

Consistency + Partition Tolerance
(Availability Compromised)



A



B

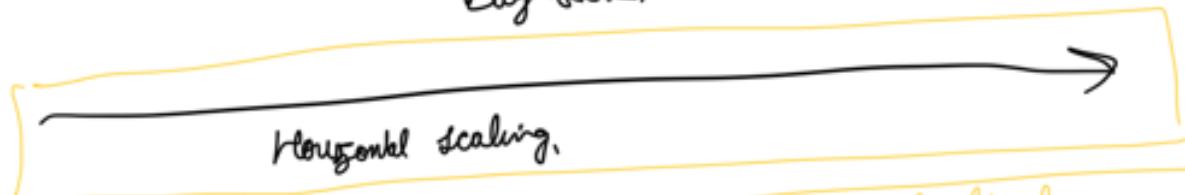
Availability + Partition Tolerance
(Consistency Compromised)

System Design Lecture Nikhil Lohia

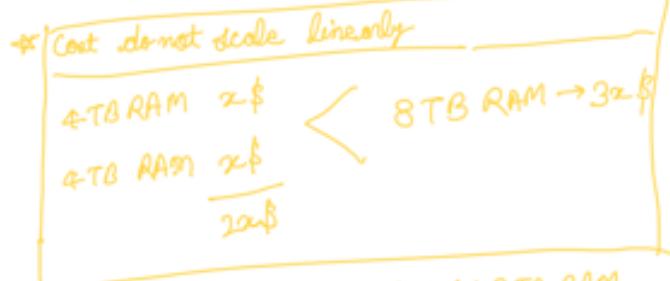
Nikhil Lohia →



What is Scaling → A way to expand and accomodate more customers.



Creating a Basic Application →



cost of (2) 4TB RAM < cost of (1) 8TB RAM
that's why we use horizontal scaling



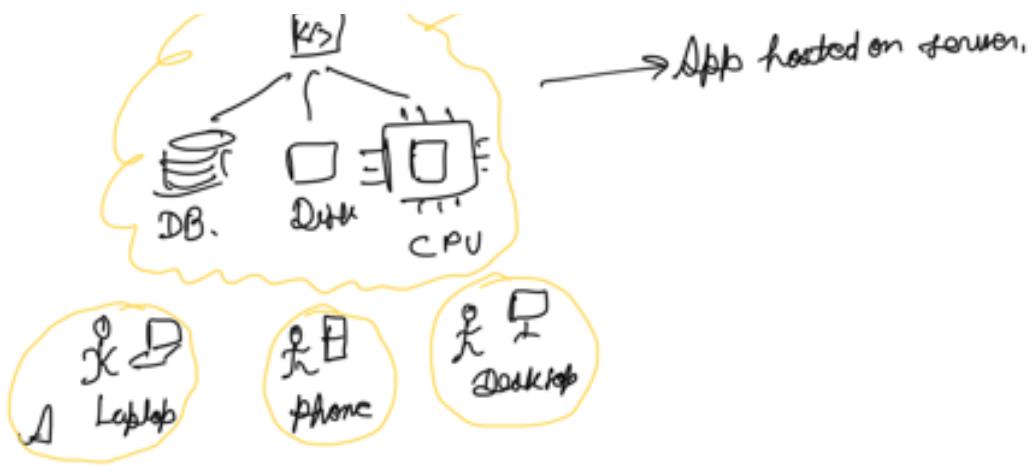
Course Structure →

- Basic of system Design
- going overall terminology used.
- start with Basic design.
- Understand each component and use cases
- explore some complex design.
- Based on how it goes, look at some real system too.



Client-Server Model →

→ server



Advantages →

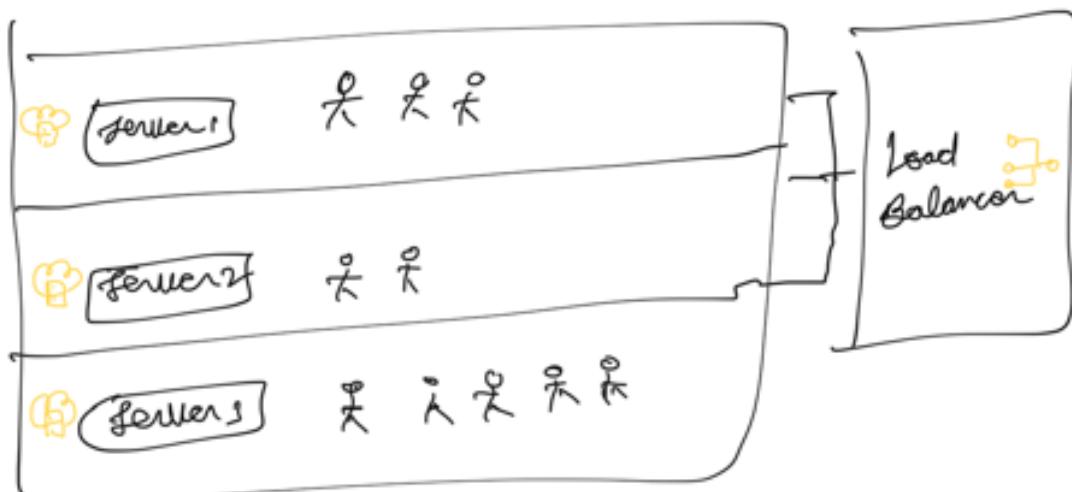
- ↳ central resource management.
- ↳ scalability.

Disadvantages →

- ↳ single point of failure.
- ↳ Network Congestion.
- ↳ security Risk.

③ Load Balancer →

Balancing a lot of request for a particular server.



Types of Balancing →

- (i) Least Connection.
- (ii) Round Robin Scheduling
- (iii) Weighted Round Robin.
- (iv) IP Hashing.

Advantages →

- Traffic is distributed and processed parallel
- Your design becomes easy to scale up.
- High availability.

Disadvantages →

- A load balancer can get very complex and not easy to set up.
- If load balancer goes down → The load balancer itself becomes a single point of failure.
- If a load balancer is compromised, system can go down.

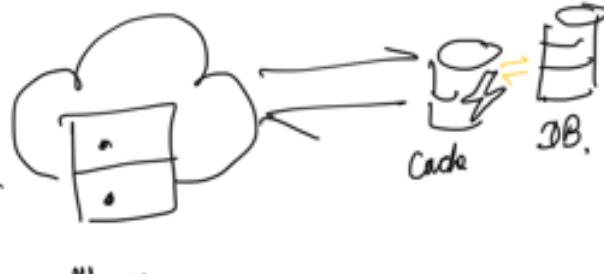
④ Caching →

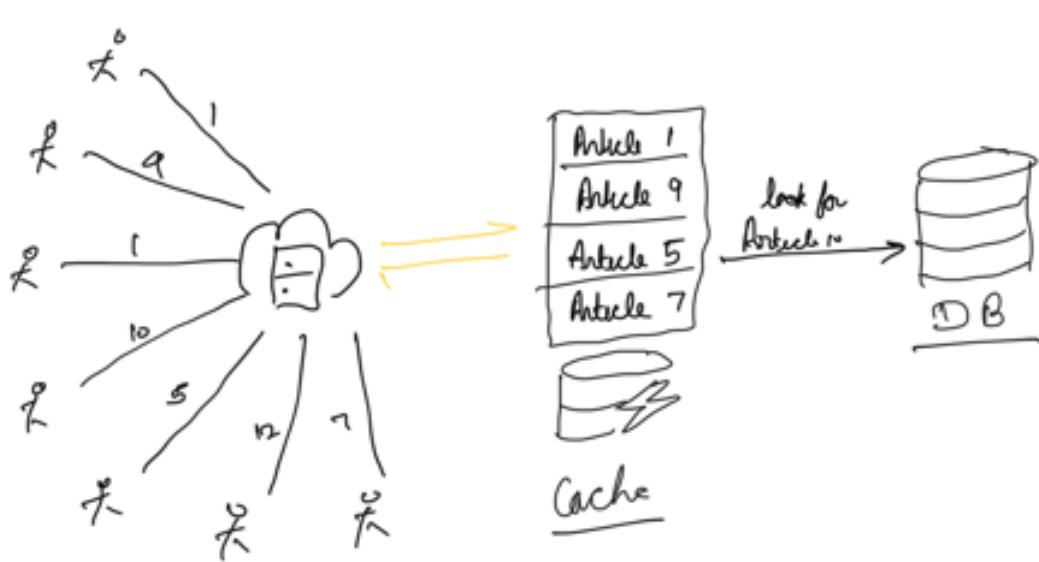
- Increase efficiency of server
- Saves time of client.
- Increase speed of transaction

Symbol of Caching →



- Kind of a memory / small DataBase
- Very fast.
- Usually small in capacity.





→ Higher the value of $\frac{\text{Cache hit}}{\text{Cache miss}}$, better is the system.

Cache Eviction Method →

- ↳ Random Eviction
- ↳ FIFO → (First in First Out)
- ↳ LFU (Least Frequently used)
- ↳ LRU (Least Recently used)

Things to be careful about →

- ↳ consistency of resource → प्रत्येक article के बारबार Field का लकड़ा हो और दिनों में article change कर सकती consistency effect होता
- ↳ Coherence: Cache in different regions could get out of sync.
- ↳ Security: Sessions and credentials should be cached with caution.
sometimes login credentials get cached तो कैसे होते हैं?

⑤ Databases →

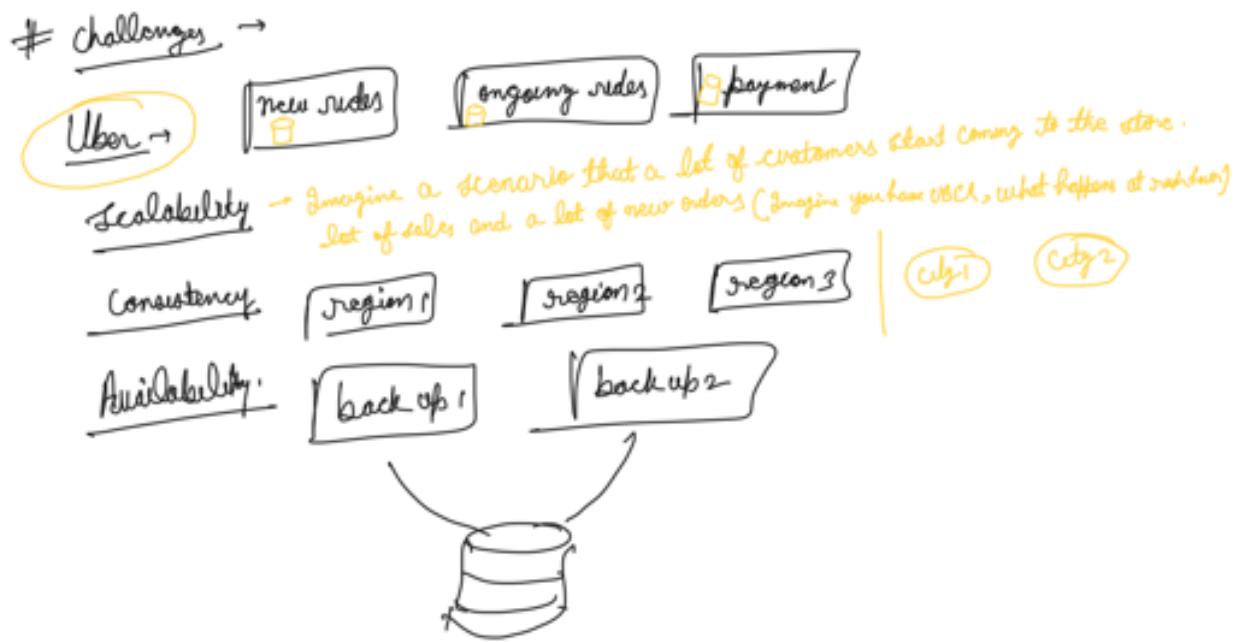
→ Efficient Data Retrieval

- ↳ Quick queries for any information that you need. You don't need to scan entire storage.

