

Design a chatting application like WhatsApp, Telegram etc.

Example of Chatting Application :-

- ① WhatsApp
- ② Telegram
- ③ Facebook Messenger
- ④ Slack
- ⑤ Teams
- ⑥ Slack
- ⑦ Web EX

Functional Requirements :-

- ① User can do one-one communication or can do group messaging
- ② User can see the old messages
- ③ Messages should be send in proper order.
- ④ Send message in form of text, image, videos etc.
- ⑤ Notification about the message
- ⑥ Message should not be lost, if recipient is temporarily unavailable.

Additional Features :-

- ① Message Encryption
- ② Delivery & Read status checks.

Non-functional Requirements | Page No:- 95

- (1) Real-time / Highly Performance, Low latency
- (2) Highly Scalable (3) Secure (Encryption)
- (3) Reliable

scale of this system

Total Active Users = 2 Billion
Daily Active Users = 500 million

Back of the Envelope Calculation / Infrastructure Estimation

Let, 1 person on an average send 100 messages per day,

Let average size of 1 message = 100Bytes

Total data for storage of data for 1 day

$$\begin{aligned} &= 500\text{million} \times 100 \times 100 \\ &= 500 \times 10^6 \times 10^4 \\ &= 5 \times 10^{12} \text{ Bytes} \\ &= 5\text{TB} \end{aligned}$$

Storage for messages in one year = $365 \times 5\text{TB}$
 $= 1825\text{TB}$
 $= 1.825\text{PB}$

Bandwidth Estimation

Note:- One read \approx One write
i.e. No. of reads \approx No. of writes.

Data generated in 24 hrs = 5TB

Data generated in per second = $\frac{5}{24 \times 60 \times 60}$
 $\approx 460\text{MB/s}$

Number of servers / Machine Required

$$\text{No. of Machines} = \frac{\text{Total Number of connections}}{\text{Number of connections per server}}$$

Since, Daily active users (DAU) = 5 Million

Let, To make a connection, server take 100Byte of data to make connection

Let, RAM of Server Machine = 16 GB
Also, if 10 GB is available to make connection :-

$$\text{Number of connection} = \frac{10 \times 10^9}{100} = 10^8$$

So, A single server machine is capable of making 10⁸ connections.

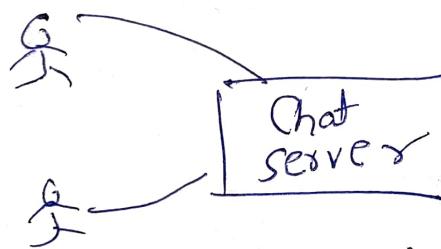
Note:- Chatting is not a CPU intensive task

$$\text{No. of servers required} = \frac{10^8 \times 10^6}{10^9} = 0.5$$

So, servers just accept the message and delays to another, hence we can handle 10⁸ connections in single server machine itself if it is capable to do that

(But we do not do this)

As it can lead to SPoF, single point of failure



Why using multiple servers instead of single Chat server which can do work alone itself

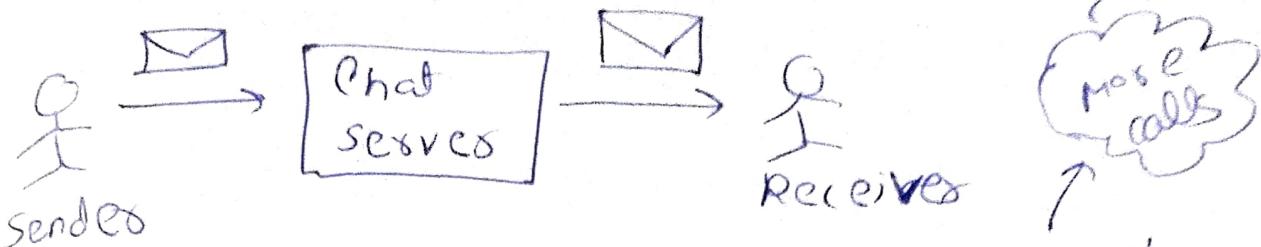
- ① Easily Scale
- ② Higher Availability
- ③ Locality Advantage

↳ means if a person is talking in India to another person in India, their connection is maintained in a chat server located in USA, so g will try that whenever he user communicates in different region their connection is made by chat server so that they are closer to each other because, g have to design low-latency application

Deep Dive into High-Level design

Page No:- 97

We need to talk about components and then communication.



