

SYSTEM

DESIGNING



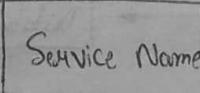
Introduction

- To Prepare You for System Design Interviews in top Companies
- To help you to make better architectural Design Decisions at your Job.
- To improve your Skills of analysing architectural Diagrams.

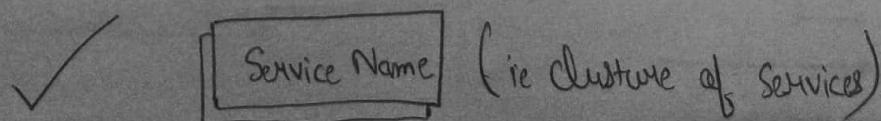
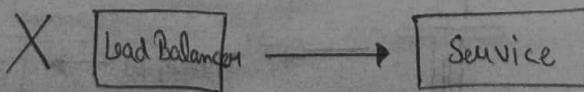


Diagram Building Blocks :-

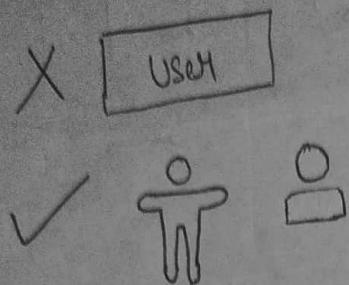
① Service :-



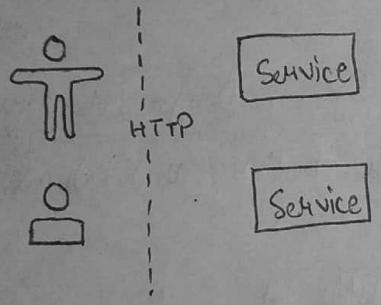
② Load Balanced Service :- Sometimes We don't run a Single Service but multiple Services at same time by putting Load balancer behind it, but it is wrong.



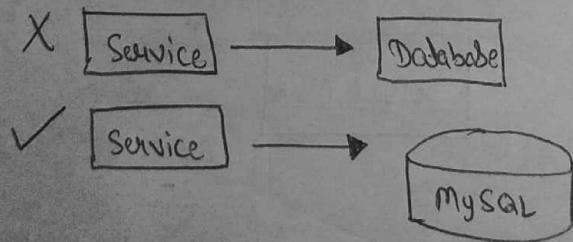
③ User :- feel free to draw them, but different symbols for different user categories.



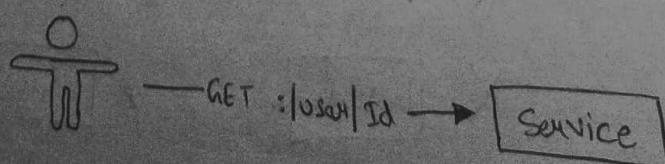
④ Create Network Dividers :- Now users generally communicate over a service and it is good to put divider b/w user and service by giving protocol name, that which protocol we use.



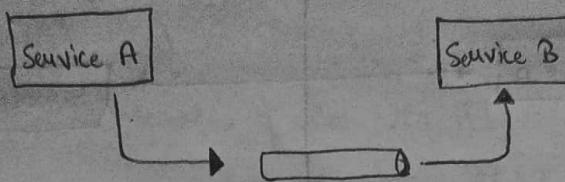
⑤ Database :- It is better to draw the diagram of database with Name.



⑥ Synchronous channel :- These are http calls drawn by Solid Arrow.

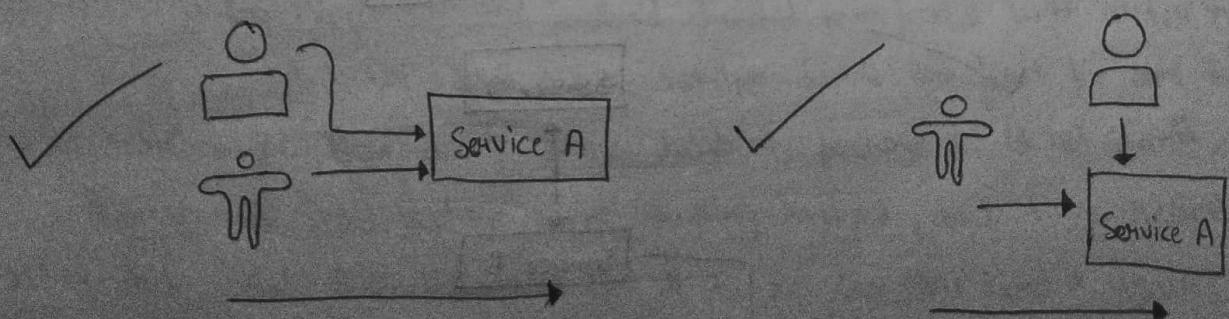
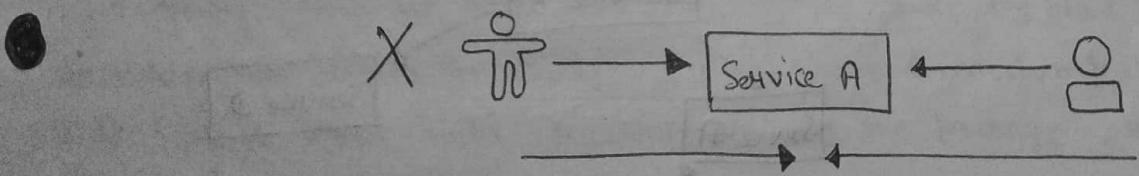
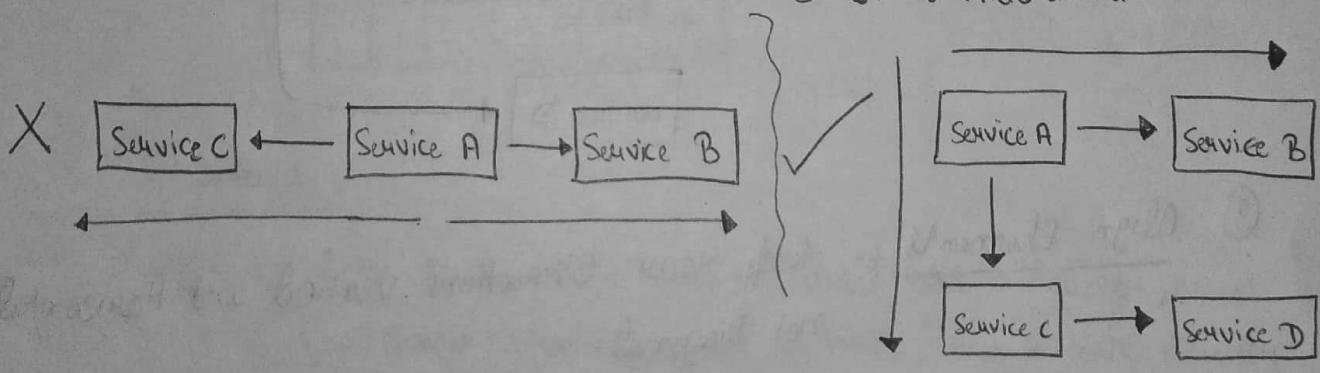


Queue :- It can be implemented via dotted line or using barrel.

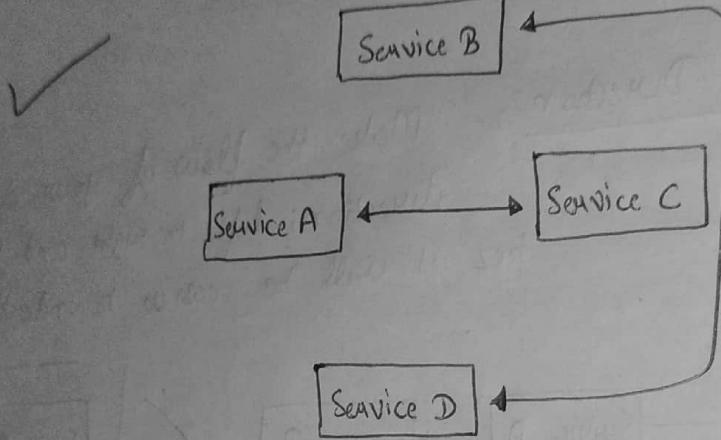
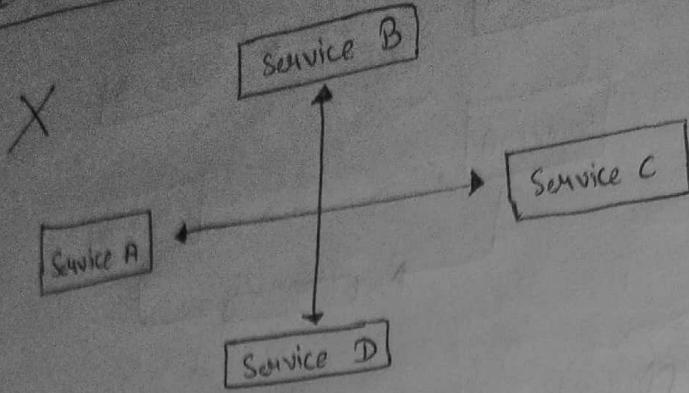


★ Diagram flow :-

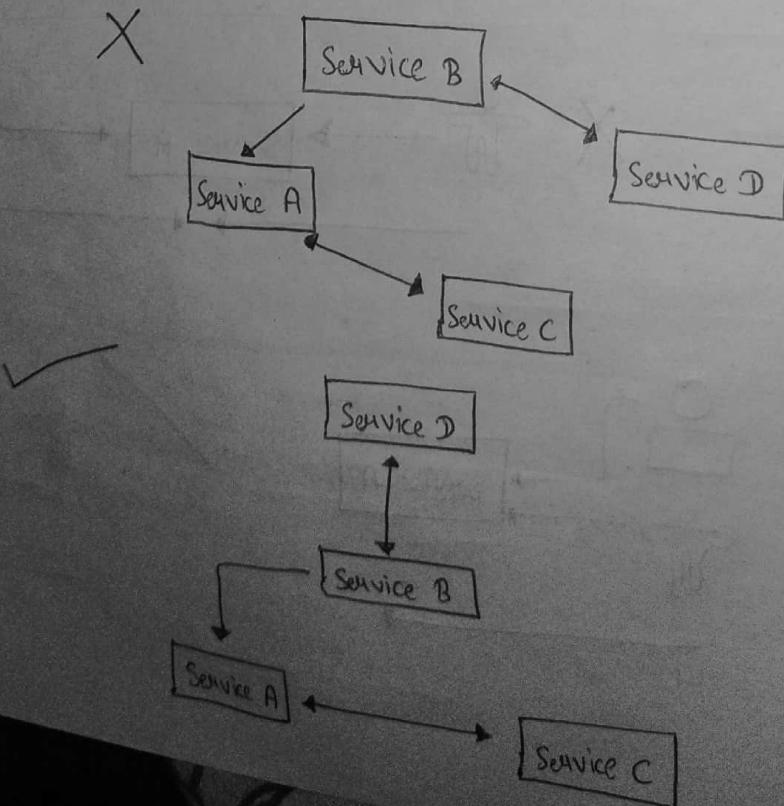
- ① Maintain Direction :- Make the flow of your diagram one direction, left to right and top to bottom bcz it will be easier to understand.