`Select name, marks from students. class 6`

→ Data Integrity Or Consistency: All operations make sure that records remain accurate.

→ Handle Complex Queries → Find out performance, usage metrics, correlations, advanced insights, data mining.



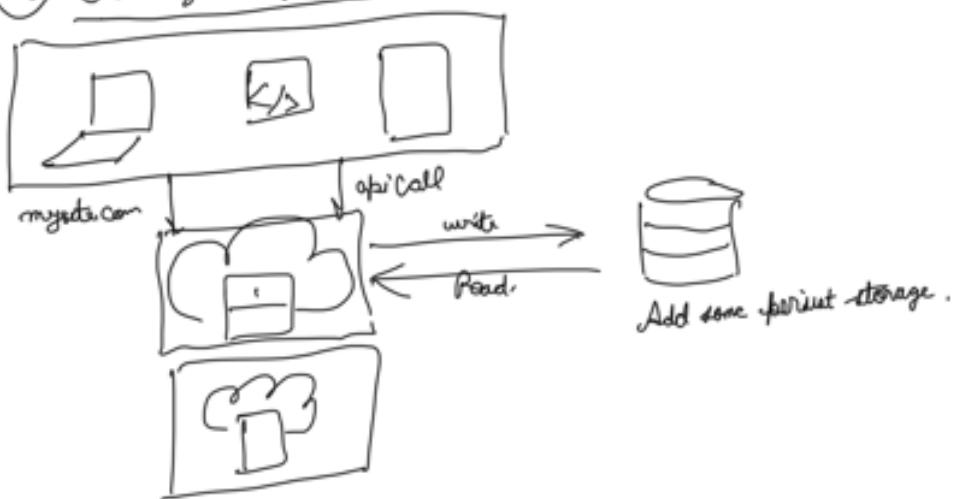
SQL vs NoSQL →

date	transaction	credit	debit	balance

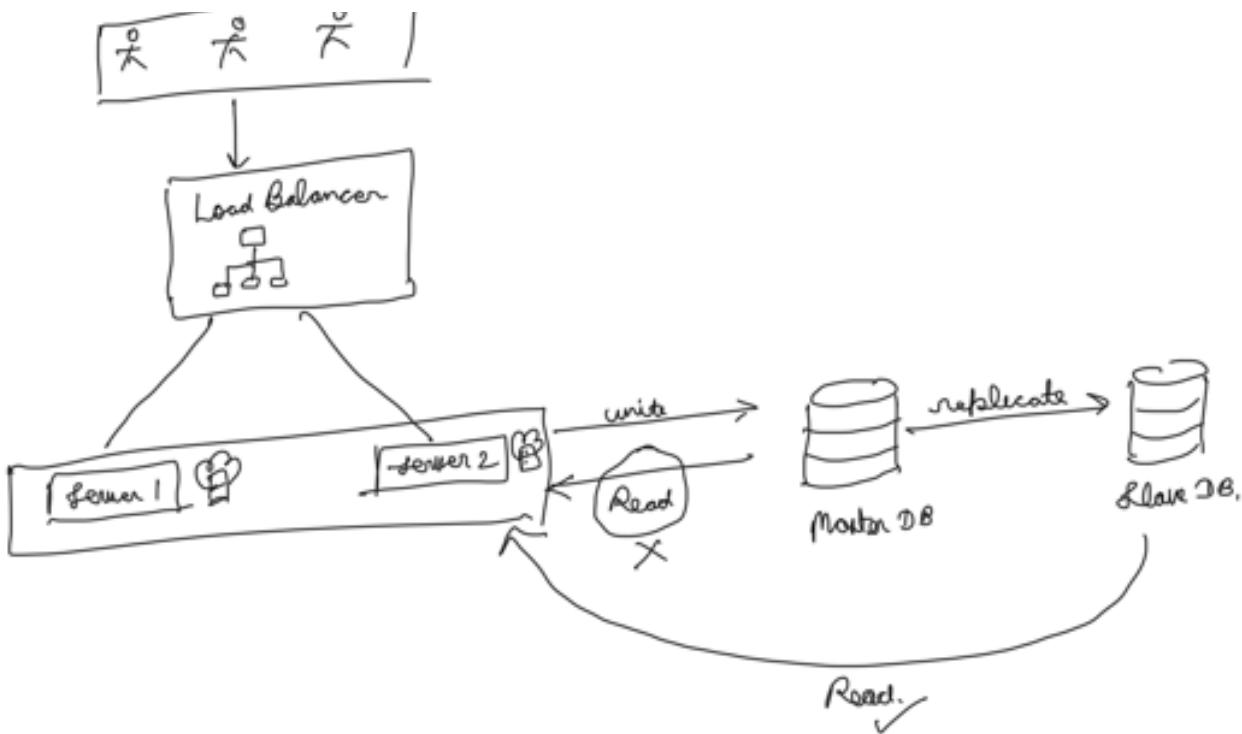
SQL → PostgreSQL, MySQL, Microsoft SQL Server

NoSQL → MongoDB, Cassandra, Redis

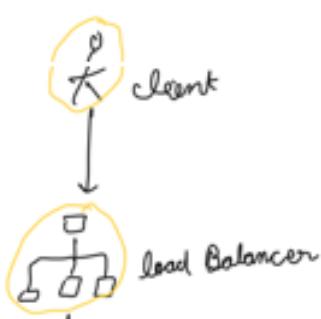
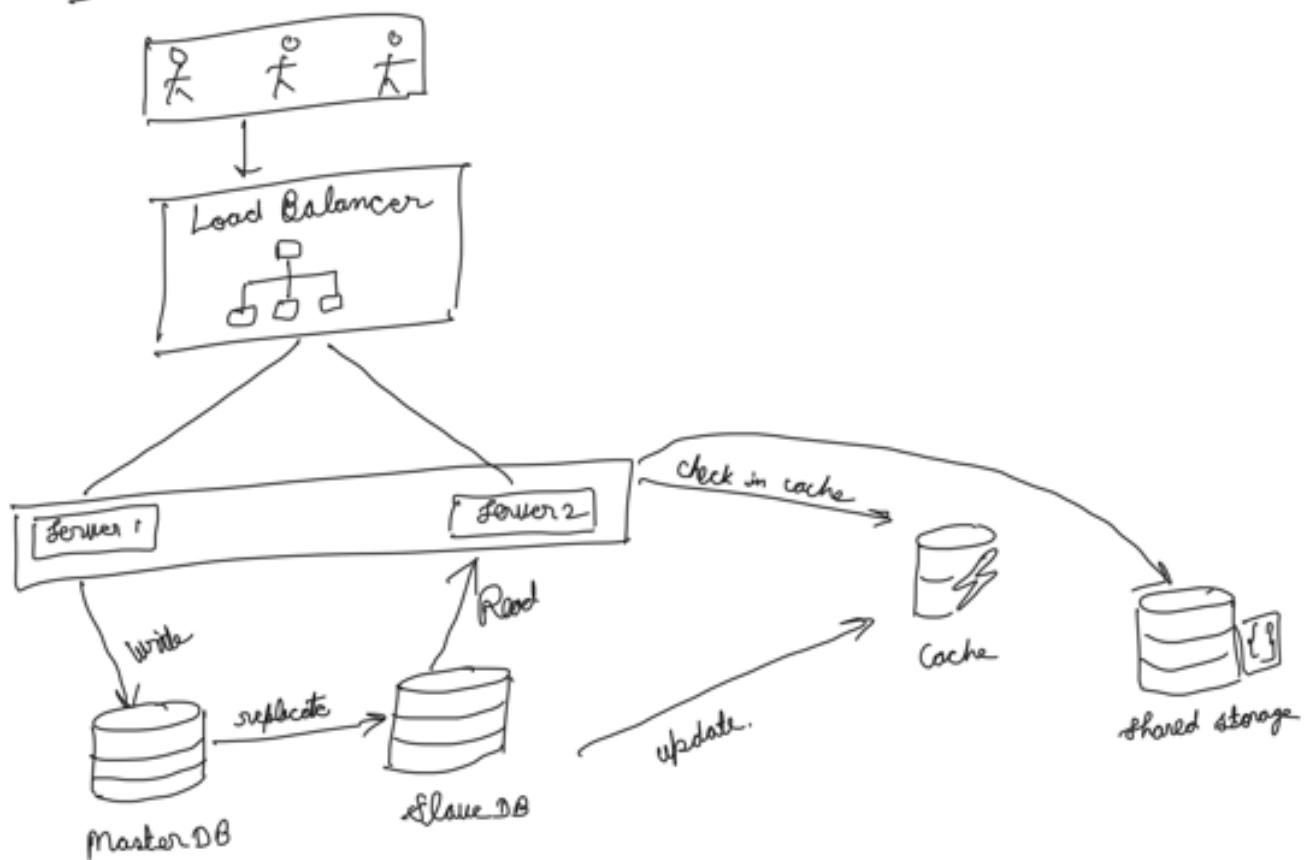
⑥ Creating Diagrams →

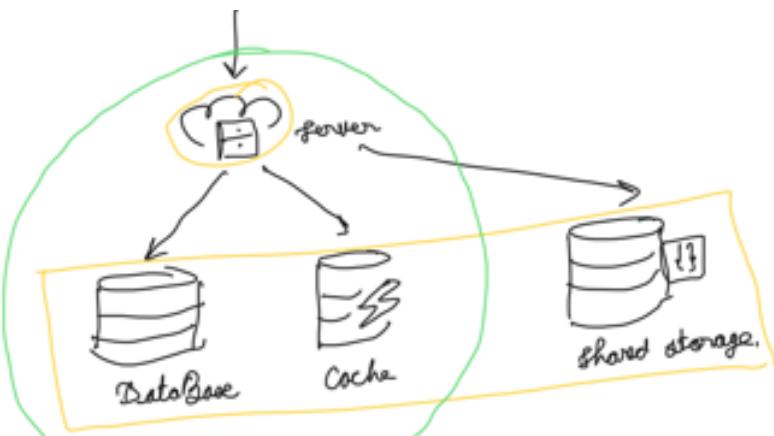


Load Balancer →

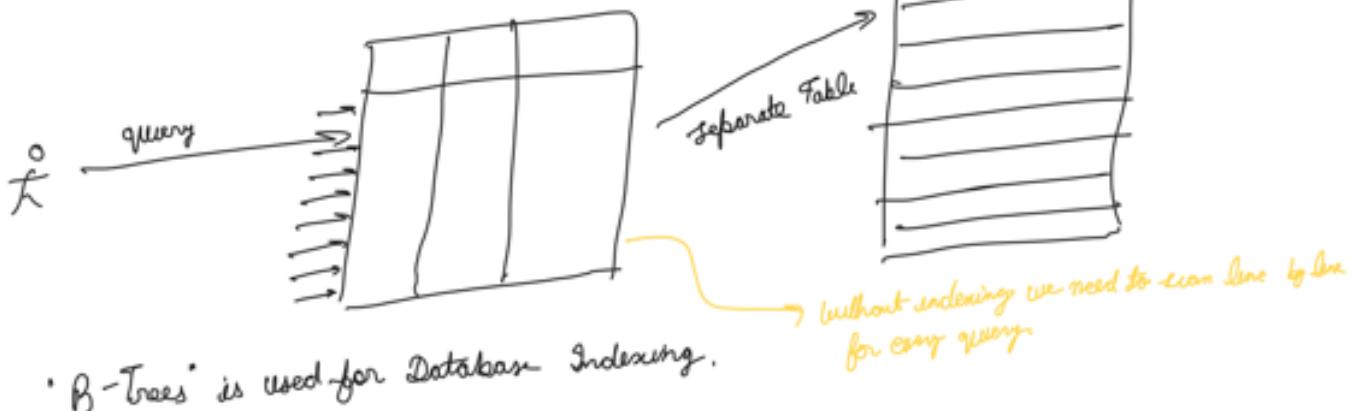


Caching →





⑦ Indexing in DataBase →



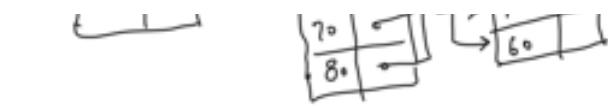
Types of Indexing Methods ⇒

P.K. Date	Transaction	Credit	Debit	Balance.