If http connection are made in messaging system

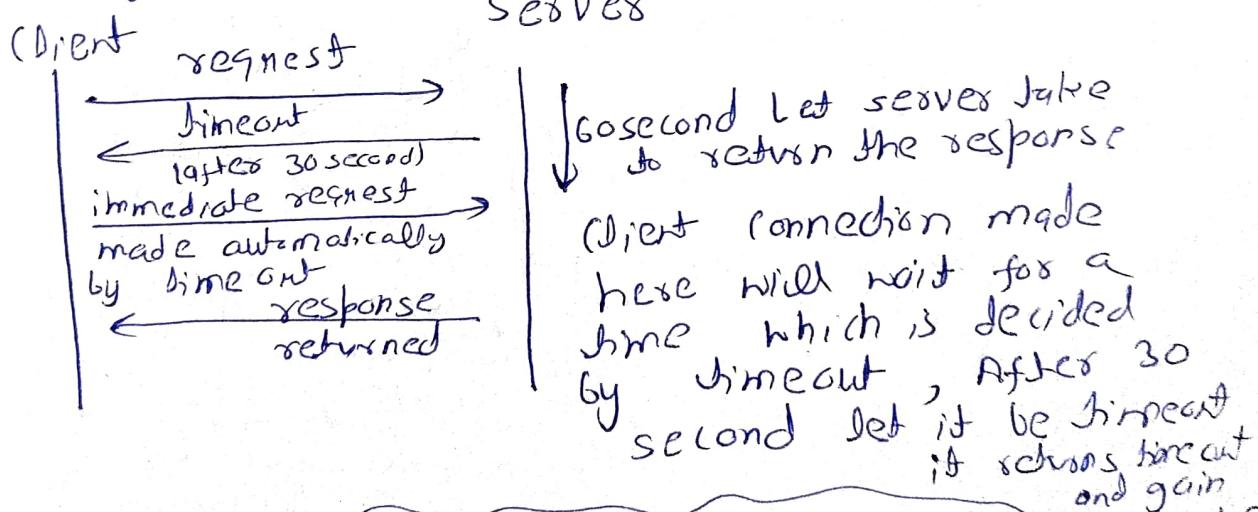
To connect the users :-

B/w client & servers , there may be multiple http calls that result in blank response if receiver has not sent the response or not available to send the message .

This strategy is called **Polling** in which you are calling the server multiple times , checking if servers has the data you are looking for and then getting the data back from the user

To make it optimized :-

Do Long Polling



Client connection made here will wait for a time which is decided by timeout , After 30 second let it be timeout if schools timeout and gain make

Long Polling helps in reducing the number of new http connections

immediately , if response is made by server , it return the response

Disadvantage of Long Polling

Thread is blocked.

Note:- If you are using multi-tasking, try to avoid long polling

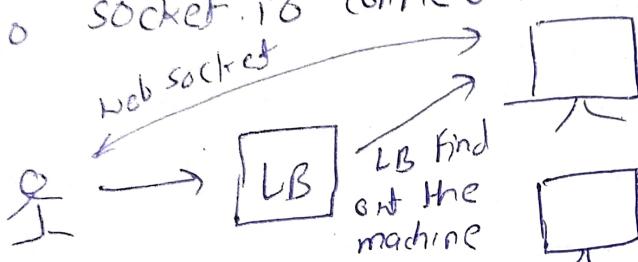
Important Point about Web-Socket Connection
communication also called statefull

① 2-way communication.

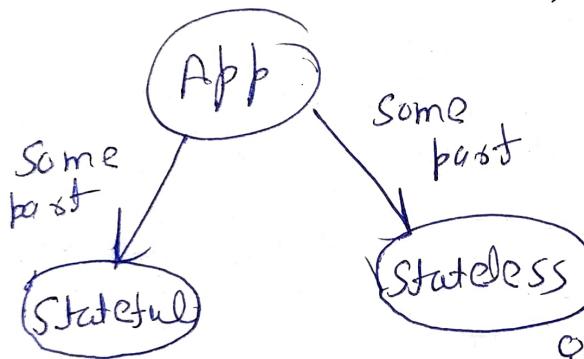
② Why Statefull Communication :- Since you are making connection

for a long period of time, we make consistent connection i.e machine is fixed for that connection so so socket.io connection are involved

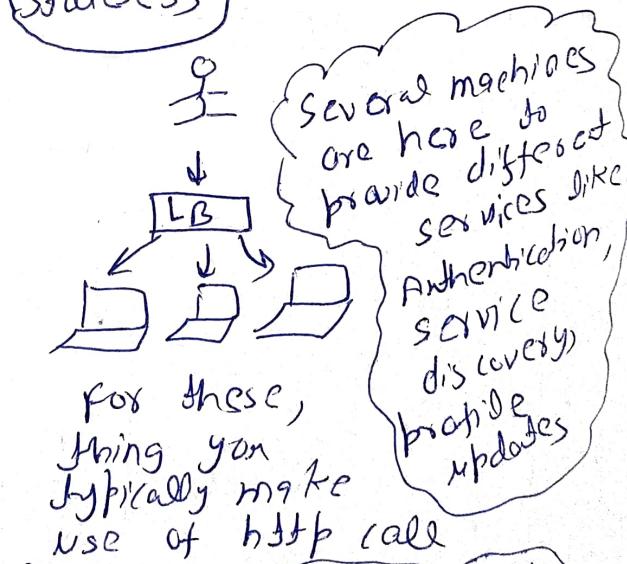
③



After finding machine, LB goes out of the feature and web-socket connection is done b/w user & machine/ chat server directly



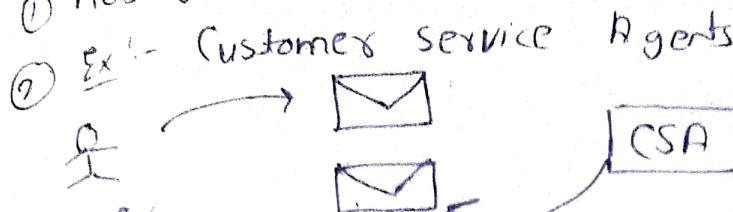
Text, img, video are also passed through web-socket connection