② Avoid Intersections :-



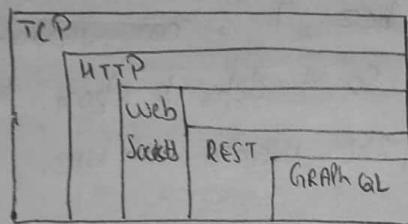
③ Align Elements :- keep your Connections Vertical and Horizontal not Diagonal.



* Networks :-

① TCP :- Transmission Control Protocol and it is a most common protocol used in Distributed Systems and we mostly used TCP, it means, if we use HTTP, we use TCP
• if we use REST, we use TCP
• if we use GraphQL, we use TCP
• if we use WebSockets, we use TCP

that means TCP is everywhere.



TCP Advantages :-

(i) Reliable :- Let's suppose we have a Sender and a Receiver, now if a Sender wants to send an image, then it can't directly send the image, bcz first the payload will break into pieces (i.e. image), now we will send the image piece by piece, and if Receiver does not receive the package due to some connection problems, TCP will send it again until Receiver receives the package, that makes it reliable.

(ii) ordered :- Now let the payload divided into pieces with order no. 1, 2, 3, 4 and when Sender sends the first pic and then second pic then due to some reasons package 2 is not sent, then it will send package 3 as nothing happened, but then Receiver asked that i got #3 but don't get #2, so can you send it again?, hence it helps us to recognize which packets we missed, which makes it ordered.

(iii) Error- checked :- TCP is also has a checksum for each packet.
Now let us suppose we have 4 packets and we send them sequentially. Suppose the #2 package of size 21B but due to some reason it got corrupted, then TCP checks the 'checksum' of both the packages (i.e. before and after sending) and if it is corrupted then receiver will say, I got 8B message, but you told me it should be 21B, Can you send it again?, This is how it helps in error checking.

② UDP :- User Datagram Protocol, it is fast as compared to TCP, bcz it is not reliable and ordered b/c it is fast, So it depends upon the System designer which protocol he/she wants to use.

Where UDP Used :- It is used for constant stream of data or sending a lot of data fast in less time.

- (i) Monitoring metrics.
- (ii) Video Streaming.
- (iii) Gaming.
- Imp (iv) Stock Exchange.

* Note :- Can you afford losing some data?, then surely you can use UDP, bcz it is fast as compared to TCP.

TCP

- ① Expects confirmation from receiver.
- ② Numbers each message (i.e. ordered).
- ③ Add checksums for each message.
- ④ Slow as compared to other protocols.

UDP

- ① Just sends packets of information.
- ② Unordered
- ③ Good, when there is constant stream of data.
- ④ It is fast.

Load Balancing :-

- Load Balancer is a machine that runs a service known as software.
- Goal of the software is to distribute the requests b/w multiple servers that host the actual application.

Interview :-

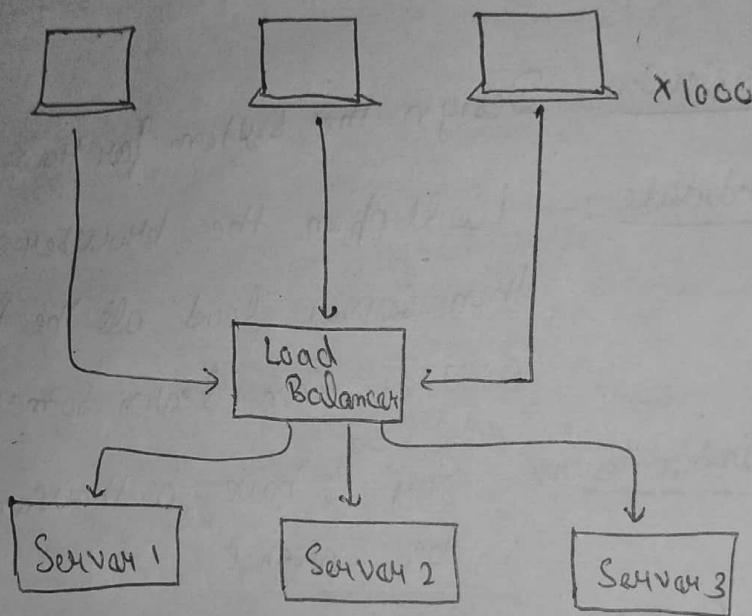
- Interviewer :- Design the system for how Udemy works?
- Candidate :- I will open the browser and type Udemy.com and then server load all the pages and videos and then a user can learn something new!
- Interviewer :- Say U have a thousand requests, will one server be enough?
- Candidate :- Of course not, we have multiple servers.
- Interviewer :- Then how user knows the tasks of each server?
- Candidate :- Yeah well, i will place one server in the middle of client and all the servers which will accept all the requests from the client and it is known as Load Balancer.

Now, Is Load Balancer knows which request sent to which server? and for this, Strategies come, and every load balancer different strategies.

- Round Robin
- Least Connections
- Resource Based
- Weighted variants of the above.
- Random

Now we will understand the most common
i.e "Round Robin"

⇒ In this strategy, Load Balancer will distribute the load
among different servers in Round Robin format.



⇒ Types of Load Balancers :-

- ① Level 4 :-
- Transport Layer Load Balancer
 - Has Access to: TCP or UDP, IP, Port
 - Low Level

- ② Level 7 :-
- Application Layer Load Balancer,
 - Has Access to everything layer 4 has
 - In Addition, has Access to HTTP Headers, Cookies, Payload.

Examples of open source load balancers :- NGINX, HAProxy,
HAProxy, etc.

Advantage :- (i) Resilience :- when if our one of the server is down (i.e. fails to load) then load

Balance will route the incoming request to another server.

(ii) Scalability :- Load Balancer makes your system Scalable, it means suppose we have 1000 clients and 3 servers, but suddenly our clients gets double (i.e. 2000) so, you can simply add more servers connecting to same load balancer, which we called Horizontal scaling.



:- We will understand by giving an example, let suppose we are going to open Udemy.com, but it will take different time in different countries to open it, bcz like Udemy Server is located in USA, then users from USA will experience good quality but like from India, users can't get good quality, bcz data has to travel more for India.

- So, we can solve this issue by using more servers, like servers in every country where we want to give our service but it will be very expensive.
- So we will use CDN, CDN stands for "Content Delivery Network". Its ideology is very simple, that a user will get data / content from the nearest server ASAP.
It is used to deliver static content :- Images, CSS, HTML, JavaScript.

⇒ Two Types of CDN :-

(i) Push :- It is used to just upload a data to all servers and then CDN distribute to all the servers.

(ii) Pull :- When user opens your site and wants to download some data then if CDN server has data, then it will serve it to user.

Note :-

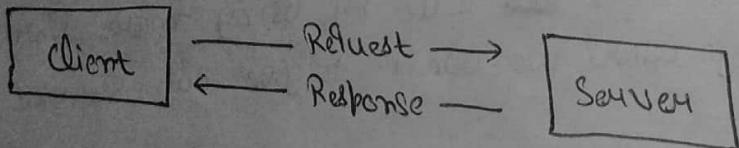
- CDNs allow placing your static assets
- Reduce Costs
- Decreases Latency
- Increases Complexity of your System.

Sometimes, it increases the complexity of our System bcz before using CDN, we will directly send data from the Server, but now we have to push it to CDN or make sure CDN approved these data files, so it will increase complexity of our System.

★ Protocols :-

① HTTP :-

- (i) HyperText Transfer Protocol
- (ii) HyperText → text with links to other Documents
- (iii) Based on TCP
- (iv) Request-Response Protocol
- (v) Layer 7 - Application Layer



⇒ Request :-
(i) Method
(ii) URL
(iii) Headers
(iv) Body (optional)

⇒ Response :-
(i) Status
(ii) Headers
(iii) Body (optional)

⇒ HTTP Methods :-

GET | superhero | 123 ≠ DELETE | superhero | 123

- (i) GET → Head
- (ii) DELETE → delete
- (iii) POST → Create
- (iv) PUT → Update
- (v) PATCH → partial Update

GET :- | superhero | 123



{ "id": 123, "first name": "Bruce", "Last Name": "Wayne" }

Post :- | superhero | 123

{ "firstName": "Bruce", "Last Name": "Wayne" }

DELETE :- | superhero | 123

PUT :- | superhero | 123

"for Updation"

Note :- PUT is similar to PATCH, but PATCH is used for partial updation, like if we just want to update first name, then PATCH will work.

⇒ HTTP Status Codes :-

- 100 - 199 Informational
- 200 - 299 Successful
- 300 - 399 Redirection
- 400 - 499 Client Error
- 500 - 599 Server Error

- ⇒ 200 - 299 Successful Codes :-
- 200 ok
 - 201 created
- ⇒ 400 - 499 Client Error Codes :-
- 401 Unauthorized (unAuth)
 - 403 Forbidden (no Permission)
 - 404 Not Found (no Page)
- ⇒ 500 - 599 Server Error Codes :-
- 500 Internal Server Error
 - 503 Service Unavailable
- ⇒ 300 - 399 Redirection Codes :-
- 301 Moved Permanently

② **REST** :- It is not actually a Protocol but it is just a guideline to how to structure your API.

Representational
State Transfer

- Built on HTTP
- Standard URL Structure
- Standard Use for HTTP Verbs.
- HTTP Status Codes to indicate errors
- JSON as body format.

⇒ Before REST :- So many ways to represent same action
Example for Deletion :-

GET | ?id = 123 & action = delete & type = user

Post |

Body = "delete", type : user;"

Cookie = "userId : 123"

DELETE | ? id = 123 & type = user

⇒ After REST :- DELETE | user) 123

\Rightarrow RESTfulness :- RESTful Verbs

- GET
- POST
- PUT
- PATCH

RESTful URL's :- METHOD | [resource|id]

GET /users/123 and not /user/123 or /getuser/123

RESTful URL's :- Nested Resources

METHOD | [resource|id] | [resource|id]

DELETE /users/123/books/567 \rightarrow return book that was borrowed by user 123

GET /users/123/books \rightarrow fetch all the books borrowed by this user

REST full URL's :- state

PUT | [resource|id] | [action]

Example :- PUT /users/567/enable

PUT /users/567/disable

RESTful URL's :- idempotency

PUT /users/123/online \rightarrow we can set it true/false, so it will be idempotent (i.e. action)

POST /users/123/likes \rightarrow it is not idempotent

So, we can see where to use PUT and where to use POST

RESTful URLs :- Pagination

GET /[resource]/id? limit = x & offset = y

Example

GET /books? limit = 50 & offset = 100

We can provide pagination in our application using Query in API

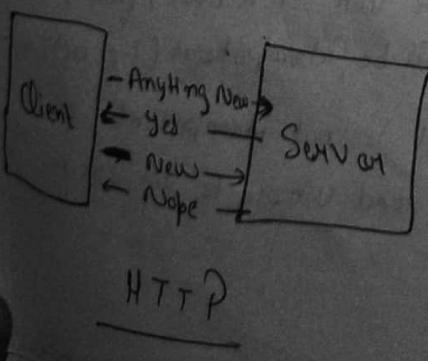
- Note :-
- Use pair of resource-name / resource-id to structure your URLs.
 - GET shouldn't change Anything
 - Use PUT to switch State.
 - Provide Pagination.