→ Primary → Main Value for uniqueness and fast access.

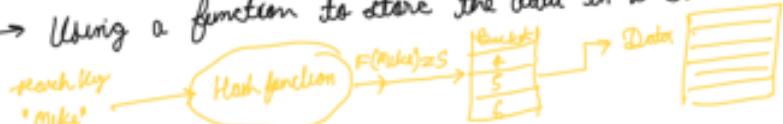
→ Secondary Indexing → Use another column with the primary key





→ Composite Indexing → Indexing Based on multiple columns to identify results.
e.g. → Amazon → multiple filter

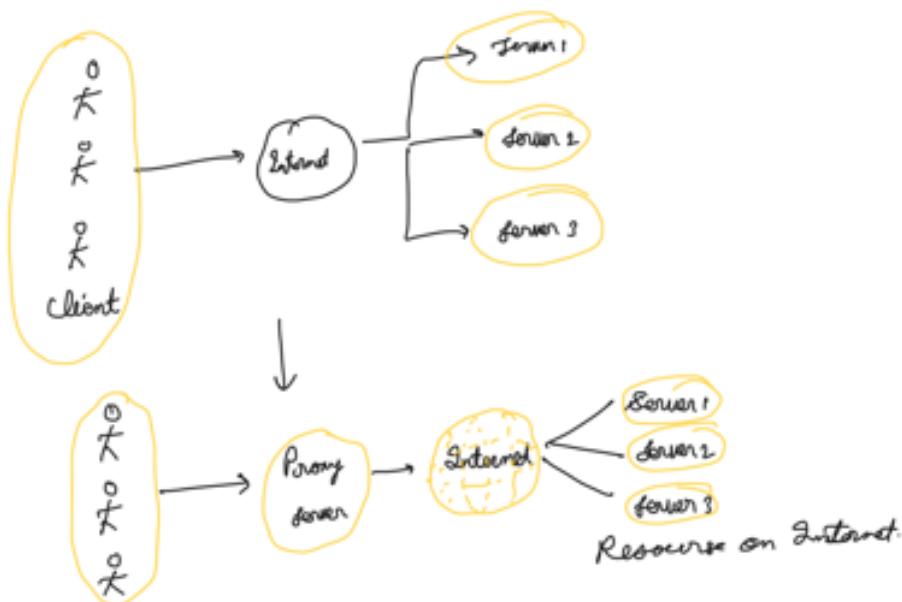
→ Hashing Indexing → Using a function to store the data in a calculated bucket.



Challenges →

- The no. of write operation will increase. (Update index table too)
- Adding index is not free, we need more space → more cost increased.
- cleanup unused indexes periodically. (Audit regularly).

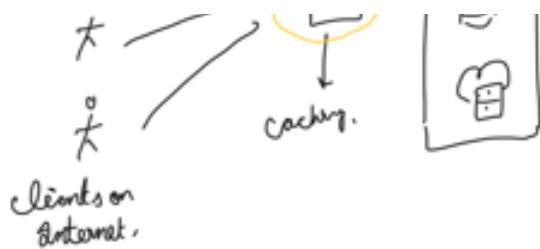
⑧ Forward and Reverse Proxy →



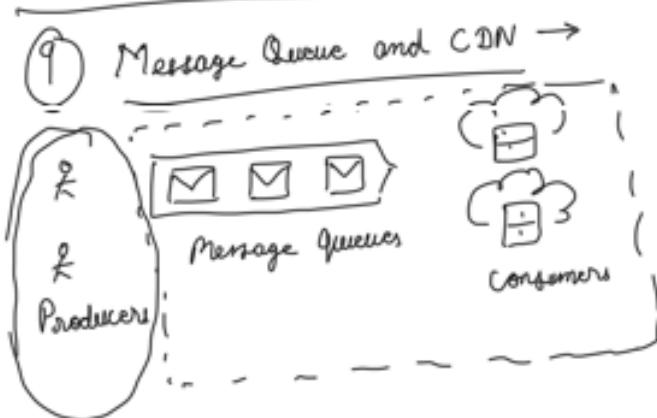
Advantages →

- Logging → Allowing you to log each request made by the user.
- Traffic Control → You can govern what traffic to pass and which to Block.
- Encryption → Make sure that all the traffic going through the proxy server is encrypted.
- Client Anonymity → Server don't know who the Client are?

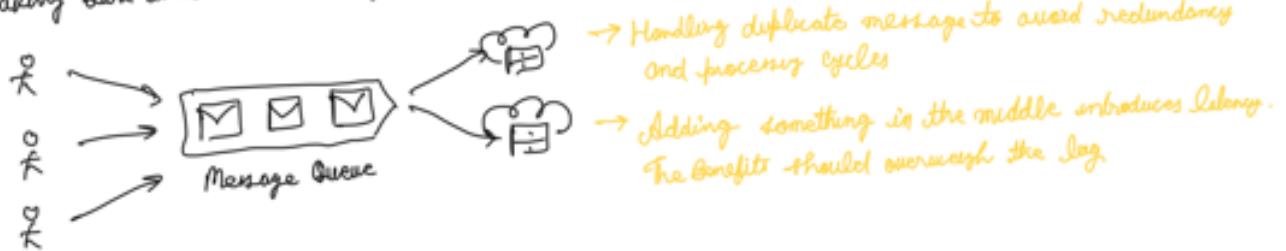




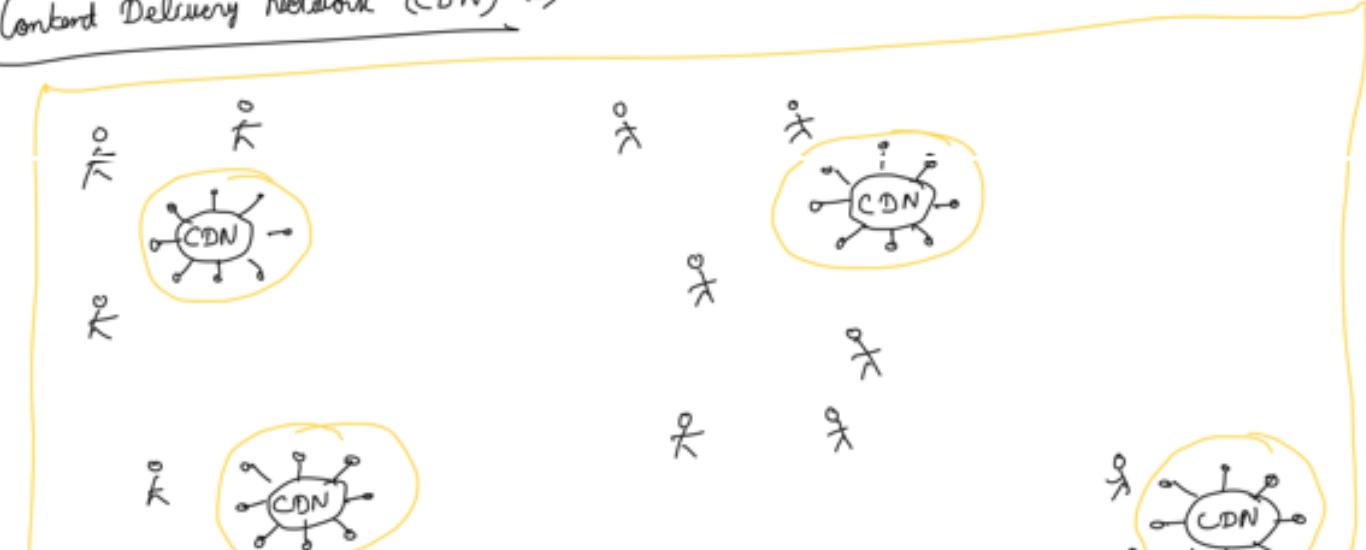
- (i) Server Anonymity → clients do not know what servers are they connecting to.
- (ii) DDoS Attack Prevention → without information about the server, we can't have a DDoS attack.
- (iii) Optimization Caching → can act as a load balancer (based on use case).



- Decoupling allows you to separate components and make them free sooner
- The system becomes more scalable as more clients can send in requests asynchronously
- you can process the message again, if the first attempt fails. This makes the system fault tolerant.
- Making sure that the ordering of the messages remain same in the order they were received



Content Delivery Network (CDN) →





Cons # Challenge →

- (i) your data should always be fresh.
- (ii) All the CDN infrastructure is costly.

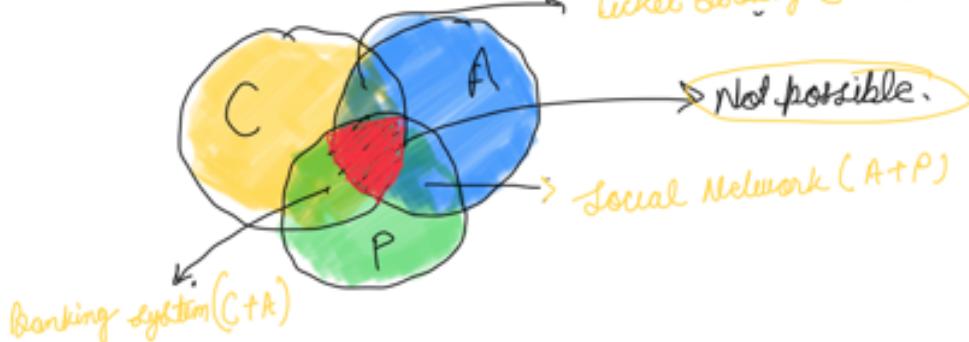
⑩ CAP Theorem

C → Consistency → System should deliver the same result no matter how and from where we query.

A → Availability → The system should always return some response or result.

P → Partition Tolerance → The system should be functional even if one partition is disconnected.

Ticket Booking (C+A)



The Paradox →



II Design Rate Limiter →

Q → Is this rate limiter client side / server side?

A → It will be server side implementation.

Q → How should I throttle the request?

A → There can be multiple conditions / rules to throttle a request.