④ I make use of LB to manage the state of the system

Concept of Refresh Token / Access Token

① Mostly used in Banks



Let, customer do not go & send message for 2 minute
CSA become expired, all this happen from the concept

of Refresh Token / Access Token

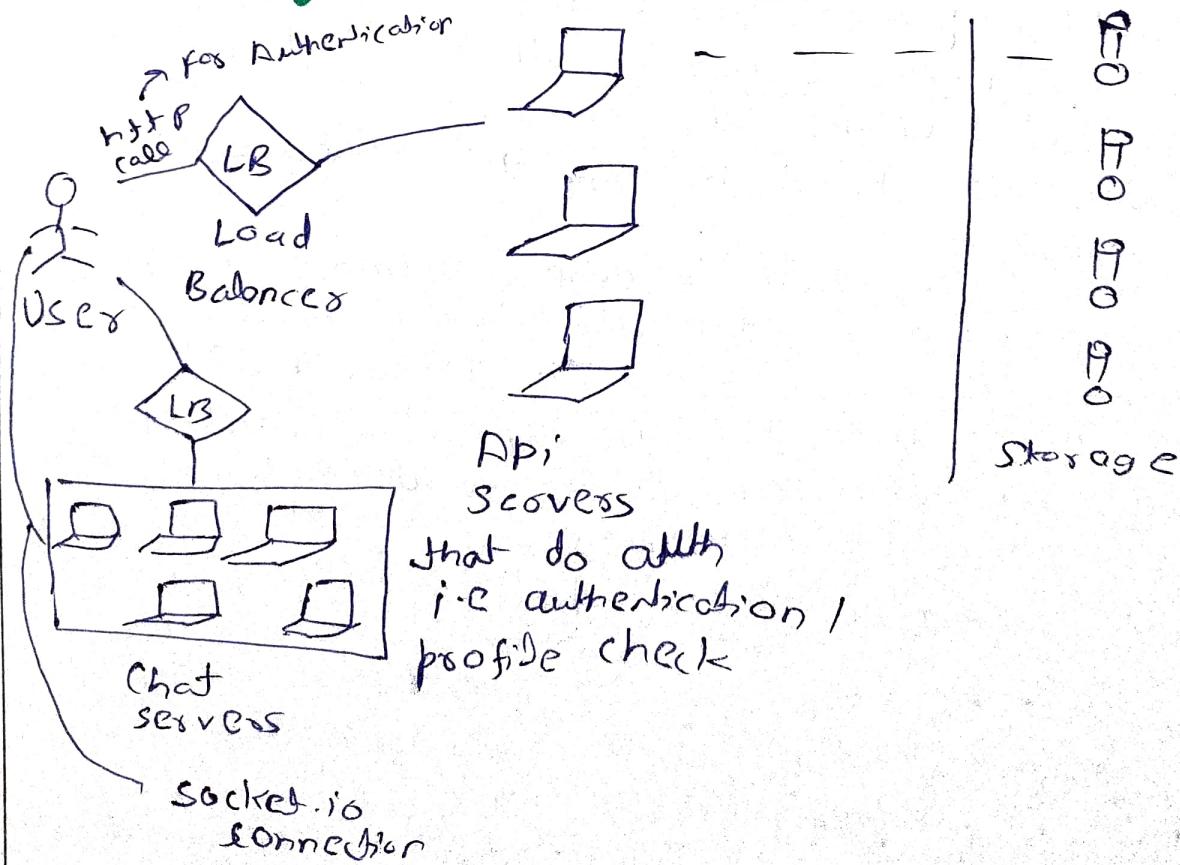
② Online Banking Portal like HDFC, ICICI, if you do not use it for few minutes, you will realize

that login is expired, you have to relogin again
why it happens:- Everytime you do authentication, you

will get a token, token has TTL +

TTL i.e Time to live) so if everytime you send a message, there is a new token that gets generated which has TTL say of 60 seconds, so within the 60 seconds if you don't send the another message, this token does not get refreshed and if token expires your whole session i.e the web-socket connection which was created was disrupted so, basically you are able to manage because of the refresh token

How things look at high level design



To store message

① Database should be scalable.

② Message :-

- (a) All message must be stored
- (b) Historical messages → Load rates
- (c) people watch the latest message

③ Schema of Message :-

- ① Message id
- ② time stamp
- ③ from
- ④ To
- ⑤ content

Make it as primary key i.e. Message id
 To make message unique in chronological order
 as unique identifier.

Message ID :-

- ① Should be made unique identifiers
- ② sorted based on time i.e. chronologically sort the messages

How to create GUID / UUID?

Used to create unique ID

Note! - RDBMS provide you feature of Sequence Generator that will generate unique ID for your new data. but NoSQL world, this feature is not available.

Since, data is huge we cannot rely over RDBMS, we have to dive into NoSQL world.

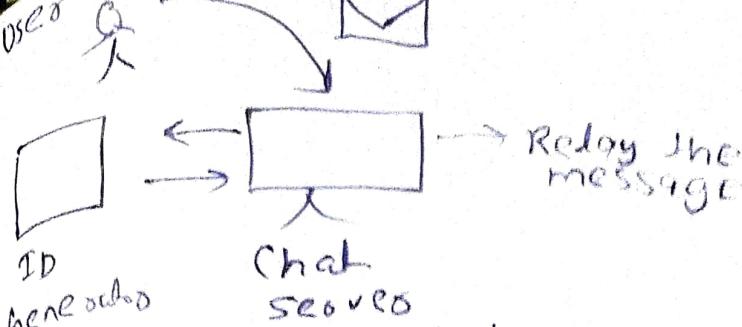
In NoSQL world,

We can make use of ID generator system that will do the same task as done by Sequence Generator in RDBMS.