③ Web Sockets :-

Now, we know what Protocol aims to solve ?, it gets the request and then give response, (i.e- User -> Request -> Response model) that means Client have to give request and then wait for the Server to give response to terminate the connection.

Now, HTTP is great when we want someone to send a message and then get a response from the Server that msg is received but what about receiving messages., ~~but what about receiving messages.~~

So, To solve this, Web Sockets came to an action, it is a duplex protocol, When the Connection establish b/w Client and Server then messages can be simultaneously.



Web Sockets

- Advantages :-
 - Connection is established only once.
 - Real time message delivery to the client.

- ⇒ Disadvantages :-
 - More Complicated to implement than HTTP.
it means if there is network Error then we have Write code to make connection established
 - May not have best support in some languages.
 - Load Balancer may have some troubles.
 - Unlike REST, need to Implement "the Protocol" every time.

Imp

Note :-

- Built on top of TCP, same as HTTP
- Duplex (two way) Protocol.
- Persistent Connection.
- More efficient than polling with HTTP.



Concurrency :-

Parallelism ≠ Concurrency

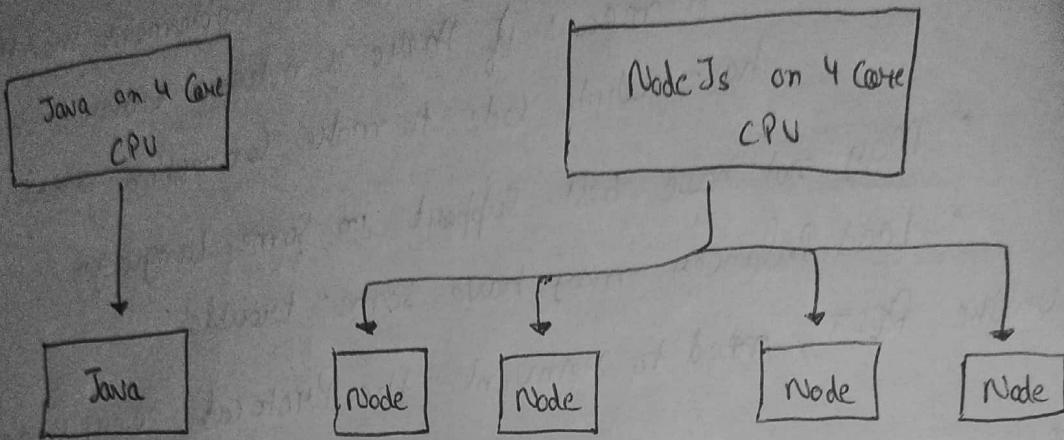
- Parallelism :- Doing more than one thing at the same time.
- Concurrency :- Providing an illusion of doing more than one thing at the same time.

Example :- ① Suppose you and your friend are in a kitchen, you are making a cup of tea and your friend is making a Sandwich on same time, both works are different and can do parallelly

② And if you are alone, you will start making coffee, then go to make sandwich then again go back to check the coffee and you will do this in parallel way by switching yourself from coffee to sandwich. Just like that our CPU works by switching b/w the processes.

Processes

- Each has Separate Memory Space



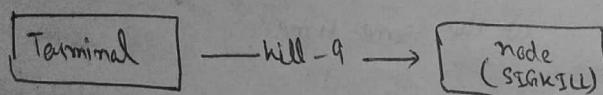
⇒ Interprocess Communication :- if processes cannot share memory, then communication built

• file (also memory mapped file)

When Servers don't want to share, we will create a single file and put all the data coming from servers in that file.

• Signal

\$ kill -9 818



There are other more methods but we will learn that too

Threads

:- We already done this topic in OS & Java in very detail.

DATA BASES :-

⇒ Indexes :- Why do we need indexes?

⇒ To get a information at particular Column/Row.
else we have to Search entire table row by row.

What type of data structure an index is?

⇒ A index should be a hashmap but not usually, bcz in a hash map, we can find value at particular level/point. but if we want to search b/w two points then like,

Select * from videogame_characters where age between 25 and 30,

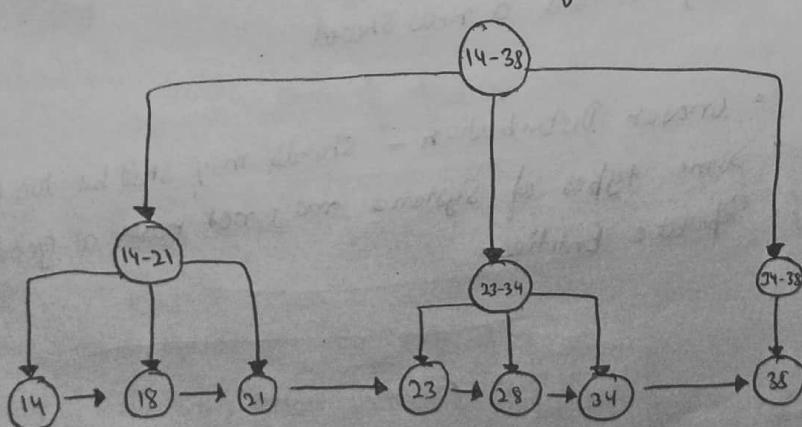
So, We can say hashmap is not suitable, then

"Tree can be a data structure for it", but which tree?
it is B^+ tree,

Why B^+ tree?

⇒ Select * from members where age between 20 and 30?

What makes B^+ tree suitable for it, Is Arrows Using in it.



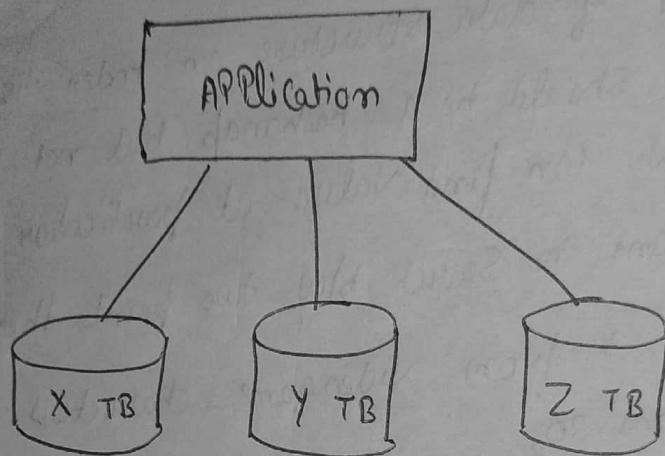
All leaves are clustered and inter-connected.

Use of these Arrows :- Like if we want 21 and 28, then we will go on Left Side, search for 21, then we don't have to go back,
we can use these Arrows to get 28

⇒ Sharding :-

Suppose when our projects grows and then we have to do sharding (i.e. break big database into small pieces).

Different Database for different Services



Now, the problem is, how our application knows to use which database. For this we have two strategies.

(i) Tenant Based Sharding :- it is used for that systems which have partition b/w entities.

Eg:- We do in Taxi-Hailing APP.

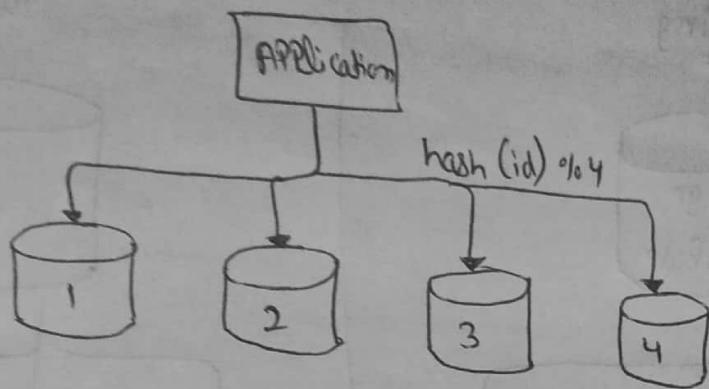
Pros :-

- Easy to understand.
- Easy to add a new shard

Cons :-

- Uneven Distribution - shards may still be too big to handle.
- Some types of systems may not have a good way to separate entities.

Hash Based Sharding :- We will perform Hashing before distribution into the Databases.



- Pros :-
- Even Distribution
 - Works well for key-value Data

- Cons :-
- Adding new shards is difficult, bcz each time we have to change the hash function
 - Weakened Consistency (No foreign keys)

Shard LookUp :- To do this, we need a Locator, that will be known as Shard Locator, in which we program it to do sharding.

Pros :- Simplifies adding new shards, bcz we don't need to keep the logic in our application.

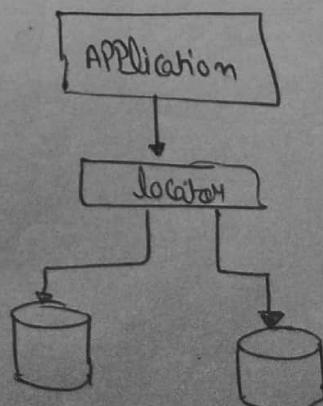
Cons :- Single point of failure

Down Sides of Sharding :-

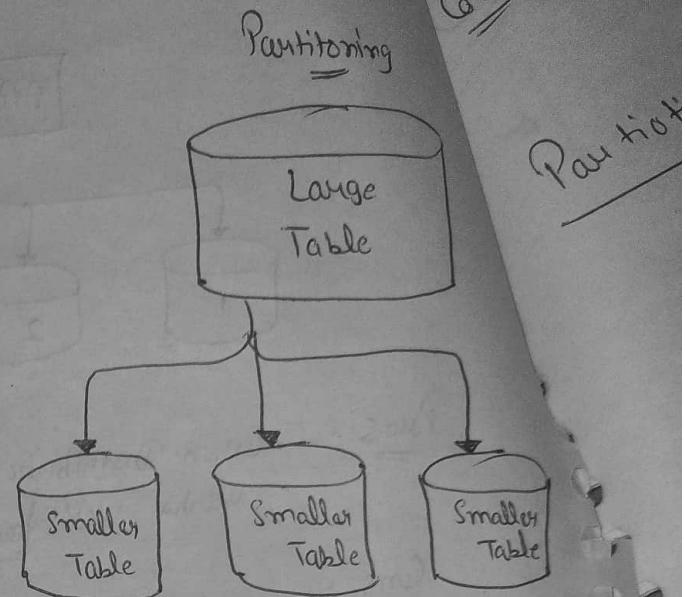
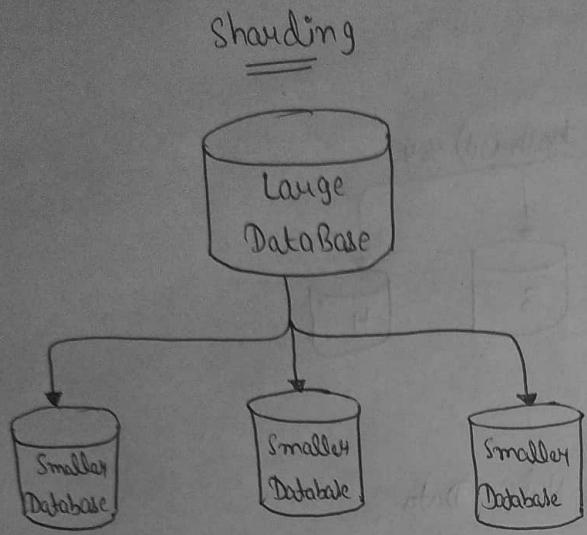
- (i) Added Complexity
- (ii) Cross shard operations are expensive
- (iii) Data still may become unbalanced.