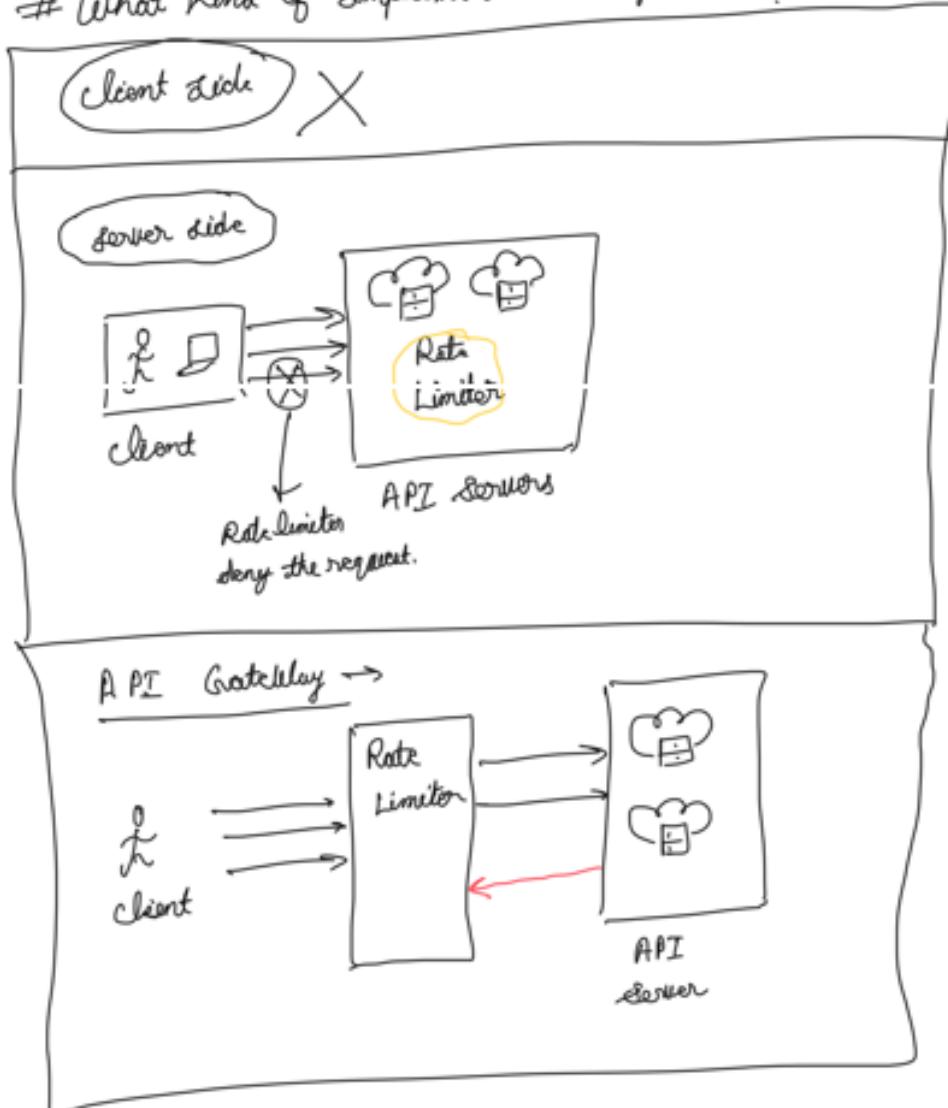
Q → Do you want to inform user that requests got denied?

A → Yes, we should inform the clients.

...

Q → Do you want a design a distributed system?
A → A Distributed Style.

What Kind of Implementation is possible?

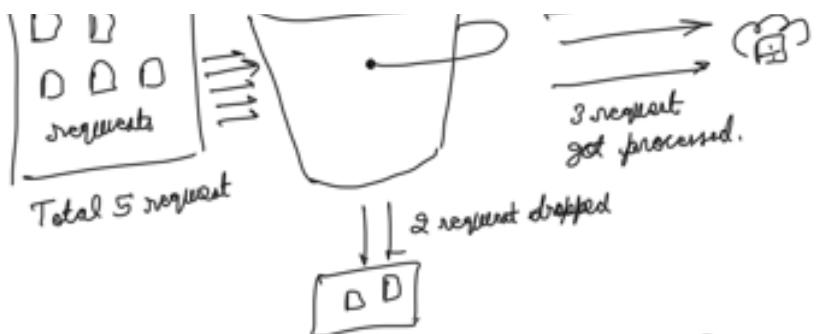


Different Algorithm \Rightarrow

- (i) Token Bucket Algorithm
- (ii) Fixed Window Algorithm
- (iii) Sliding Window Log

(i) Token Bucket Algorithm \rightarrow





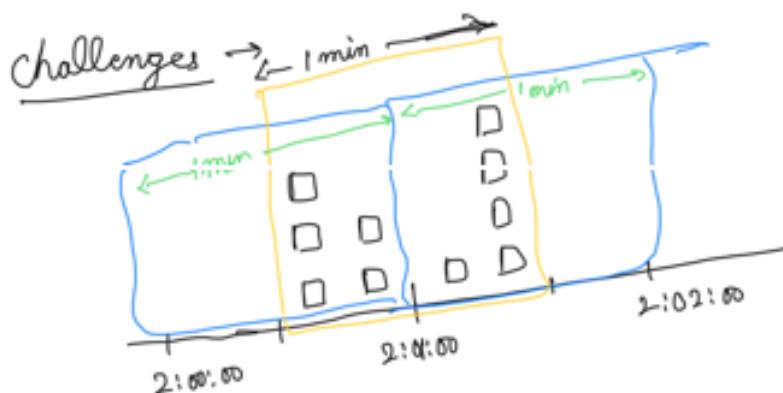
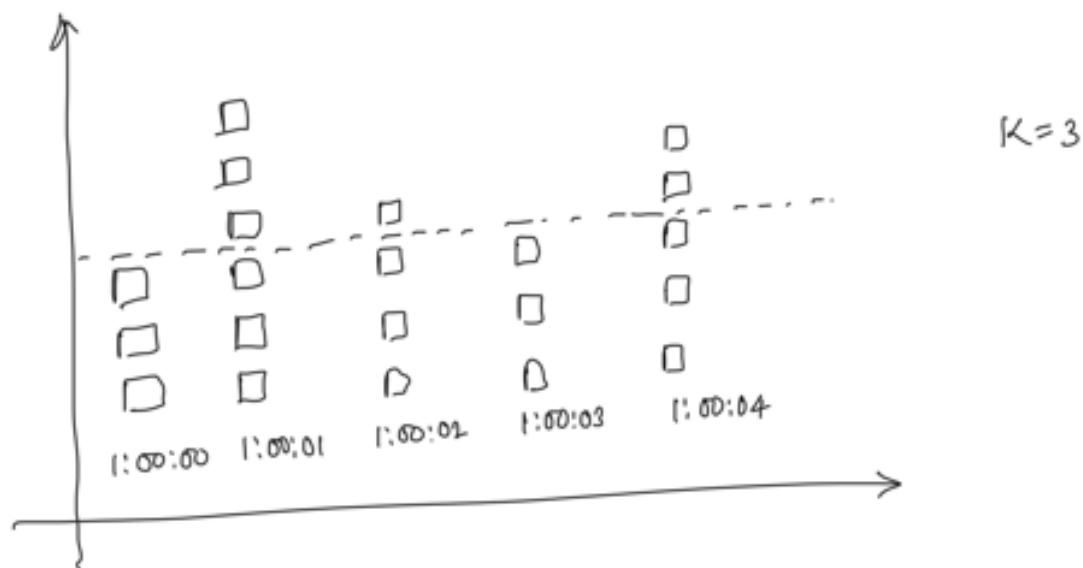
Advantages →

- (i) Very Easy to implement
- (ii) Memory Efficient.
- (iii) Allows Burst of traffic

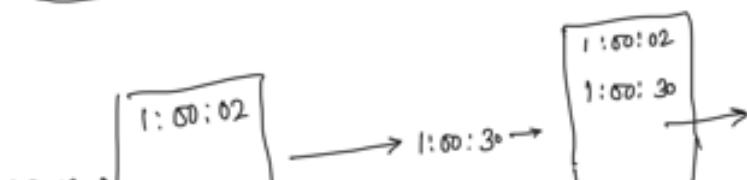
DisAdvantages →

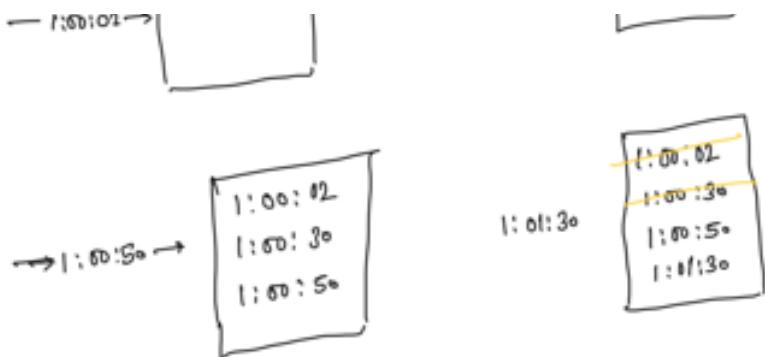
- (i) Tuning the refill rate and bucket size is tricky

(ii) Fixed Window Counter Algorithm ⇒



(iii) Sliding Window Log Algorithm →





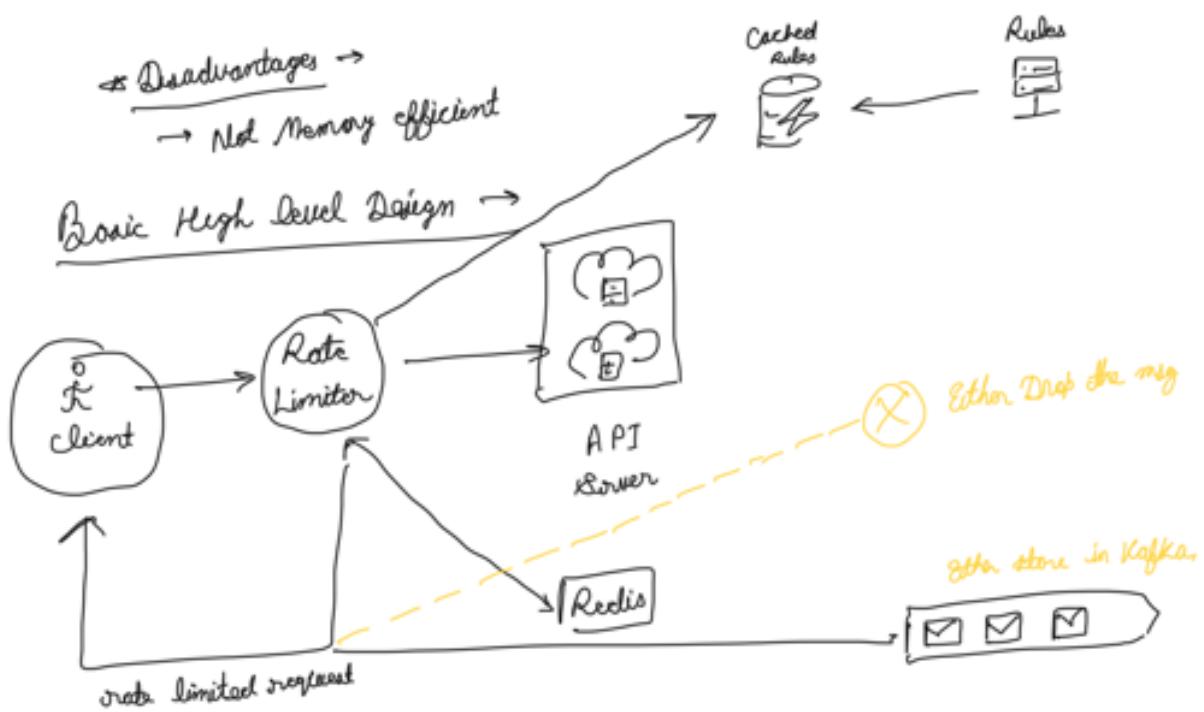
* Advantages →

→ very accurate and requests are handled correctly in a rolling window.

