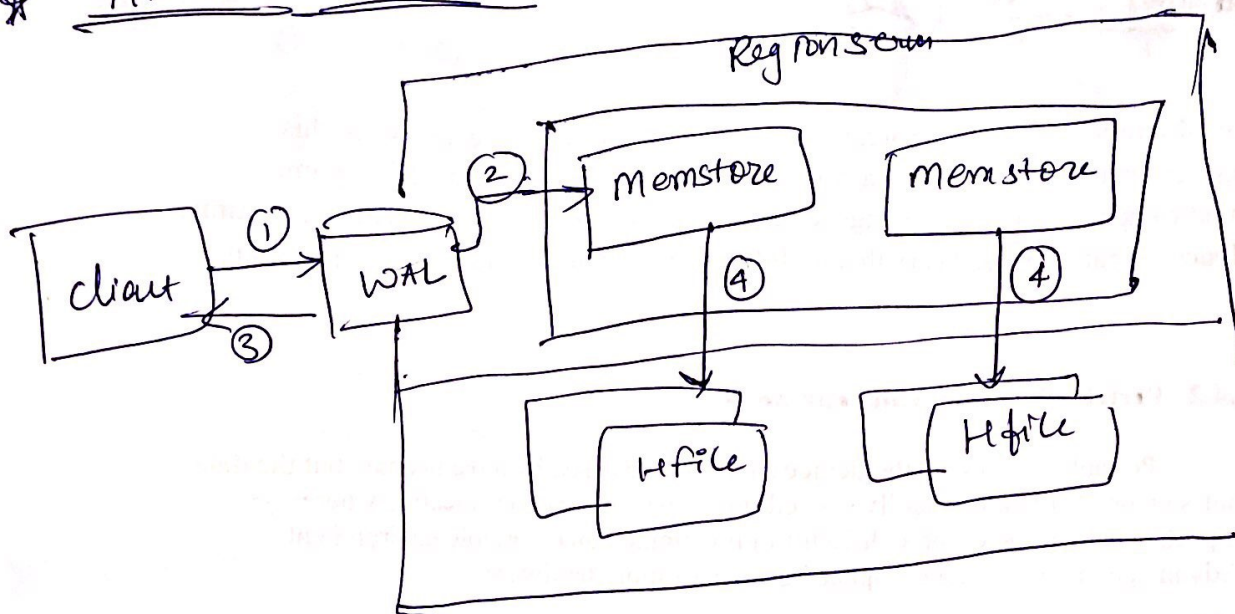
* → Zookeeper stores the META table location.

- ① client retrieves the location of the META table from zookeeper
- ② client then requests for the location of the region server of corresponding row key from the META table to access it.
- ③ Then it will get the row location by requesting from the corresponding Region server.

→ For future references, the client uses its cache to retrieve the location of META table & previously read row key's Region server.

→ client ~~is~~ will not refer to the META table, until & unless there is a miss.

* HBase Write Mechanism!



- ① Whenever client has a write request, the client writes the data to the WAL
 - The edits are then appended at the end of the WAL file.
 - This WAL file is maintained in every Region server uses it to recover data which is not committed to the disk.
- ② Once the data is written to the WAL, then it is copied to the memstore
- ③ Once the data is placed in memstore, then client receives the acknowledgement.
- ④ When memstore reaches the threshold, it dumps or commits the data into a Hfile.

* HBase Read Mechanism:

- first client retrieves the location of the Region server from META server if the client does not have it in its cache memory.
- ① for reading the data, the scanner first looks for the row cell in Block Cache.
 - Here all the recently read key value pairs are stored.
- ② If scanner fails to find the required result, it moves to the memstore.
 - then it searches for the most recently written files, which has not been dumped yet in Hfile.
- ③ At last, it will use block cache to load the data from Hfile.

* HBase: Compaction

- HBase combines Hfiles to reduce the storage and reduce the no. of disk seek needed for a read.
- This process is called compaction.
- compaction chooses some Hfiles from a region and combines them.
- ② types

↳ ① Minor compaction

- HBase automatically picks smaller Hfiles and recommits them to bigger Hfiles
- It performs merge sort for committing smaller Hfiles to bigger Hfiles.
- helps in storage space optimization.

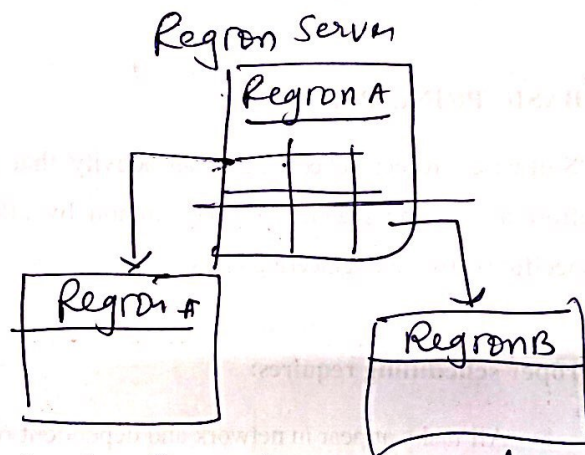
L → (c) Major compaction:

- HBase merges and re-commits the smaller HFiles of a region to a new HFile.
- Same column families are placed together in the new HFile.
- It drops deleted & expired cell in this process.
- It increases read performance.

→ During compaction process, input-output disk & new traffic might get congested.

→ known as write amplification.

* HBase: Region split



→ whenever a region becomes large, it is divided into 2 child regions.

→ The split is reported to the HMaster.

→ This is handled by the same Region server until the HMaster allocates them to a new region server for load balancing.

* HBase Crash and Recovery:

- whenever a Region Server fails, Zookeeper notifies to the HMaster about the failure.
- Then HMaster distributes and allocates the regions of crashed RegionServer to many active region servers.
- To recover the data of the Memstore of the failed Region server, the HMaster distributes the WAL to all the Region servers.
- Each region server re-executes the WAL to build the Memstore for that failed region's column family.
- The data is written in chronological order (in a timely order) in WAL. Therefore, re-executing that WAL means making all the changes that were made and stored in the memstore file.
- So, after all the Region servers execute the WAL, the memstore data for all column family is recovered.