Like :- if we have even distribution of users but users may differ

Some are writers and else are viewers which ends up with different size shards.



⇒ Partitioning :-



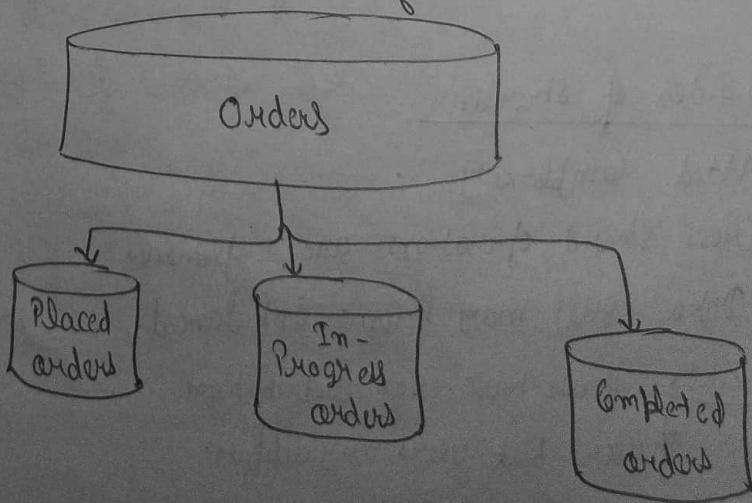
Benefits of Partitioning :-

- Smaller files == faster queries (ie. users will be happy)
- Smaller indexes can fit into memory
- Dropping Partitions is fast.

Partitioning Strategies :-

(i) List of Values :-

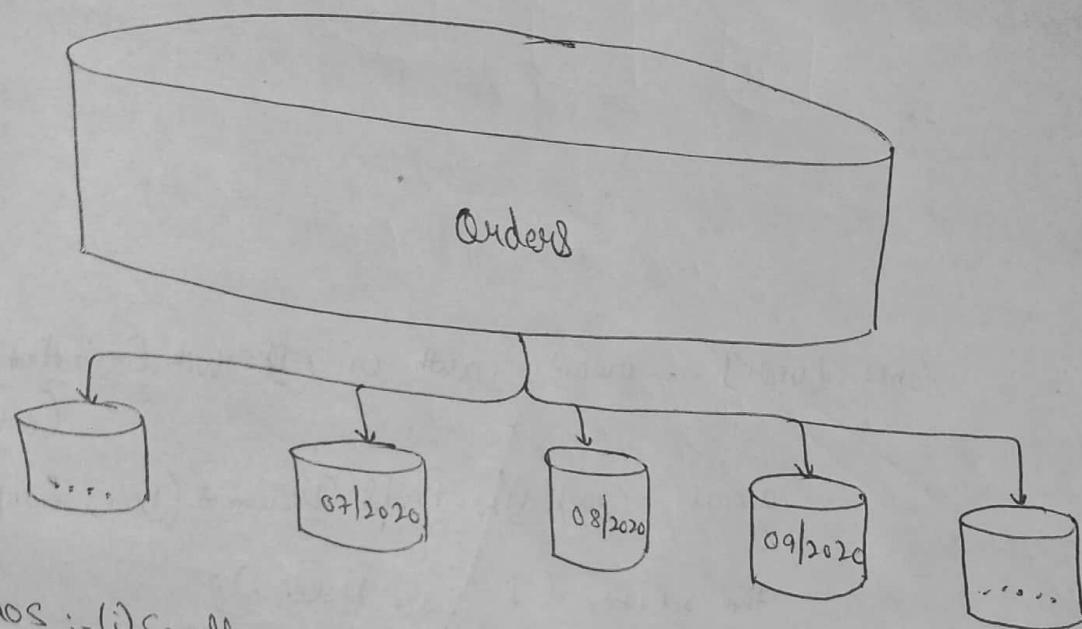
Partitioning by List of Values



smaller tables = faster queries

- Cons :-
- (i) Un-even data distribution.
 - (ii) Need to move data between tables.

Partitioning By Range of Dates :-



- Pros :-
- (i) Smaller tables = faster queries
 - (ii) Cheaper if you want to get rid from old data

- Cons :-
- Still, uneven key distribution.

Partitioning by hash of a key :- As same as Hash-Based Sharding

- Pros :- Even key distribution

- Cons :- Changing no. of partitions is hard.

Down sides of Partitioning :-

- (i) Complexity of maintenance
- (ii) Scanning all partitions is expensive
- (iii) Harder to maintain uniformity; each partition has its own separate table.

Design a Taxi Hailing APP (AKA UBER)

★ Problem Introduction :-

Main Topics :- (i) Marketplace, (ii) Geolocation, (iii) Geostamping

Interviewer :- Design a Uber App for me? (i.e. user Perspective)

Candidate :- A user will open up the application then do Authentication and then he/she can select a drive and book it to reach the destination.

Now, Marketplace Aspects :- (i) Demand / Supply
(ii) At least two very different kinds of users (i.e. Riders and Drivers).

Driver's Perspective :- We may send request to one driver or to multiple drivers on the same time or we can send request to that driver which is nearby.

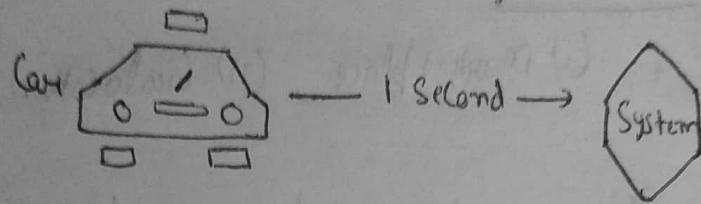
Note :- During Interview, we have to tell all the options to the interview, bcz it tell us about, how far we think.

Marketplace : Riders :- (i) 100 M Riders in total
(ii) 10 M active Riders
(iii) 10 Riders per month on Average

Marketplace : Drivers :- (i) 1 M Drivers in total
(ii) 500k active Drivers
(iii) Average shift : 6 Hours

★ Driver Locations :-

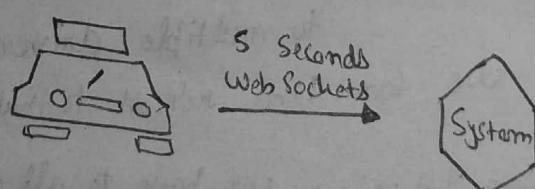
- Frequency of Location Updates :- We need the current location of Driver after each 5 sec.



- Communication :-
- REST (Not Necessarily)
- Web Sockets
- UDP (Much better)

We can use either Web Sockets / UDP, but we will stick to Web Sockets, bcz it will be more preferable for future requirements bcz it allows duplex communication.

- Payload :- Interviewer may ask, which type of data you will send

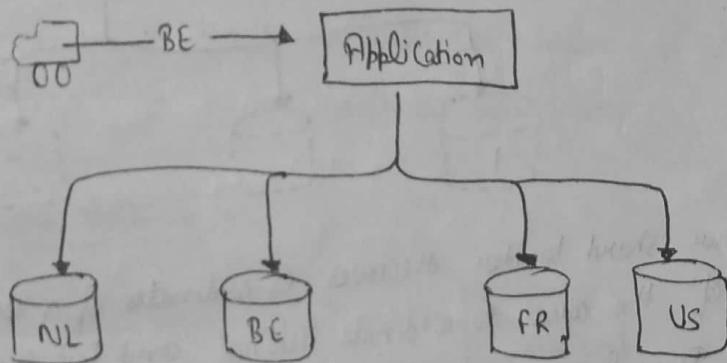


latitude	float
longitude	float
direction?	float
driver-id?	uuid

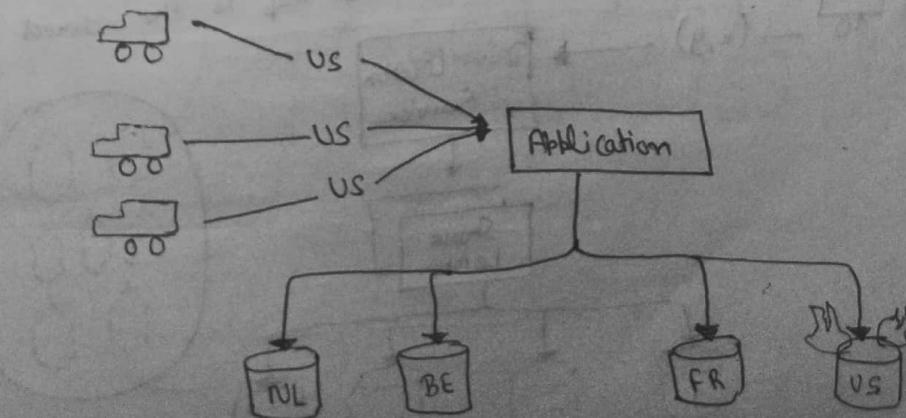
- Through Put :-
- Sock Active Drivers
 - Average Shift :- 6 Hours
 - 4 "shifts" a day
 - 125k drivers at every moment on Average
 - Message every 5 seconds (i.e. current position of driver)
 - 25k requests per second

* Sharding Locations :-

- Sharding :- Attempt #1 ("Shard By Country")