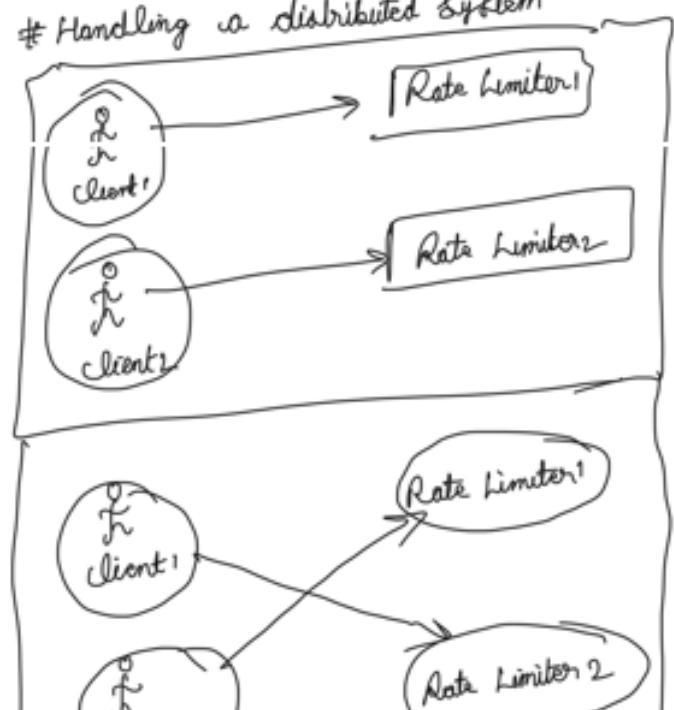
* Disadvantages →

→ Not Memory efficient

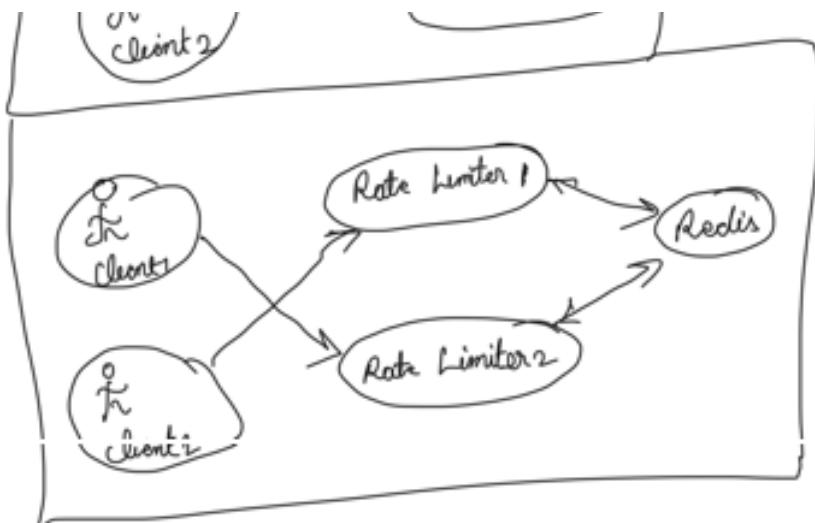
Basic High Level Design →



Handling a distributed System →



→ Rate Limiter 1 does not know about Client 2 that how many req. Client 2 has made and vice versa.



② Design a Unique ID Generator →

- # why do we need Unique IDs?
- image/social media post
- Video for youtube
- Songs on spotify
- order-ID for e-commerce
- comments on social network
- Cash transaction for Banking
- IoT device identification
- message and event driven architecture

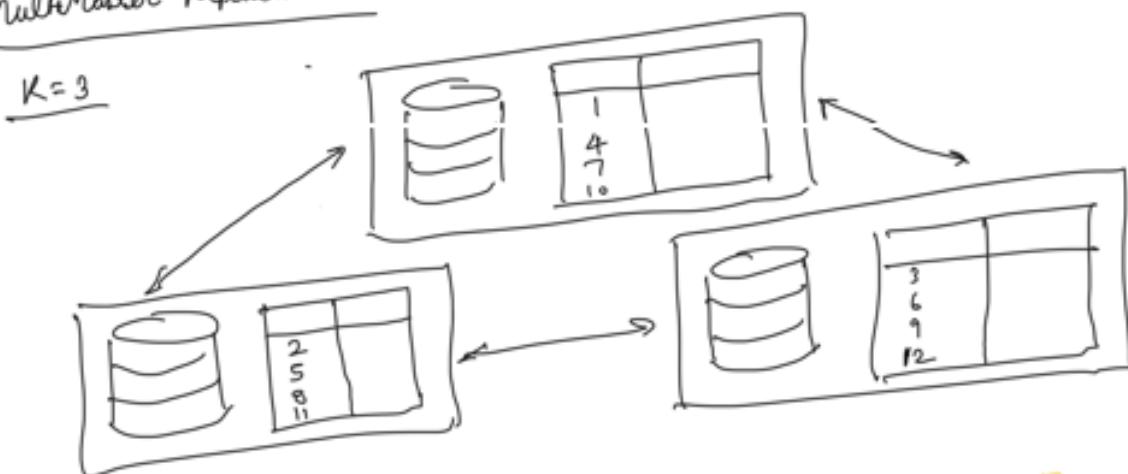
A very simple solⁿ →





- Not good for distributed system
- A single ID can get assigned to multiple resources.

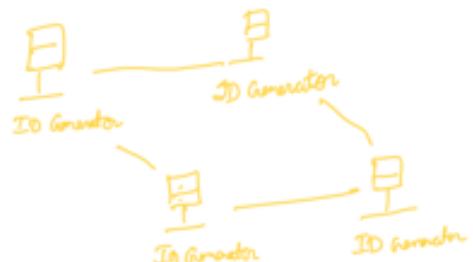
C) MultiMaster Replication →



- Easy to Implement
- Not Easy to scale.

Using UUIDs → always 128 bits

- simple way to generate UUIDs
- No Co-ordination need b/w servers
- Easy to scale and add new nodes



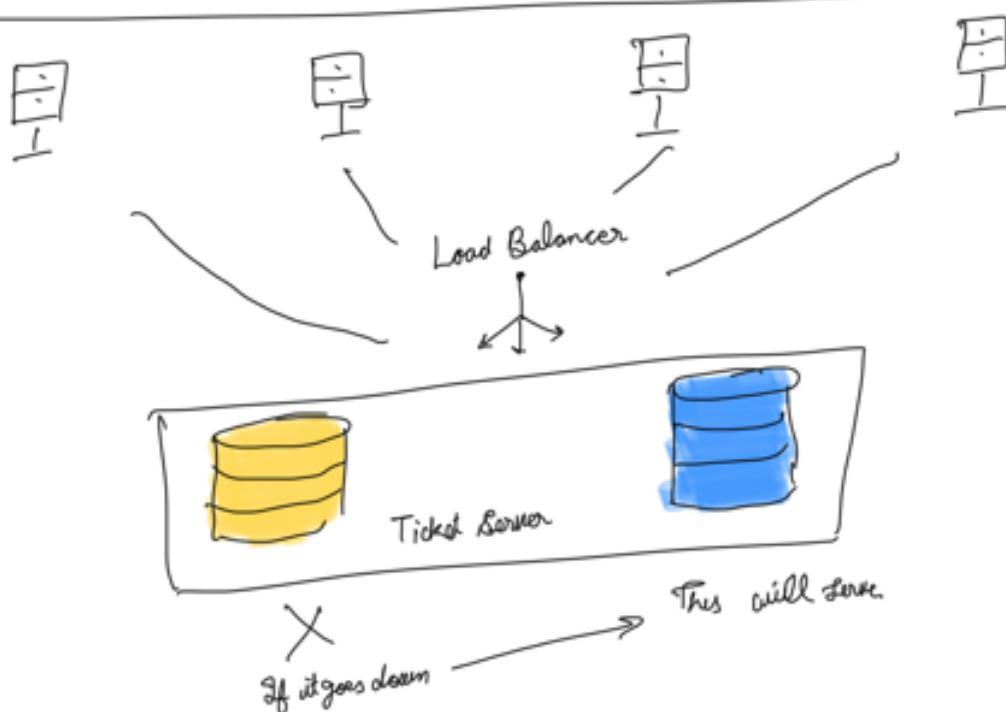
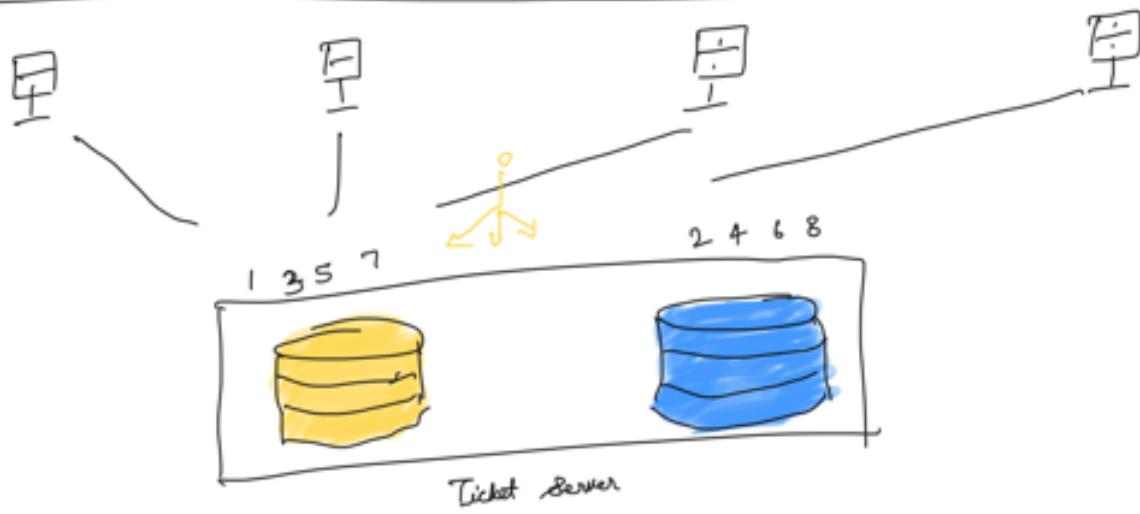
Advantages

- IDs generated are 128-bit (too long)
- UUIDs are alphanumeric
- Not Sortable

Disadvantages

Ticket server Approach →





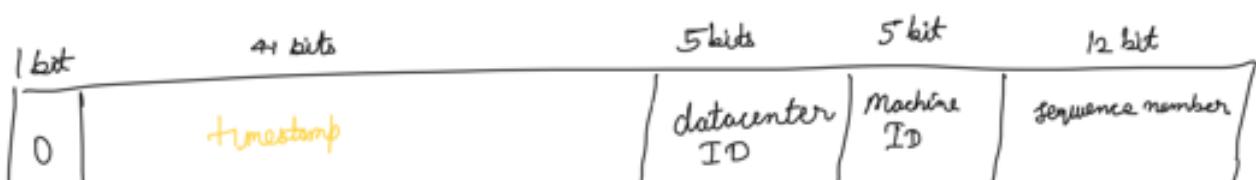
Advantages →

- Easy to implement.
- We can get unique IDs in any format.

Disadvantages →

- Not good for a huge system.
- Single point of failure.

Twitter Snowflake Approach →



Advantages

- (i) Unique IDs are generated without any collision.
- (ii) Fast, sortable IDs (also has information about timestamp).

Disadvantages →

- (i) Time synchronization across servers is important.
- (ii) The implementation can be tricky.
- (iii) Limited by time range (69 years).

⑬ Design a URL shortener → (long URL → short long URL)

Requirements for a URL shortener →

Q → What is the volume of URLs that we want to handle?

A → 100 millions / day.

Q → How short the URLs should be?

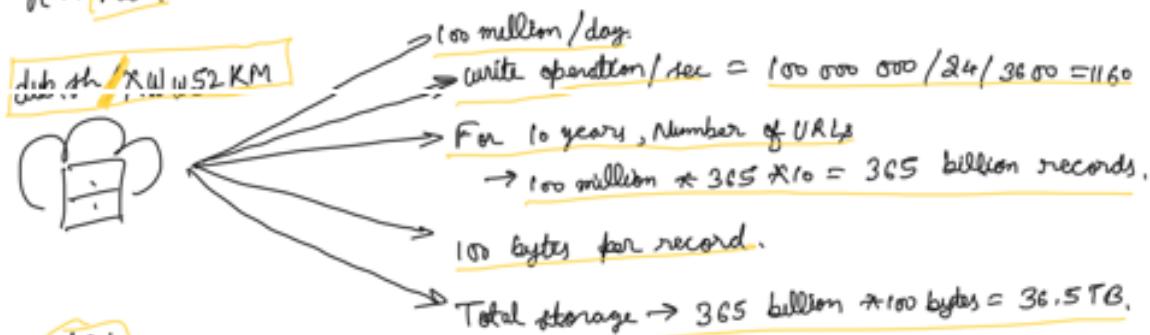
A → As small as possible.

Q → What kind of character do you need?

A → Alphanumeric without symbols.

Q → Can URLs be updated / deleted?

A → No.