How ID Generator System Works?

(a) This system uses Process ID, MAC Address, timestamp, random numbers to generate unique ID of message

(b) Each thing is of 16 bits. So total 64 bits are used to make unique ID



Note:- Why not used counter in place of the ID generator system is because inconsistency can be created & spoof lead to make our system less available & hence less reliable.

Cluster Manager:-

In distributed world, we use cluster i.e group of machines to do a task, we make use of load Balancers that know how to select a machine based on certain algorithm like Round Robin, sticky LB etc. but finding the right machine to perform task like finding right service to do chatting depends on multiple factors like how many connections are made by servers or maintained by servers, health, state, locality etc. of servers so, in distributed world, when you want to manage the cluster well Cluster manager comes into use.

What Cluster Manager can do:-

- ① Check if every machine is up or down i.e ~~health~~ health check
- ② It is able to do settings, so, use cluster manager to check how much load is taken by unique machine, how machines are connected to each other
- ③ You can do LB by cluster manager too.
- ④ You can also ensure that if one machine goes down it will bring up another machine in place of down machine.

Ex:- Zookeeper → One of the most common Cluster Managers to manage the cluster of servers.

See architecture of zookeepers & get information about it from Google or ChatGPT

Page No.: - 102

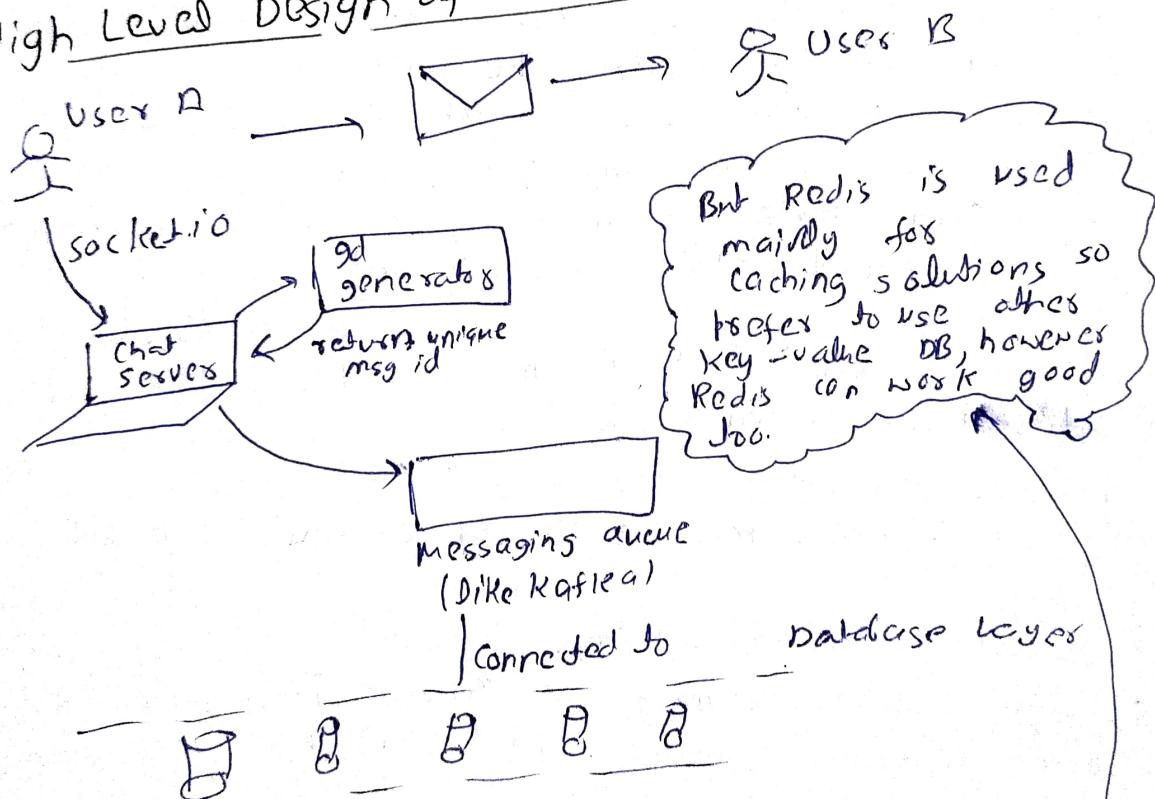
Containers:-

Example of Distributed System, containers are self-sustained process which can itself kill or run.

Kubernetes is used to manage the cluster of containers in your PC. Opposite to do task of Docker.

HTML CSS
JS CSS
If all are in some, it is containers.

High Level Design of Messaging System



But Redis is used mainly for caching solutions so prefer to use other key-value DB, however Redis can work good too.

Cassandra can be an option as data is humongous.

Key-value DB can also be used to store data like Redis because on the basis of msg id it can retrieve the content in constant amount of time.

Note:- Facebook Messenger uses Cassandra to store data

Note:- Where do you store User Data

User data is small, so typically you make use of RDBMS. Note:- User data is relational & small. That's why we are making of RDBMS i.e. MySQL. Why Not used MongoDB instead of Cassandra in messaging system.