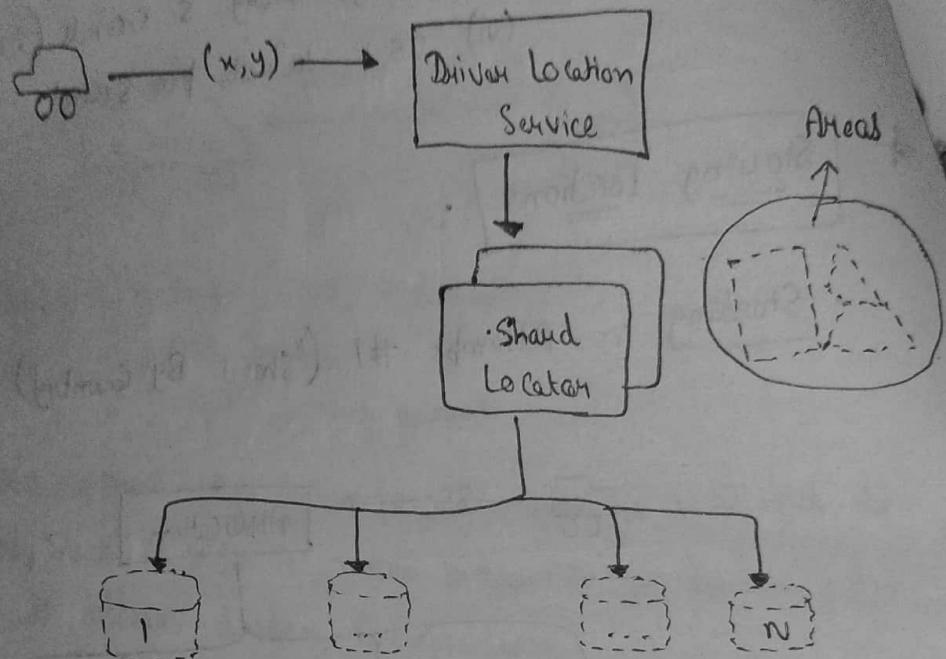


This is a most Common Approach (i.e. shard by Country), we will pick the shard for that Country and write location in it, but it also has a Disadvantage, if most of our drivers are from Single shard and it is called hotspot (i.e. most of the requests sent to one shard).



⇒ Attempt #2 :-

Alternative Would be Used Custom Areas, work if our company has lots of employees. In we will divide the map into shapes as big as countries or small as states.



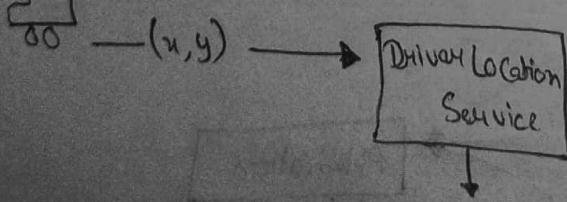
When our shard locator receives Co-ordinates of a Driver, it checks in which of the areas Co-ordinate belongs and put them in correct shard.

Problem :- It requires lots of maintenance.

⇒ Attempt - #3 :-

In the we can use Geospatial libraries like Uber h3, Google S2.

Uber h3 :- Divide the world into hexagon shapes having same area, and then shard locator will map to the correct shard.

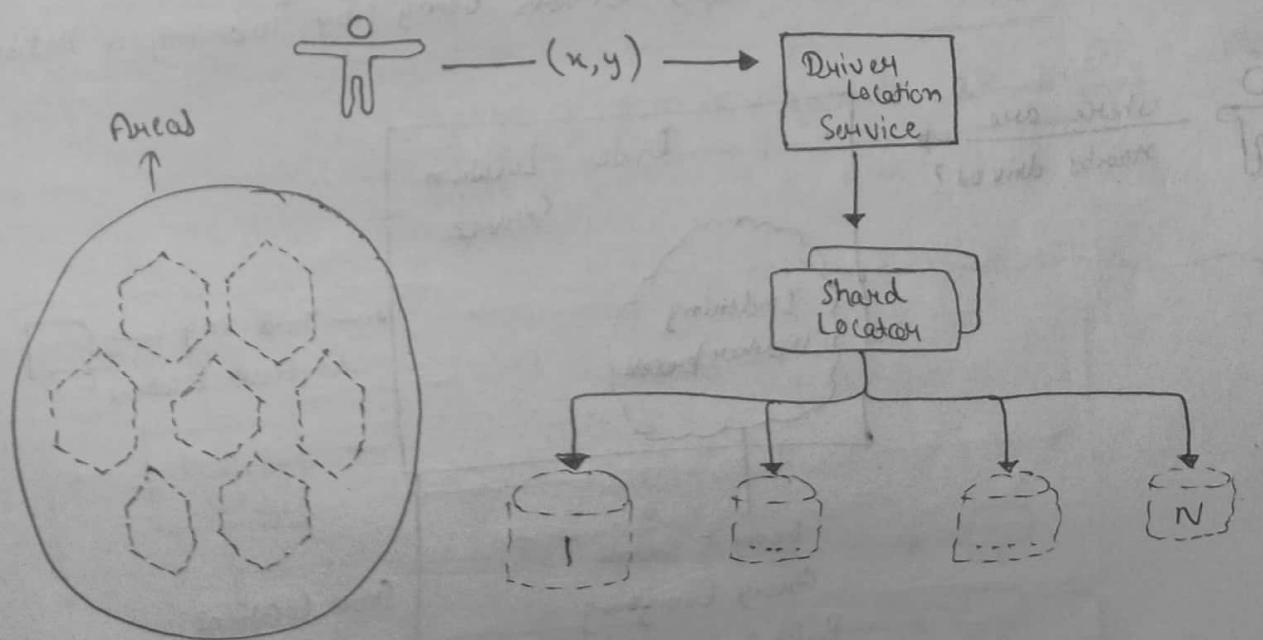


Taxi Around You :-

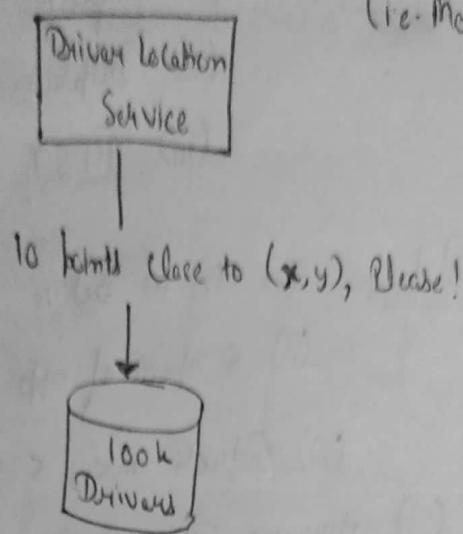
- User Perspective :-
 (i) Refresh every 5 seconds
 (ii) Display 10 drivers almost
 (iii) REST Based

- Through Put :-
 (i) 10 m active Riders
 (ii) 5 open of app a day per user = 50M opens a day
 (iii) Refresh every 5 seconds.
 Total
 (iv) Average Session time is 1 minute = 12 refreshes
 (v) $50M \times 12 = 600M$ requests per day
 (vi) $600M / 24 = 25M$ requests per hour
 (vii) $25M / 60 = 400k$ requests per minute
 (viii) $400k / 60 = \sim 7k$ requests per second

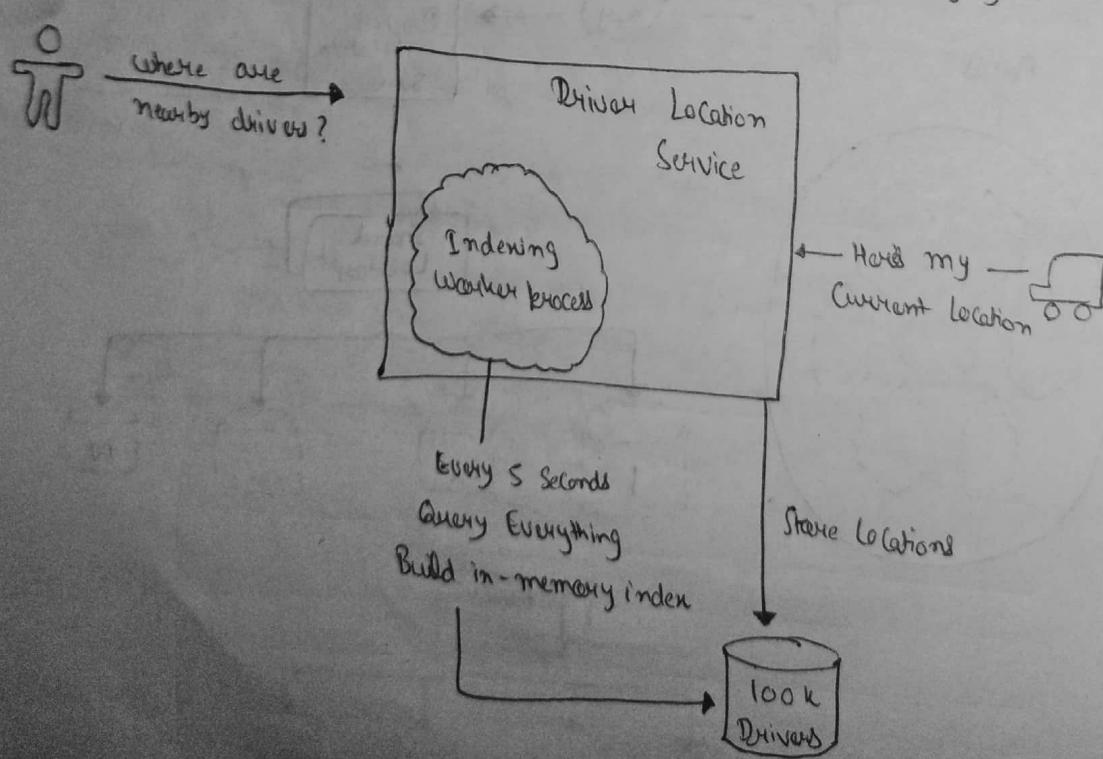
- Getting Drivers Positions :-



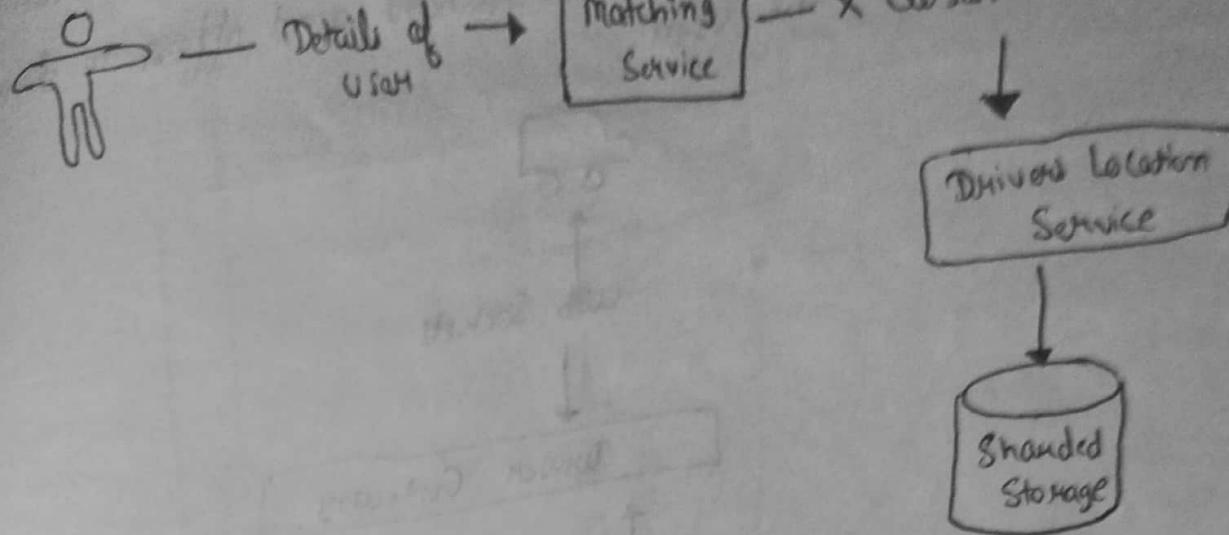
Option #1 :- Geospatial Search in DB → Database that can handle
nearby Locations Search
(i.e. MongoDB, MySQL)



Option #2 :- In-Memory Geospatial Index (Safest approach)
In this Driver Location Service will get the location of
the Driver and then store it in the correct shard but on
a particular index using Indexing Worker process (we can say)
and then User gets nearby Drivers using / by querying a Database.