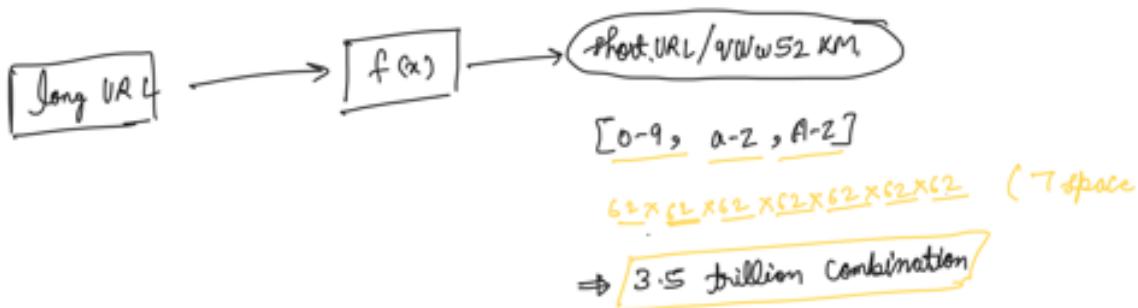


client
HTTP

301: Permanent redirect
302: Temporarily moved.



URL shortening (High level view)

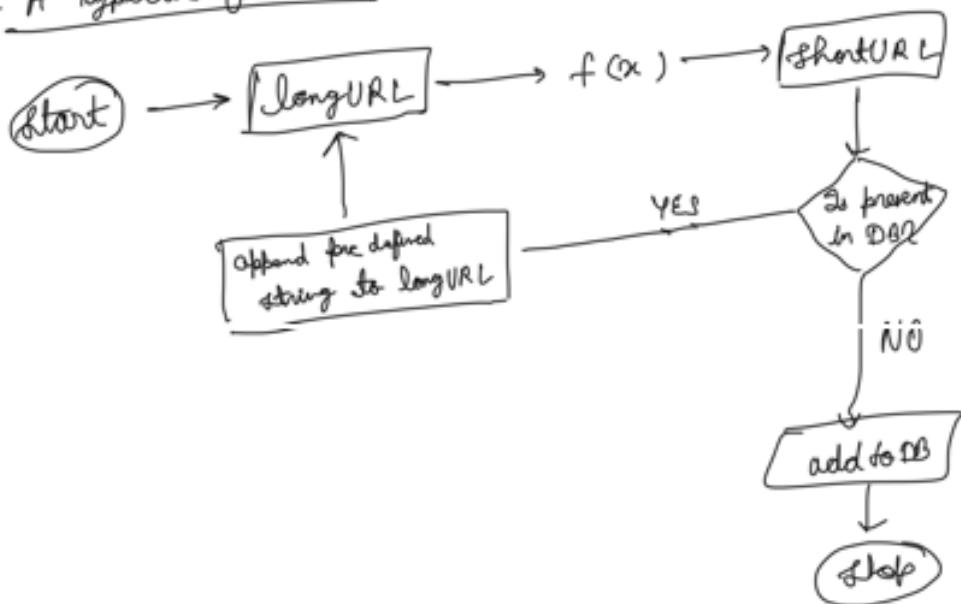


Hash + collision Resolution →

https://www.youtube.com/playlist?list=PLFdAYMIVJQH0WJgRnfv-RH-ny95B2hSON

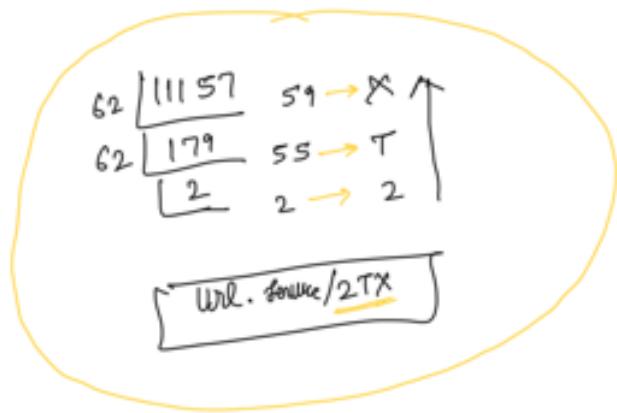
MD5	5cb54054
SHA-1	5a62509a84df9ce03fe1230b9fd8b84c
CRC32	0ecaec7916c06853901d9ccbefbfcaf4de51cd856

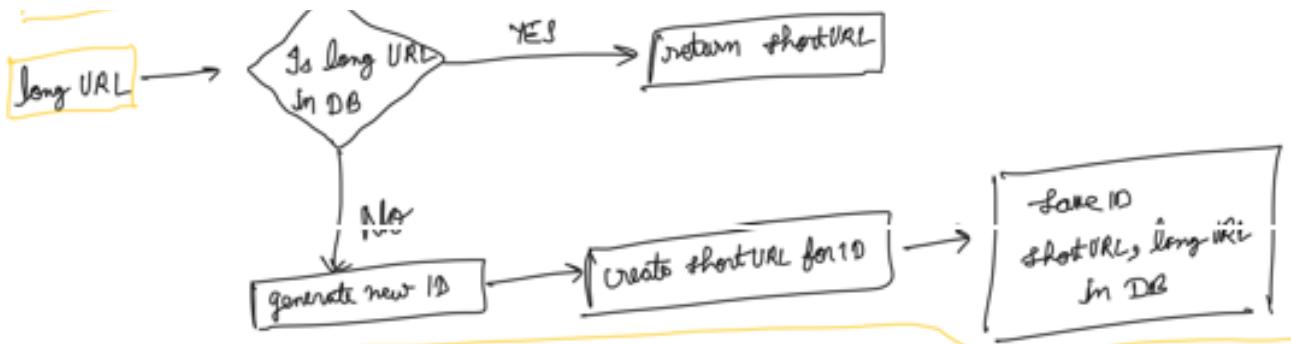
A Typical flow →



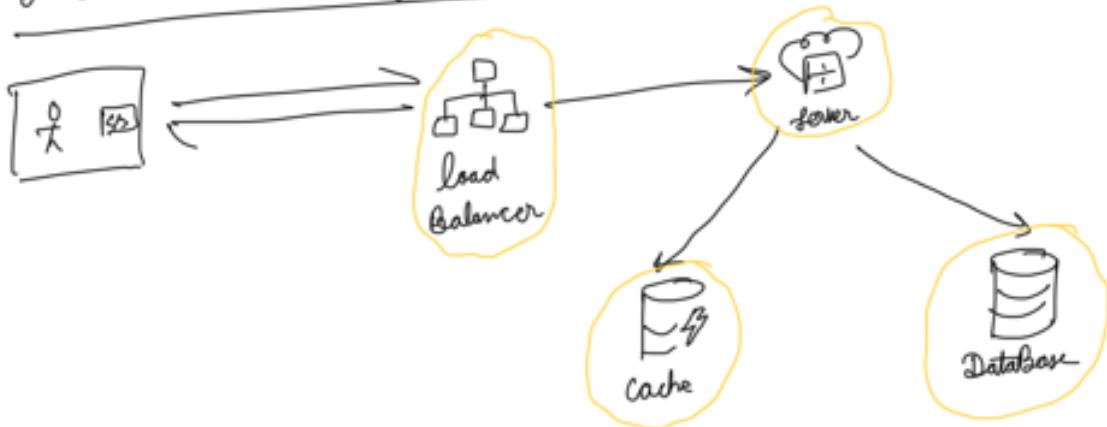
ID Based Approach →

$[0-9 \ a-z \ A-Z]$		
0-0	a-10	A-36
b-11	1	B-37
:	:	:
z-9	Z-35	Z-61





A complete overview of the Design →



⑭ Notification system →

Push Email SMS

Q → what kind of notification do you want to send?

A → All 3 type (Push Email SMS)

Q → Are these notification realtime?

A → we want them in near realtime. A small delay is fine in busy hours.

Q → What kind of devices do we want to send the notification to?

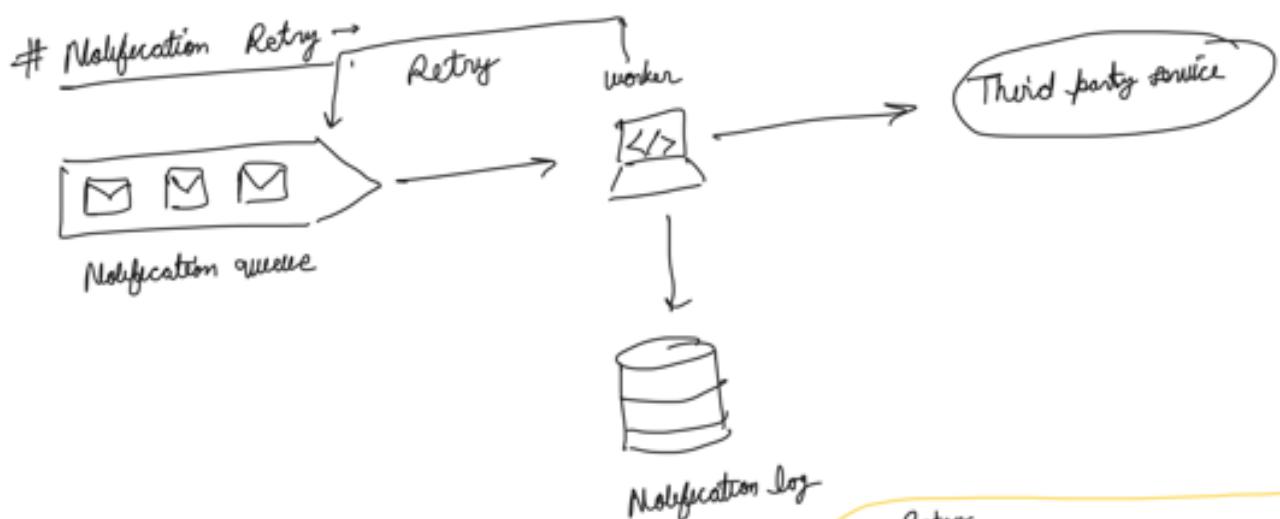
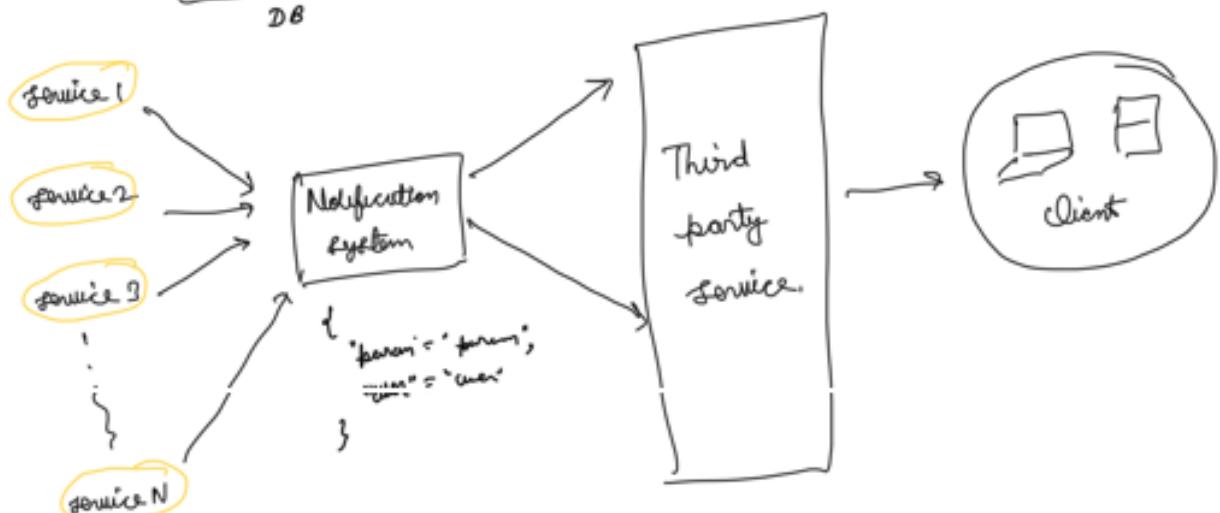
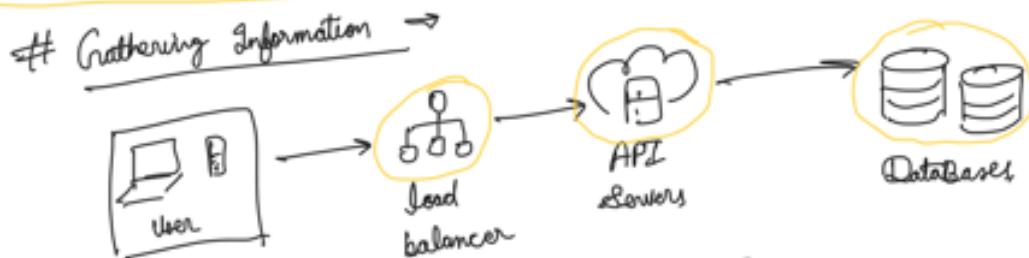
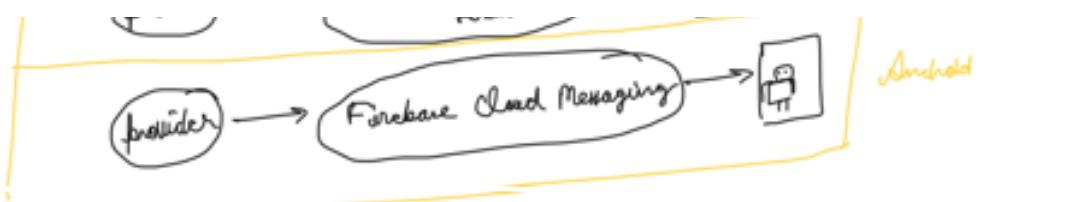
A → All device (Laptop / Mobile / iOS / Android).