Schema of message is fixed, it is not dynamic, since, Schema of message is fixed, it is not dynamic, flexible etc., we typically do not prefer MongoDB in that case. MongoDB is used where dynamic schema is available like.

In Facebook, posts has comments, comments can have comments or likes so here schema is not fixed, so in post we typically make of MongoDB to store post data. Posts will be stored in form of documents/JSON and hence MongoDB will be much better here.

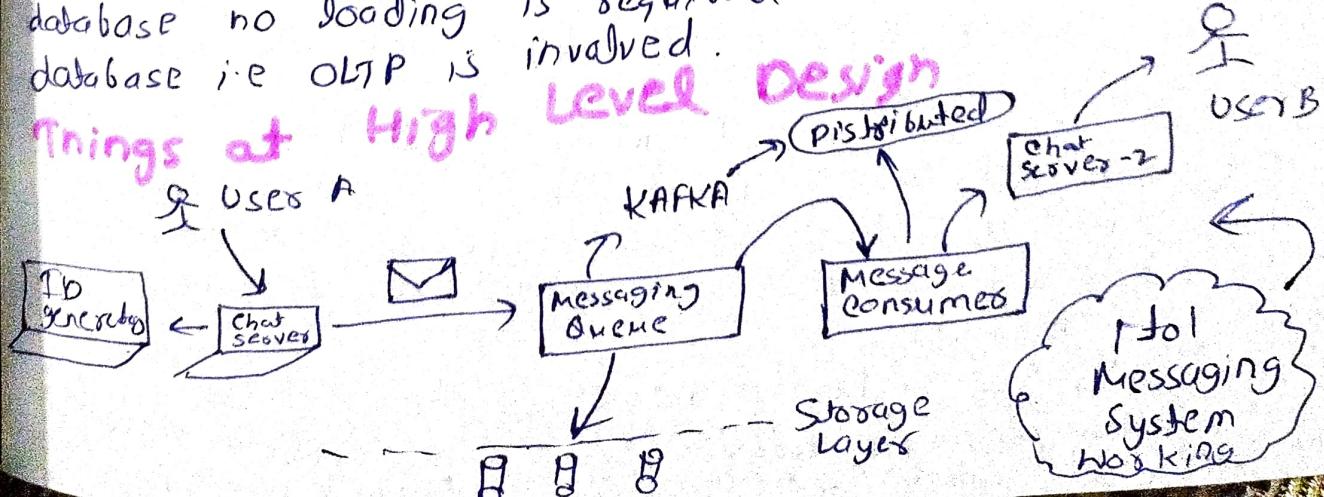
Video/Audio take time to load in messaging system while text take less time, why?

Note:- You cannot store or it is not preferred to store audio/video files in regular database, Distributed file system like cloud storage (HDFS, S3 etc.) are preferred to store audio/video file and their links are stored in regular database, so in

regular database, we store links of audio/video, hence you need to fetch content from link so it take time to load audio/video file.

Note:- In regular database no loading is required i.e. transaction database i.e. OLTP is involved.

Things at High Level Design



Note:- In one go, you cannot make a perfect system

Page No:-
104

Work of Message consumer:-

- ① It gets the message from the Queue and from message it knows which machine it is pointing now.
- ② This message consumer will have the information that to which chat server our receiver is connected to.
- ③ It will make call to chat server-2 and send it the message.

How to ensure that our bandwidth of sending the message is not choke

To prevent choke of bandwidth of sending the message we can put some restriction that any message which is to be sent must be not greater than 20 MB. Like 1 MB size you cannot send a message greater than it at a time.

Identifying weak points of HLD:-

- ① Consuming from KAFKA → All message are sent to Kafka.
- ② Find which chat server to send/pass message

Improving KAFKA Topic

Since, you are system designer, so you can choose topic for each user. If I send all messages of different users to same topic, things get tough for message consumer to do task.