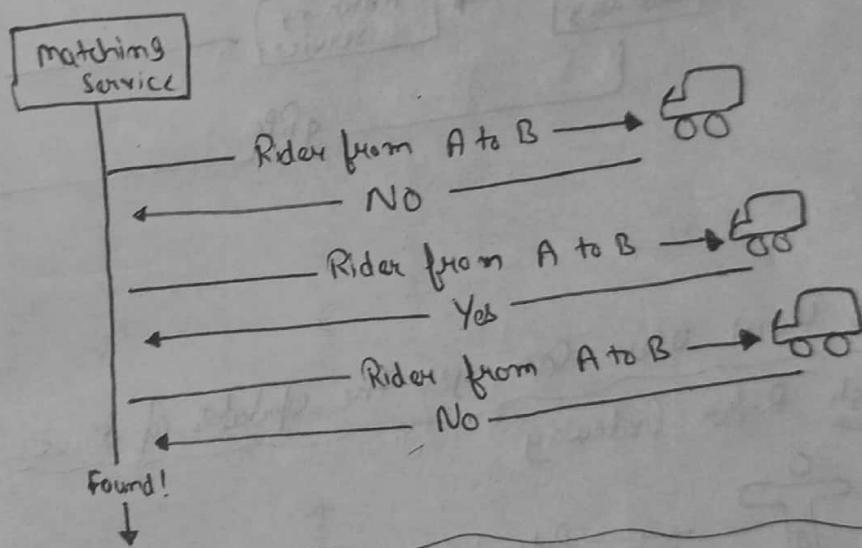


Matching :-

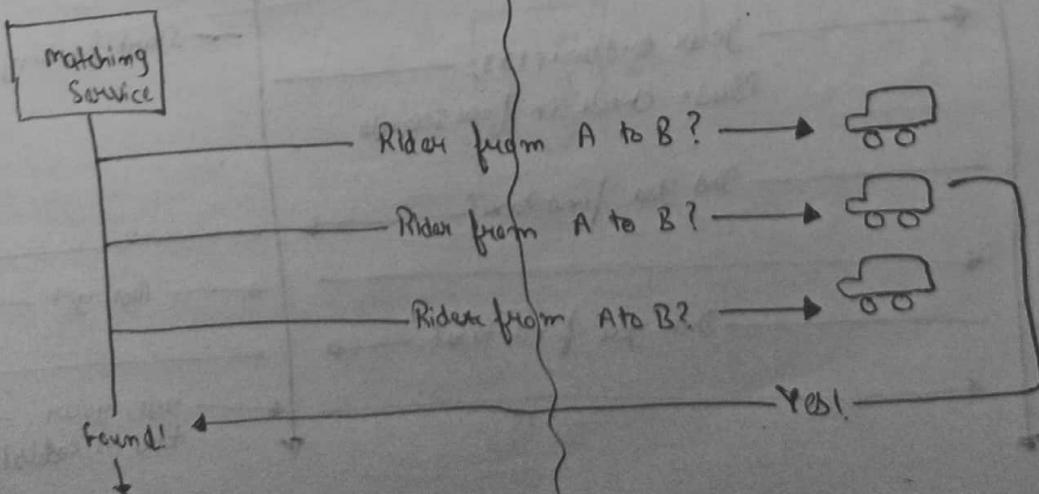


Approach #1 :-

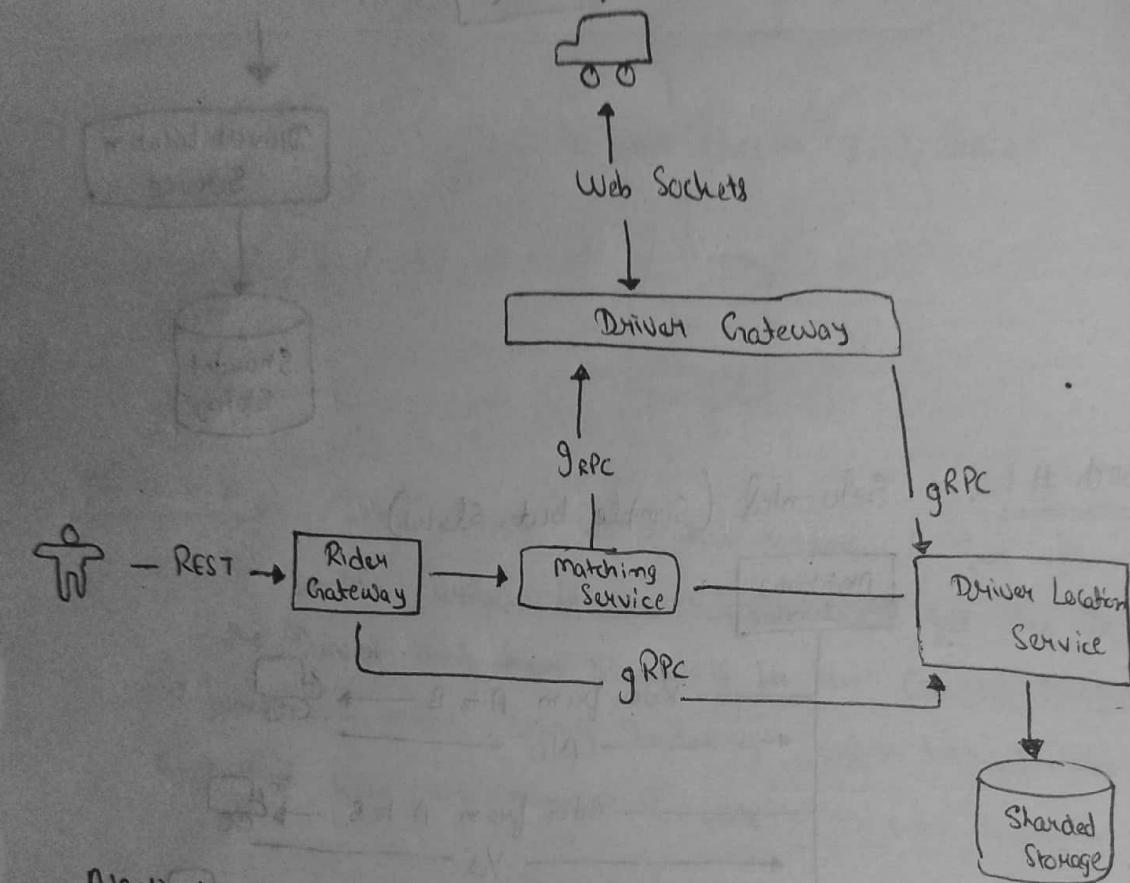
Sequential (Simple but slow)