Q → How are the notification triggered?

A → Can be automatic / scheduled / manual.

High Level design →

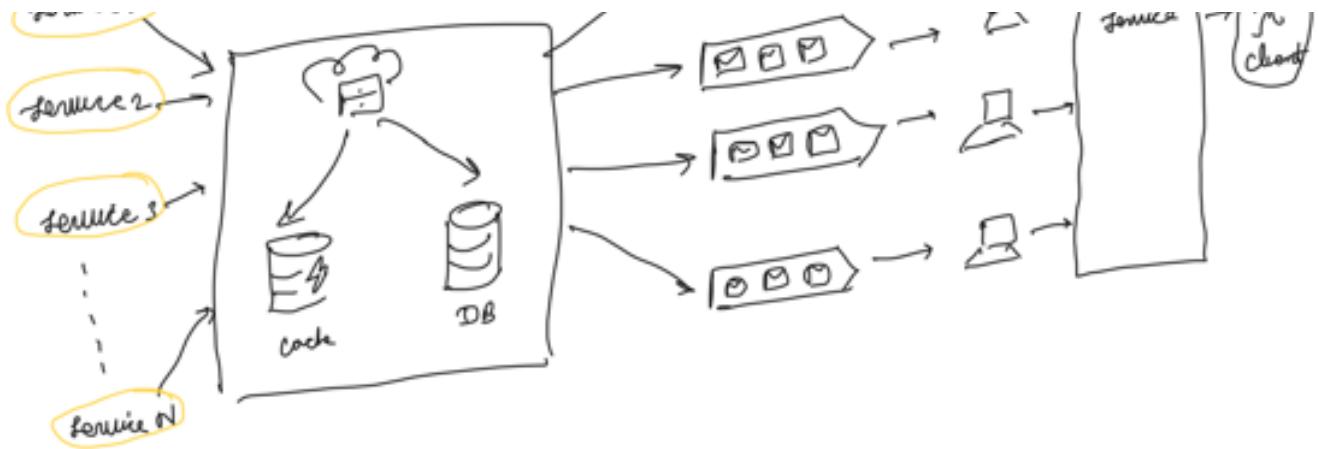




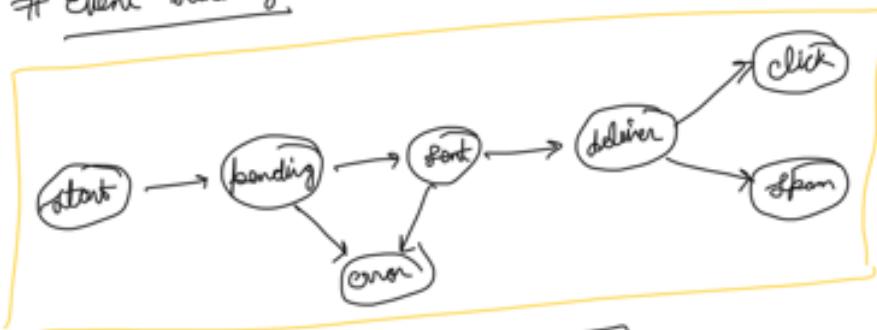
Improved Design →

service 1

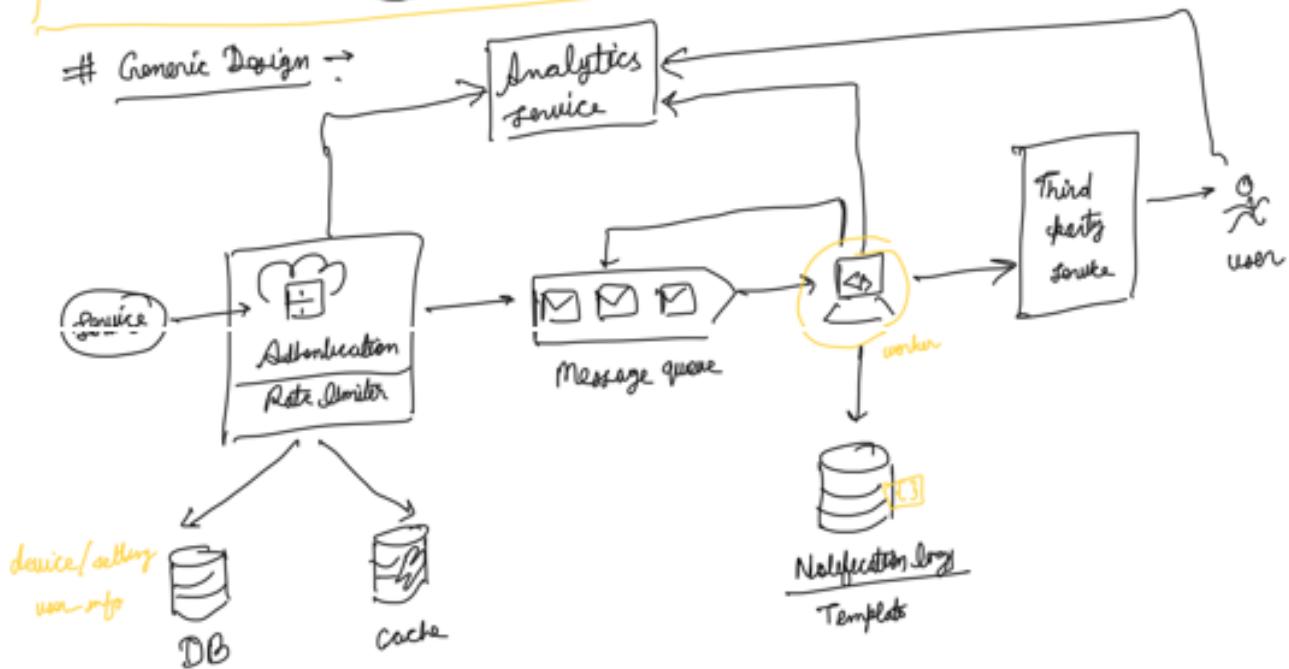




Event Tracking →

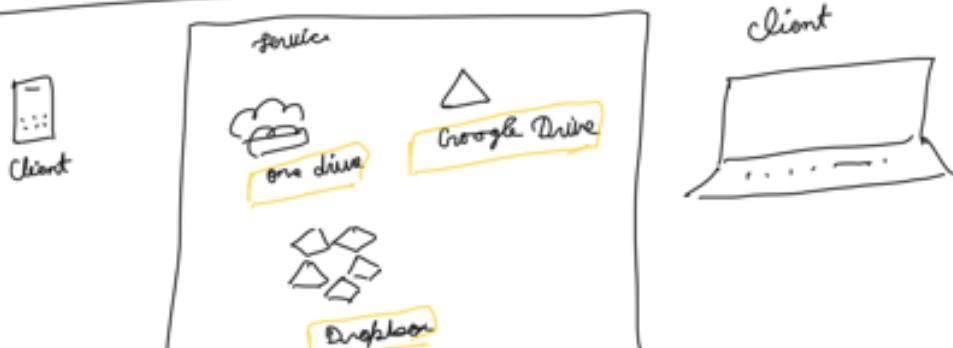


Generic Design →



⑤ Design Dropbox →

A service like Dropbox →



Advantages →

- (i) Availability → Data is available everywhere with internet connection.
- (ii) Reliability and Durability → Data is secured and remains on the cloud.
- (iii) Scalability → your needs can expand and reduce as you desire.

Requirements and Considerations →

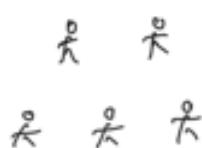
- (i) User can upload and download from all of the configured clients.
- (ii) User can share files with others.
- (iii) offline editing of files.
- (iv) Synchronise files on different clients.

→ Versioning → Restore previous versions.

→ Premium subscription with more features

→ Very heavy on bandwidth. High upload/ downloads

Idea of the Scale of System →



→ Each user has average of 3 devices



→ Each user has average of 200 files.



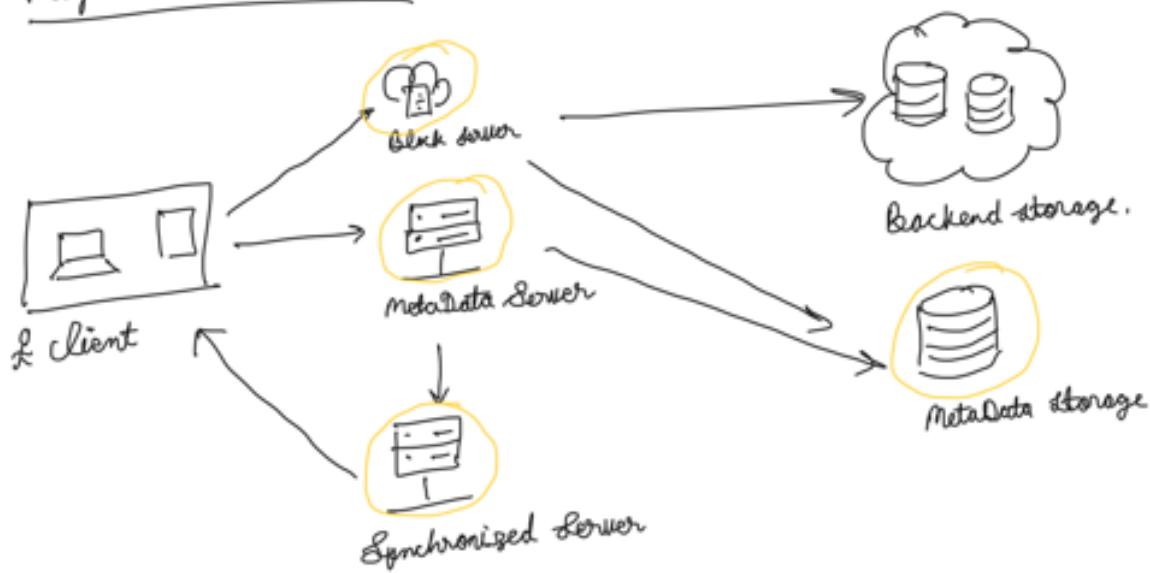
Total files $\Rightarrow 200 \times 500 \text{ million} = 100 \text{ B}$

500 million total users
(100 million active) (2:1)