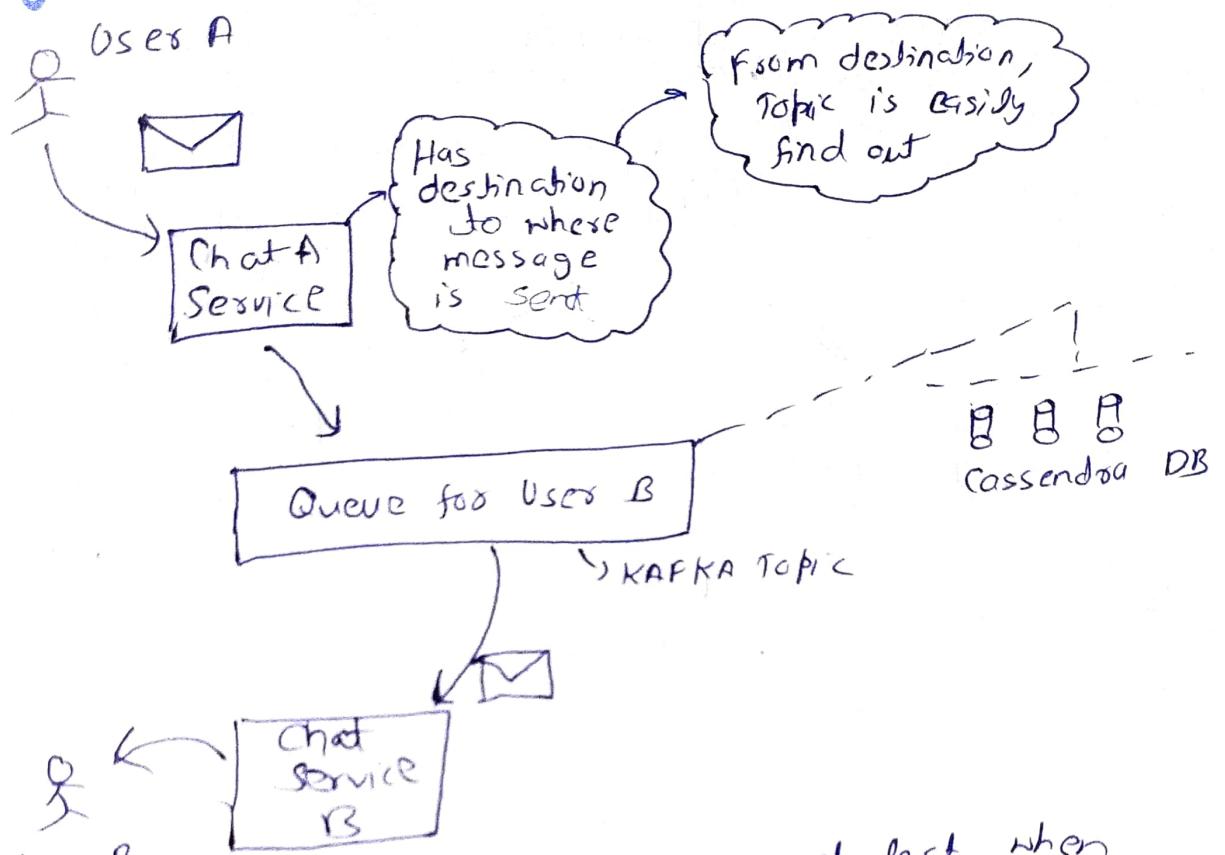
So, ensure each user is given a unique topic so message consumer will not have search for user to show its message it directly call to that unique topic, get the message and display to the user.

Normally, I have to ensure that for each user when they register, they should be assigned message queue, to be decided for that user. Message queue can be considered like a Mail Box.

Now, this leads to scalable solution as no message consumer will be required, chat server can do work of message consumer.

Final Design or Improved Design of Messaging System

Page No:-
105



When User B is offline, message is not lost, when he/she becomes online, chat servers to which User B is connected check that its topic has some message if yes, takes message and display it to User, hence if UserB is offline message is not lost also, no message consumer will be required as chat servers has replaced task of message consumer. This has led to more scalable design.

How To grab each question of System Design Easily

To Think Solution

Step :-

Start thinking on single machine, how this problem can be solved using single machine

Think about components

How components should interact with each other
Which database should I need?

Note:-

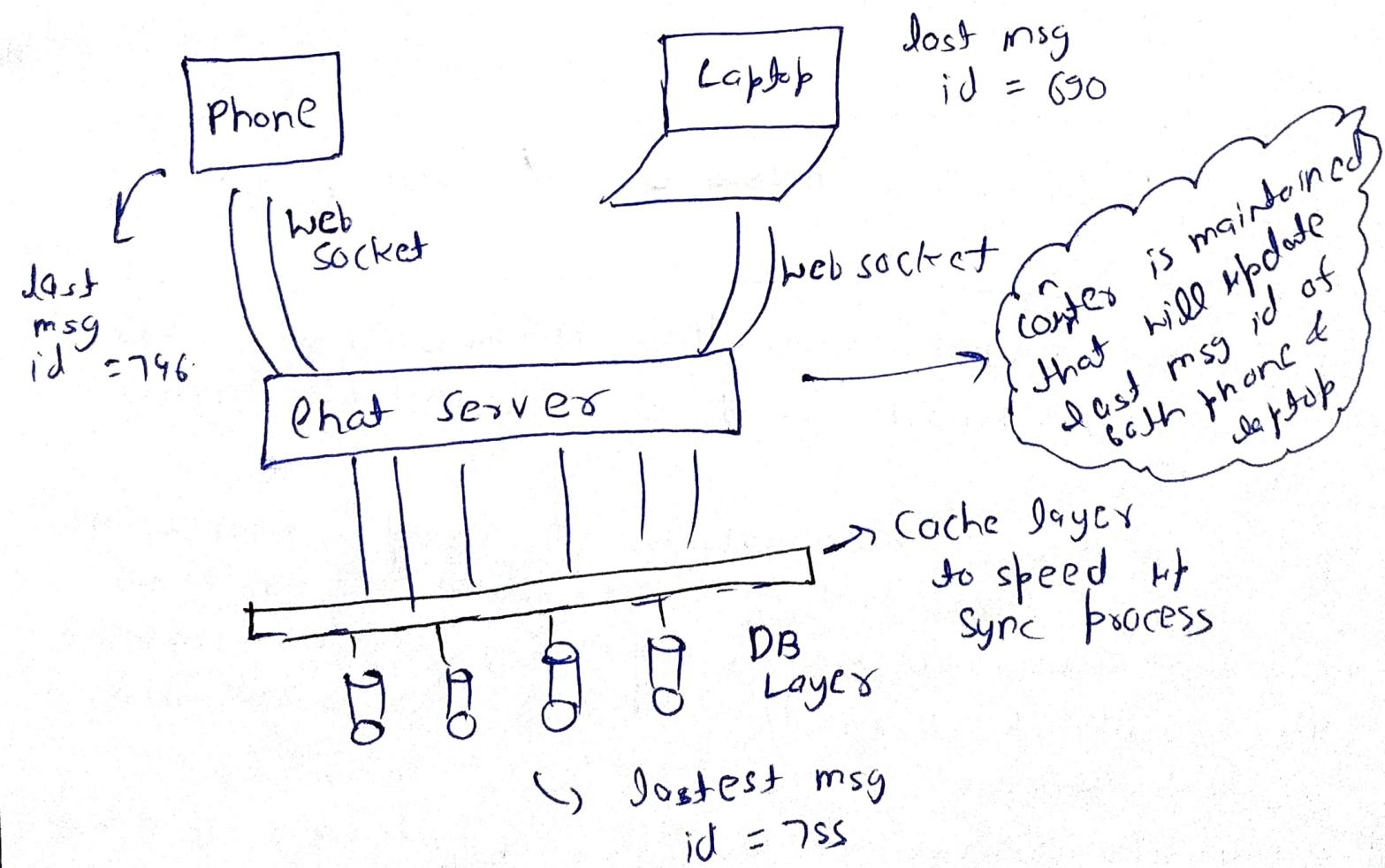
Polling and Long Polling connection are achieved with the help of http connection.