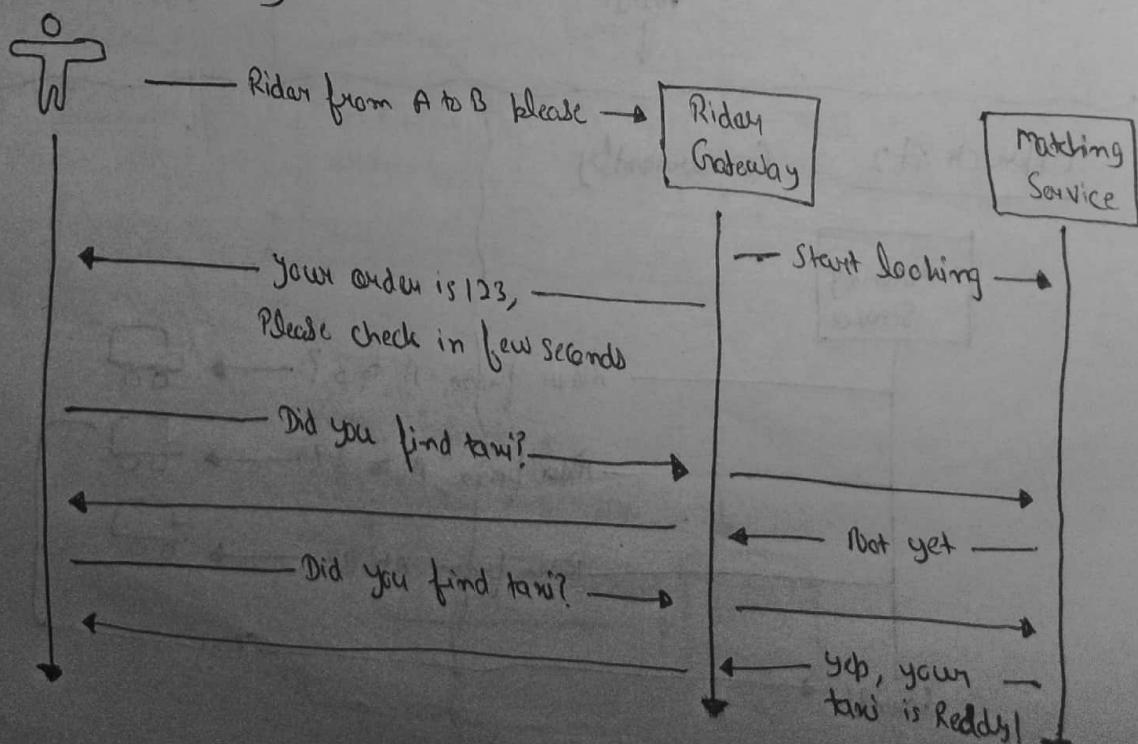
Approach #2 :- Concurrently

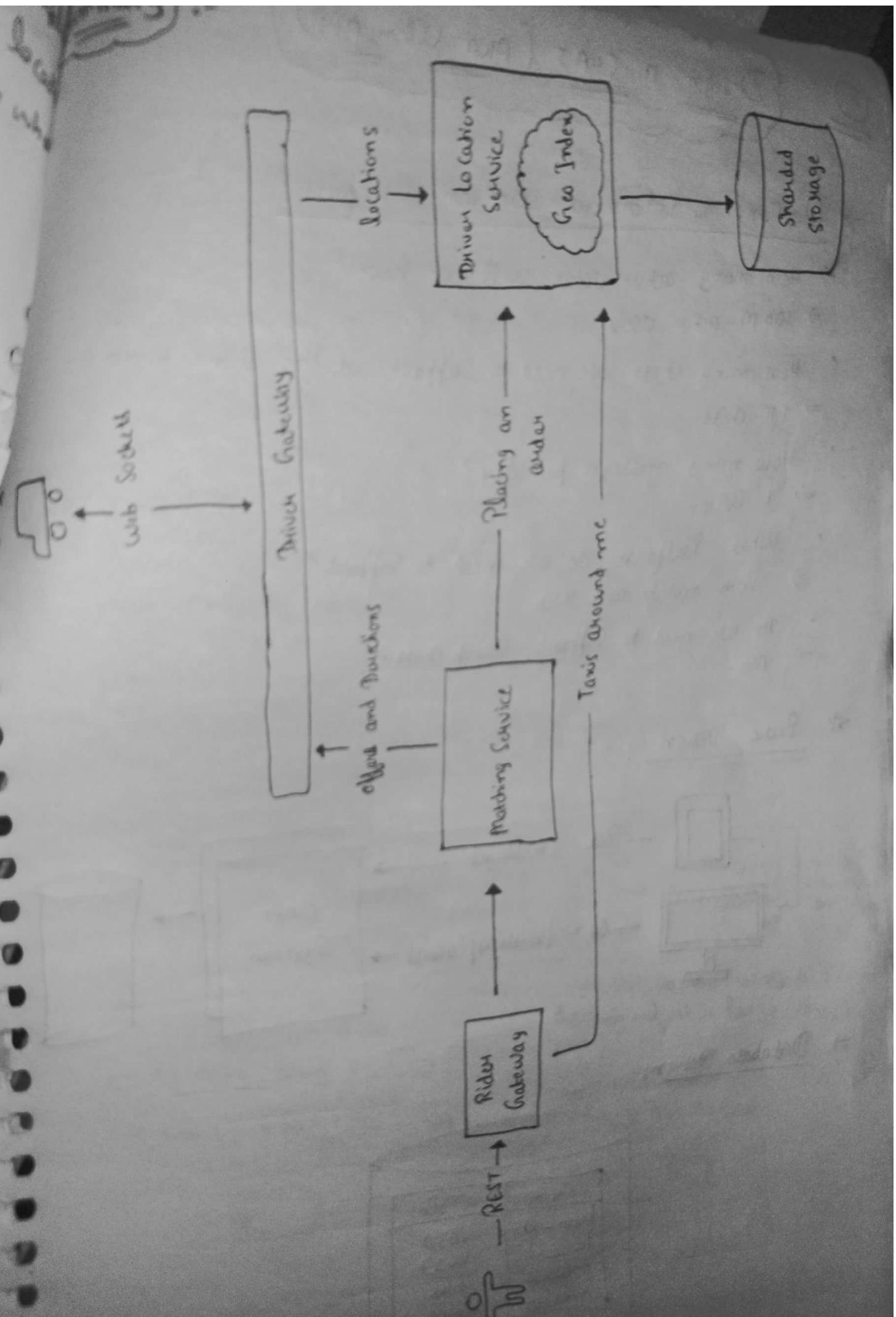


Driver Gateway:- Will use Web Sockets, bcz it provides the communication system, so we will get driver as well as give directions and offers to driver to pick the user.



Now, How user can get the updates of driver through Rider Gateway, it can be achieved





②

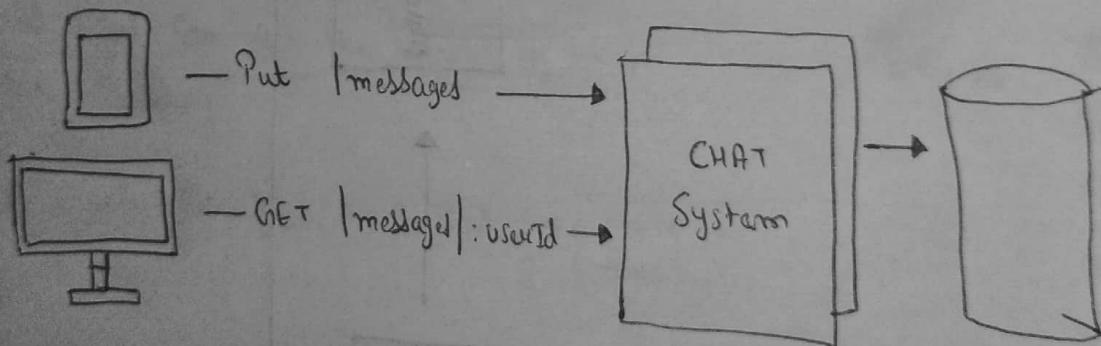
Design a CHAT (Aka WhatsApp)

Questions

★ Questions to interview we can ask first?

- How many active users will we have?
⇒ 100M active users
- How many chats we need to support at the same time?
⇒ 1m chats
- How many messages per day?
⇒ 1 Billion
- What Platform do we need to support?
⇒ Both mobile and Web
- Do we need to support group chats?
⇒ No

★ Basic Design :-



⇒ Database Design:-

Name	Type
from_id	UUID
to_id	UUID
message	text
Sent_at	timestamp

Select * from messages
 Where (from-user = 'Neo' and to-user = 'Trinity')
 OR (from-user = 'Trinity' and to-user = 'Neo')
 Order by Sent-at desc
 Limit 10

Indexing

Create index

on messages (from-user-id, to-user-id)

Select * from messages

Where (from-user-id = ? and to-user-id = ?)

OR (from-user-id = ? and to-user-id = ?)

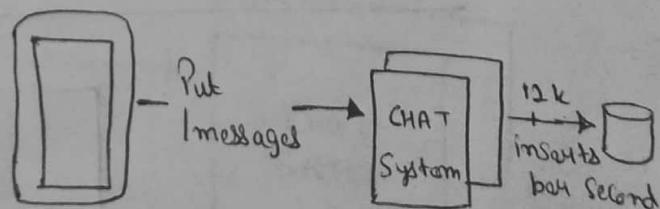
Order by Sent-at desc

Limit 10



Identifying Bottlenecks (Basic Requirement)

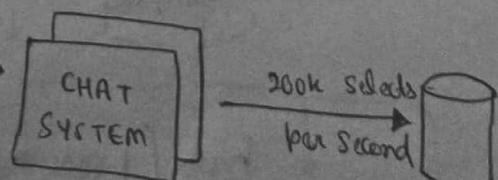
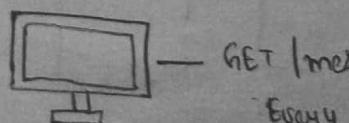
- = 1 Billion messages a day
- = 40 M messages per hour
- = 700k messages a minute
- = 12 k messages a second



- = 1 million open chats
- Refreshing every 5 seconds

$$= 1000000 * (60/s) / 60 = 200k \text{ Selects/Retrievals per second}$$

We can do this using a single database if it is large enough.



⇒ it doesn't goot, actually it is a pretty-bit, so we will look at it.

* Scaling Reads :-

Now, Earlier, we can see that there is a load on Database. So, to solve this issue we can introduce a Load Balancer and DB Replicas.

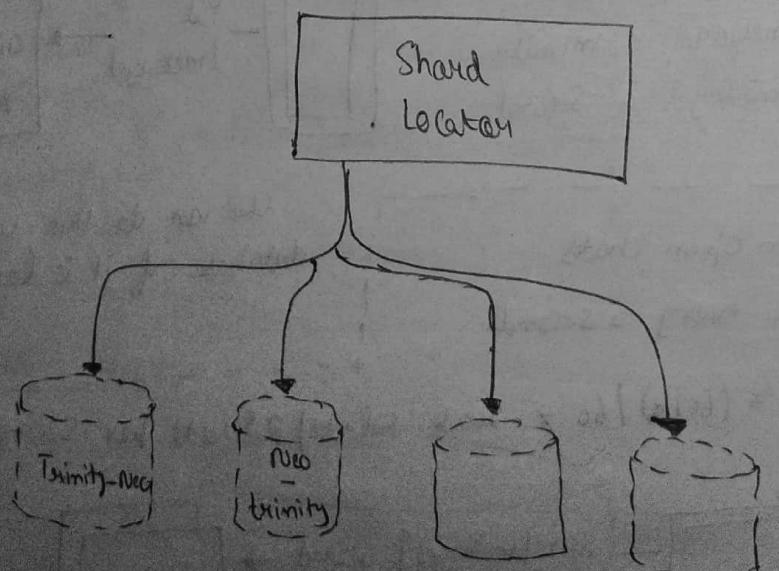
In this, we will put the DB Replicas behind the Dedicated Load Balancer.

But it also has 2 Disadvantages.