Additional features / Enhancement of messaging system

How to ensure that each device has the latest data / synced

Assume, user has laptop, mobile phone. Now, if he use Mobile and laptop remains off to use messaging system, we need to sync the activity of user to both devices so that this might not happen user get stale data at laptop.

More Modular design



Back up feature in Messaging System

Page No:-
108

System

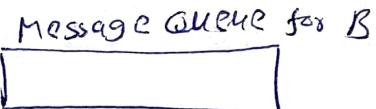
Backup is maintained in local device and cloud storage. These backups are for the seen data.

Why local cloud storage not Central server DB to do backup.
Reason is that local cloud storage is faster than backup.

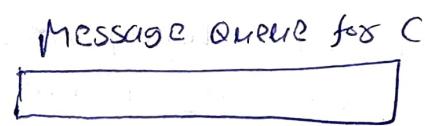
How to design Group Chat?

① Every person is given its own message queue

Q A



Q B



Q C

Works like an inbox mail for its user

Why in messaging system, group size is made finite?

This is because.

Number of events/ message & group size published

For ex- We Chat is another messaging application in which group size cannot be more than 500.

Note:- Message also consist of group information i.e. Message consist of that it will of one-one communication or group communication. And on the basis of that chat servers send msg to the required queue.

② How my chat servers decide that which msg in group should be pushed to which msg queue of member?

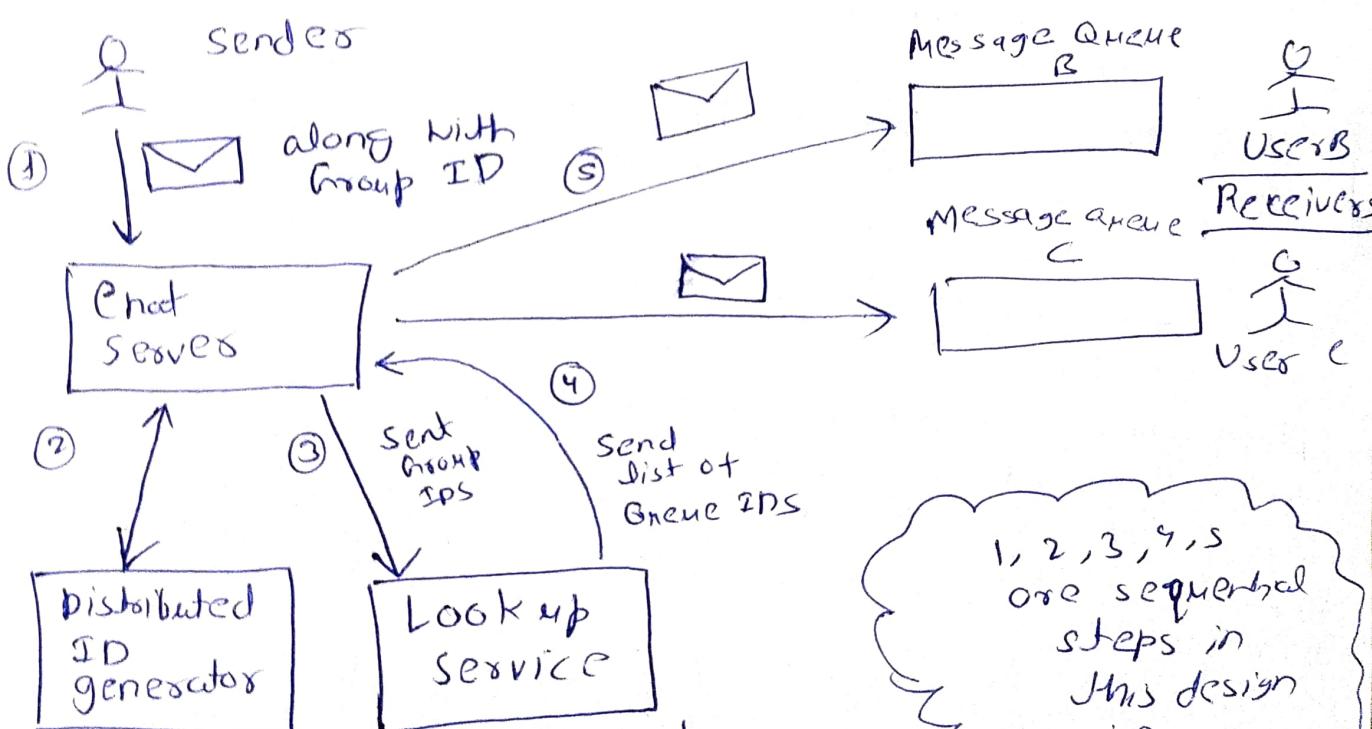
on system, there can be multiple group like G_1, G_2, \dots, G_n .