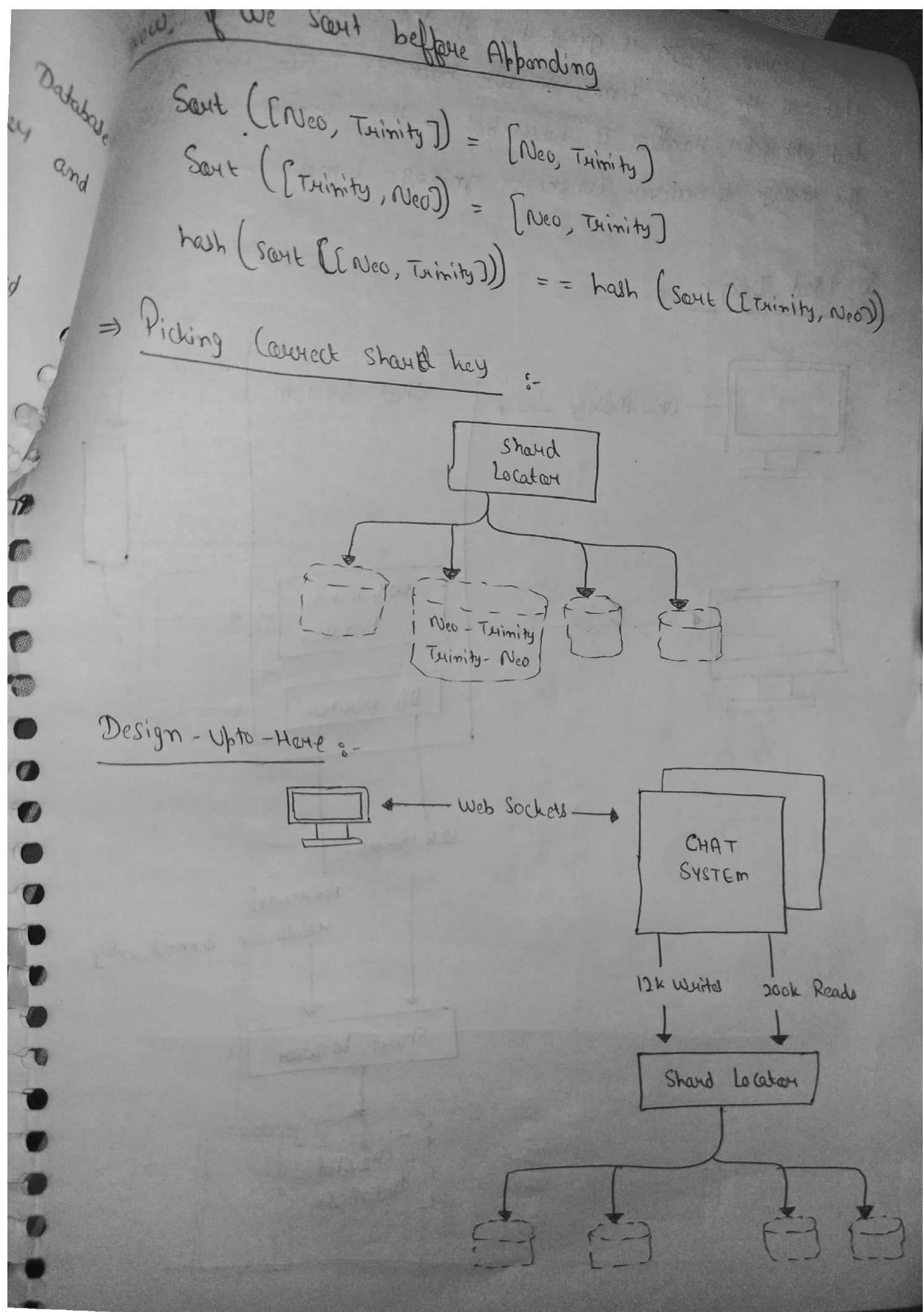
- (i) Costly, due to more Databases. (Copies/Replicas of same Database)
- (ii) As we want to deliver msg in every 5 seconds, but it will take around 1 minute bcz DB replicas may face some problems so, it will inconsistent.

* Sharding Chat Messages :-

⇒ Appending Sender and Receiver :-

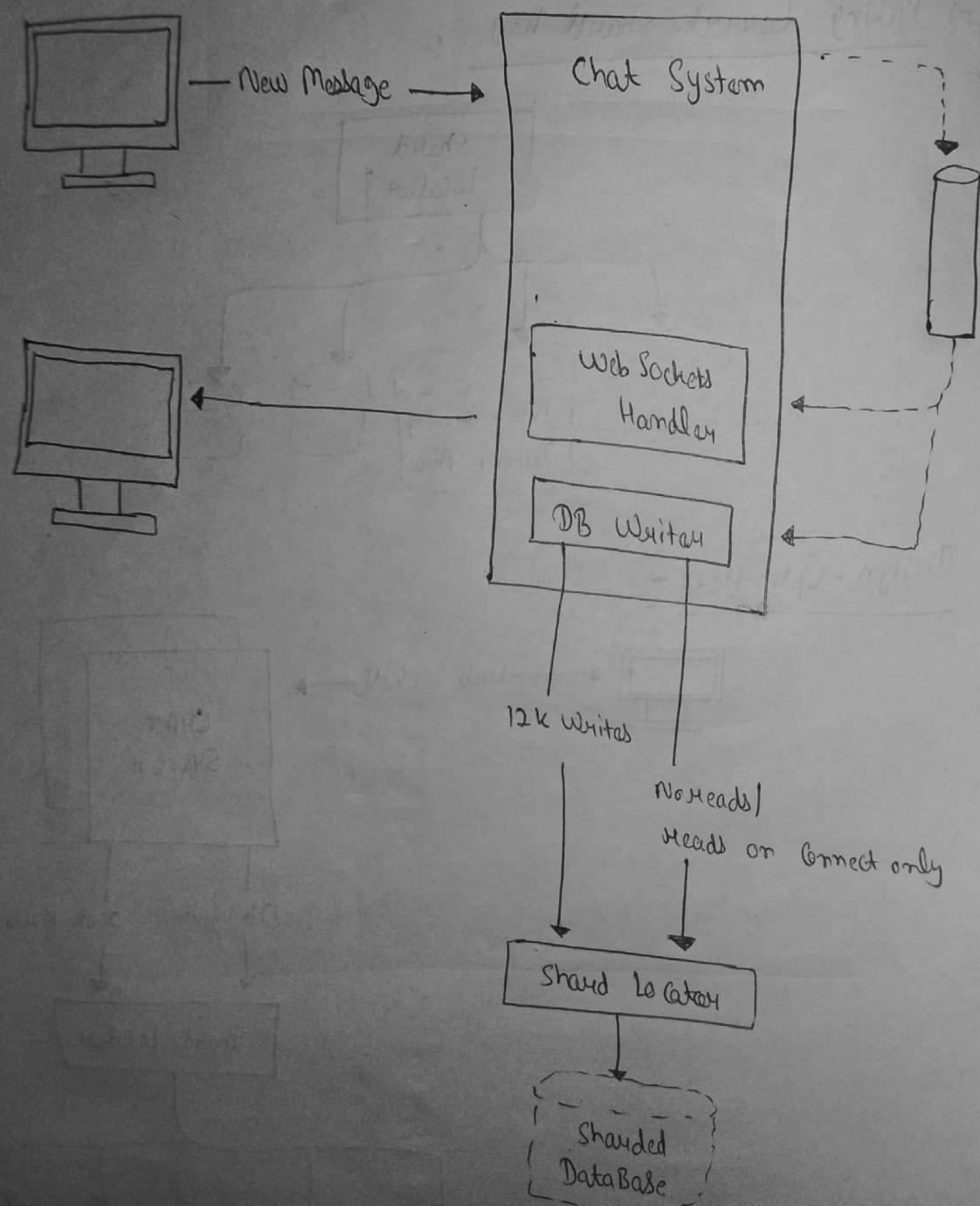


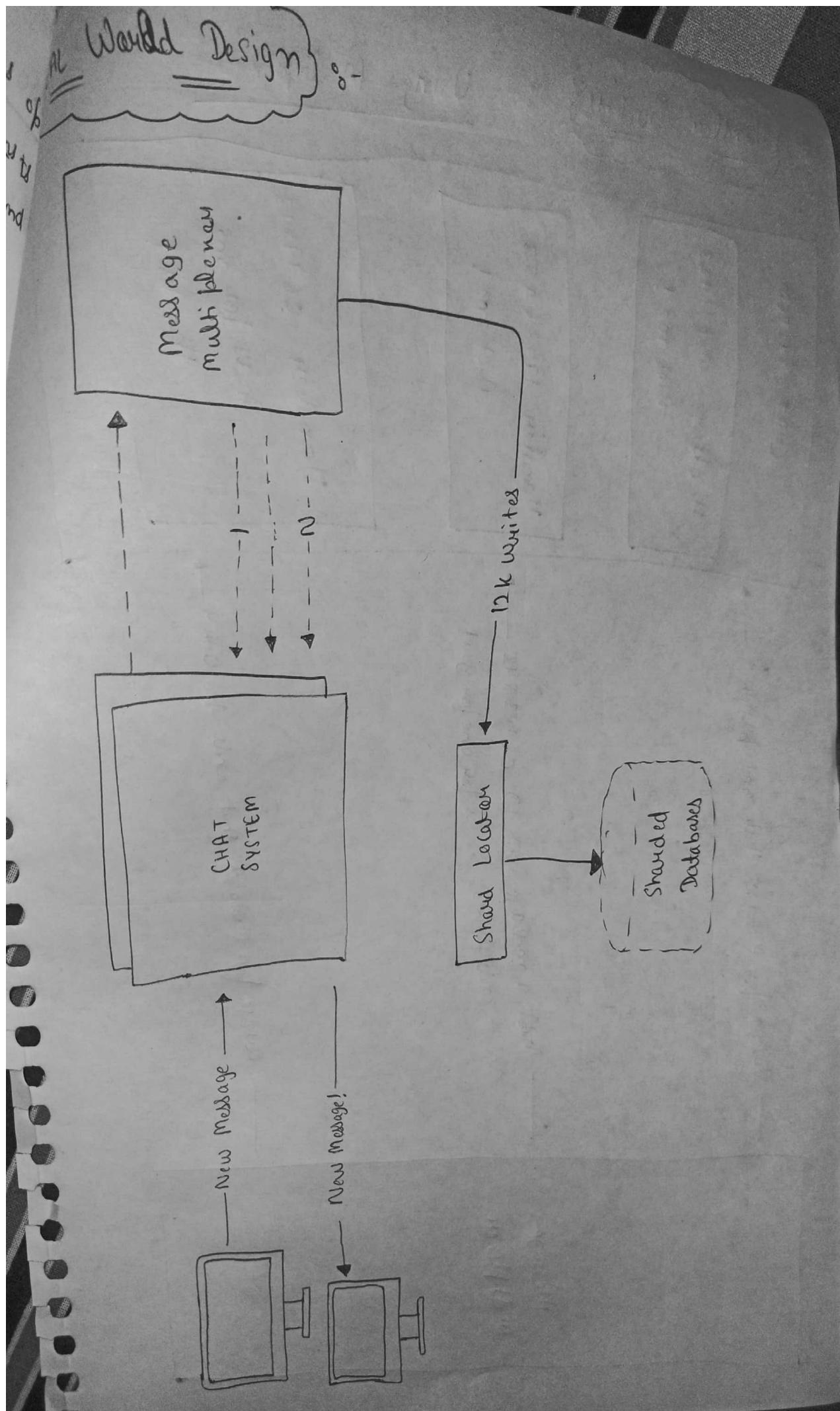
$$\text{hash}([\text{Neo}, \text{Trinity}]) \neq \text{hash}([\text{Trinity}, \text{Neo}])$$



- Our Previous Design is great but if there are million users at the same time, so we have to queue the message. Each web socket Handler is subscribed to that queue and will pass the message to receiver whenever message comes.

So, Ideal Design :-





A

Handler Buckets

:- Using Multiplexor

CHAT SYSTEM

Handler ab741423

Not me!

Handler 19936a53

Not me!

Handler c679842

Not me

Thanks!

— Queue for Messages that ends up with 2

Message multiplexor

— Hold a message with ID c6793843
and it ends with 3, So it goes one for you!

— Queue for messages that ends up with 4