But every group is assigned a unique id and we can create some kind of map which consist key as 'Group ID', value as a list of User IDs.

By maintaining a map data structure where in the **constant time** you can do **look-up** that what are the different users in the group.



More Modular Design



maintain map data structure in which Group IP corresponding to list of user IDs are stored. Based on that User IPs if return list of Queue IDs where messages need to be pushed.

1, 2, 3, 4, 5 are sequential steps in this design i.e.
 $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$

How to reflect the online/offline status of user to its friends

A system called as Online status server is added which keep an eye or monitors the web-socket connection made b/w chat server and the user.

Work of the Online Status Server :-

It is going to check that what are the active web-socket connections, by the active socket connection it comes to know that what are the two sides of socket connection, i.e. it gets the User-ID.

And for this User-ID in persisted storage, for each of the user a flag named status can be maintained that will tell that user is online/offline controlled by online status server.

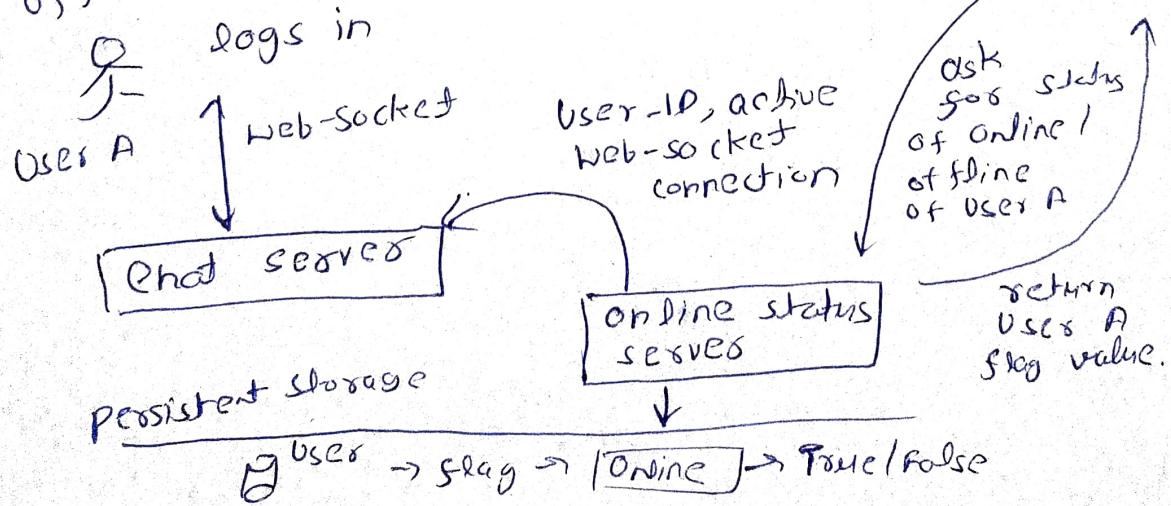
So, whenever user logs in, status is set as online.

User flags:-

When there is boolean outcome, we typically make use of flag.

How do user check his/her friend are online or not

Make call to Online status server by giving their User-ID, OSS send request to persisted User DB to ask for User-ID, if it is true OSS return online otherwise offline.



Way to ensure if a system is alive

PAGE NO:- 111

Let, we have a server and in it some application is running. Now if I have to go and check that my server is live or working or it has gone down.

Then we have to check that the health of our server/application is good or not.

If it is healthy we say that our server is alive.

For this, we make use of something called as Heartbeat Test

Heartbeat Test

In this test, we make a server which is used to test the heartbeat.

By this Test, the machine/system after every 1 second has to tell about its health to the server that dude I am healthy. just like we have to tell our health status to our mother/father regularly.

After 1 second, if no heartbeat is sensed by Heartbeat Server from system, it means that our service has gone down.

Why are we focusing on Web-socket connection to check status of online/offline. This is because when you close WhatsApp window (still it is running in background), your web-socket connection breaks up so keeping track of web-socket connection provide you **real-time awareness of that user is online/offline.**

How to see last active time of user:-

Everytime when web-socket connection is broken/terminated, we have something called log so, online status server (OSS) has two tasks it gets a list of all web sockets which is active and it also sees that the web sockets which were terminated, what was the last time it has terminated. By OSS you check:-

- ① User is Online/Offline \rightarrow achieved in almost real-time
- ② Last seen feature

Note:- sometime , it can happen that system cannot send message after α second . That's why you make checkmark that after y second if my system does not send a single response , i will consider that my system has gone down.

Ex:- $\alpha = 5$ second , $y = 30$ second.

Note:- Heart Beat sever can contain stale information but it is useful when you work in clusters environment like Big Data Cluster, kubernetes.