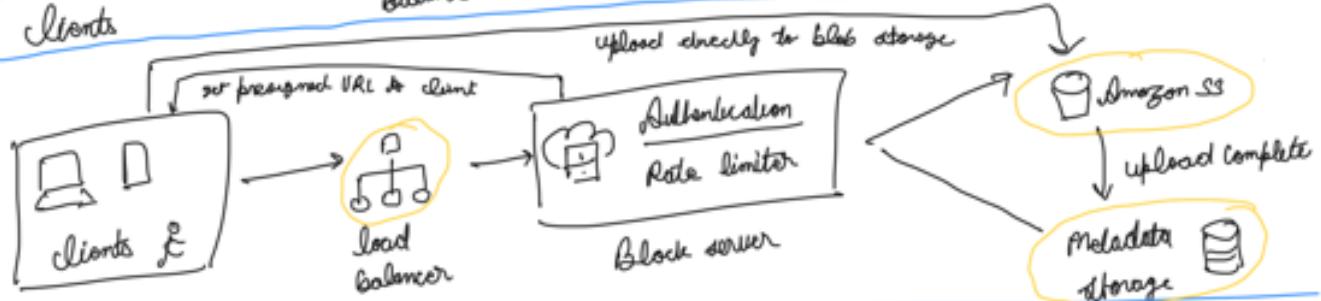
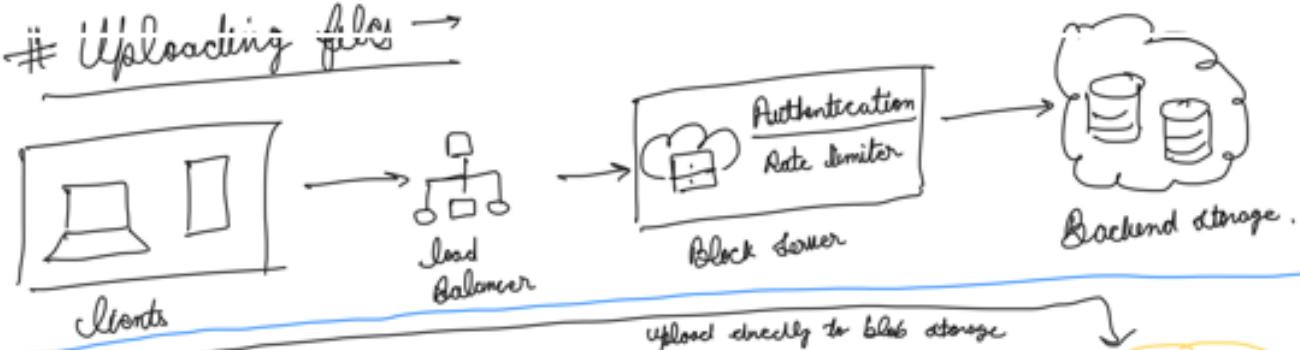
→ Average file size = 100Kb

Total storage = $100 \text{ B} \times 100 \text{ Kb} = 10 \text{ PB}$.

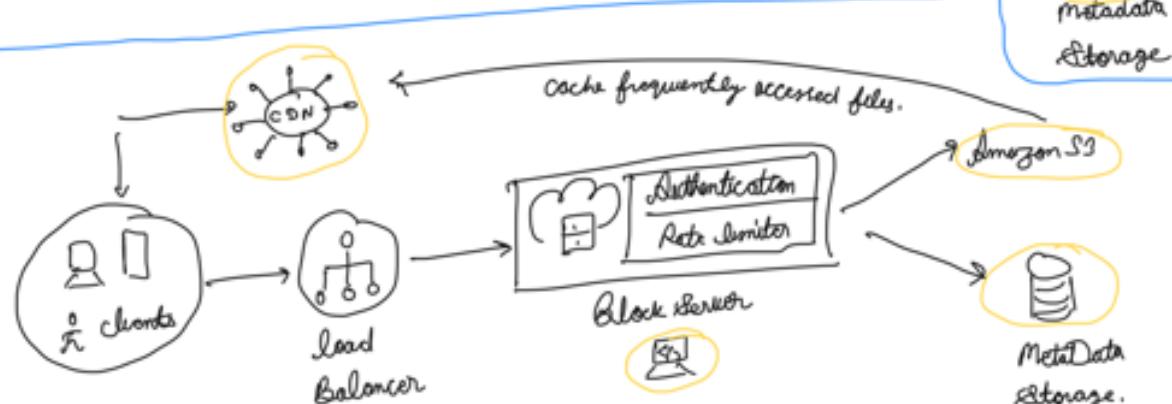
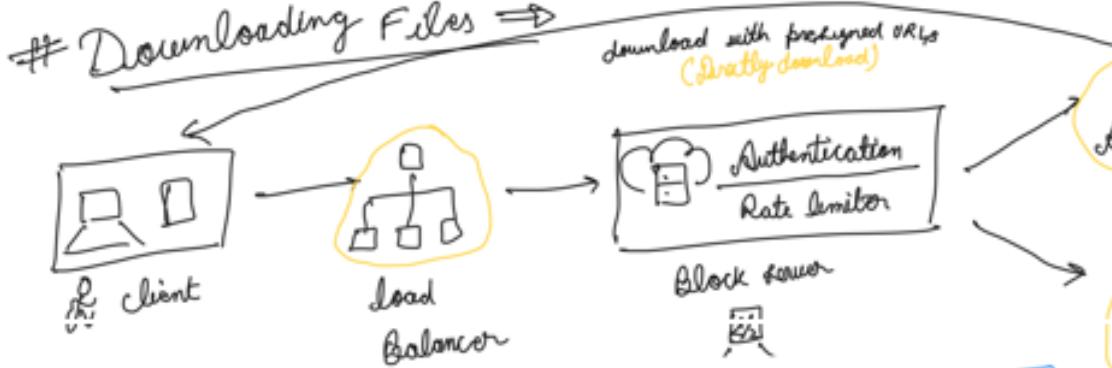
High Level Design ⇒



Uploading files →



Downloading Files →



Sharing files →



```

t id: 123
  name: file.txt
  size: 1000
  mimeType: text/plain
  uploadedBy: user1
  shareList: ["user2", "user3"]
  
```

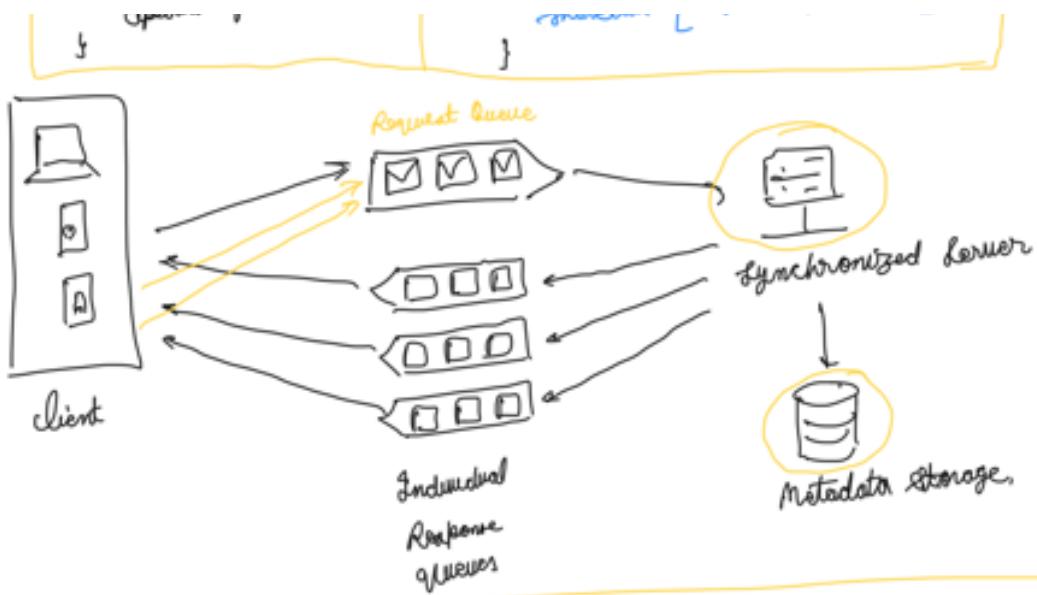
```

t id: 123
  name: file.txt
  size: 1000
  mimeType: text/plain
  uploadedBy: user1
  shareList: ["user2", "user3"]
  
```

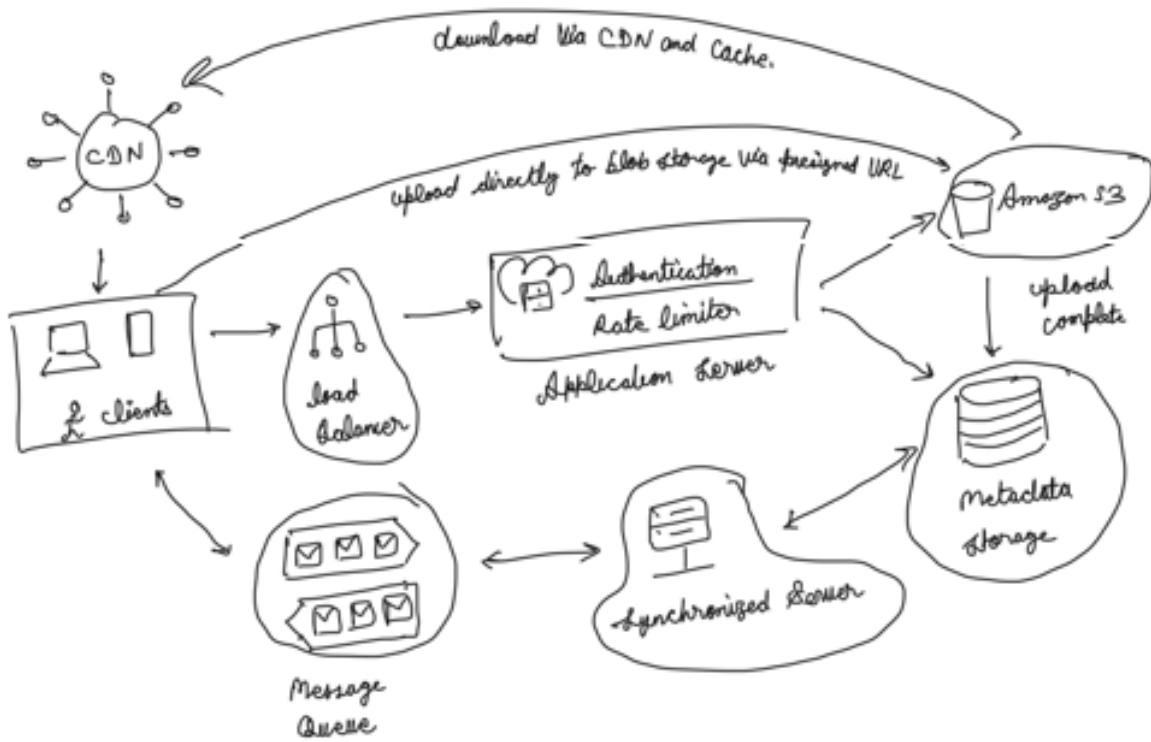
user1:[fileID1, fileID2]

```

t id: 123
  name: file.txt
  size: 1000
  mimeType: text/plain
  uploadedBy: user1
  shareList: ["user2", "user3"]
  
```



Full Picture →



⑯ Design a NewsFeed →

- Should be available for mobile and as a web application too.
- Most important feature: should be able to publish updates and push to friends.
- Update should be in near realtime.
- Have the update in a reverse chronological order (Newest first).

Q → what traffic volume is expected?

Q → what is the maximum no. of friends?

Q → what do post contain ('image | text | video')

Idea of the System →



Fanout Service →

Fanout on write (PUSH Model)

- News Feed Cache is pre-populated during write.
- generated in real time, can be pushed immediately
- fetching news feed is fast, as it is pre-computed
- if a user has many friends following and generating the newsfeed is very slow.
- for inactive users, waste a lot of resources.

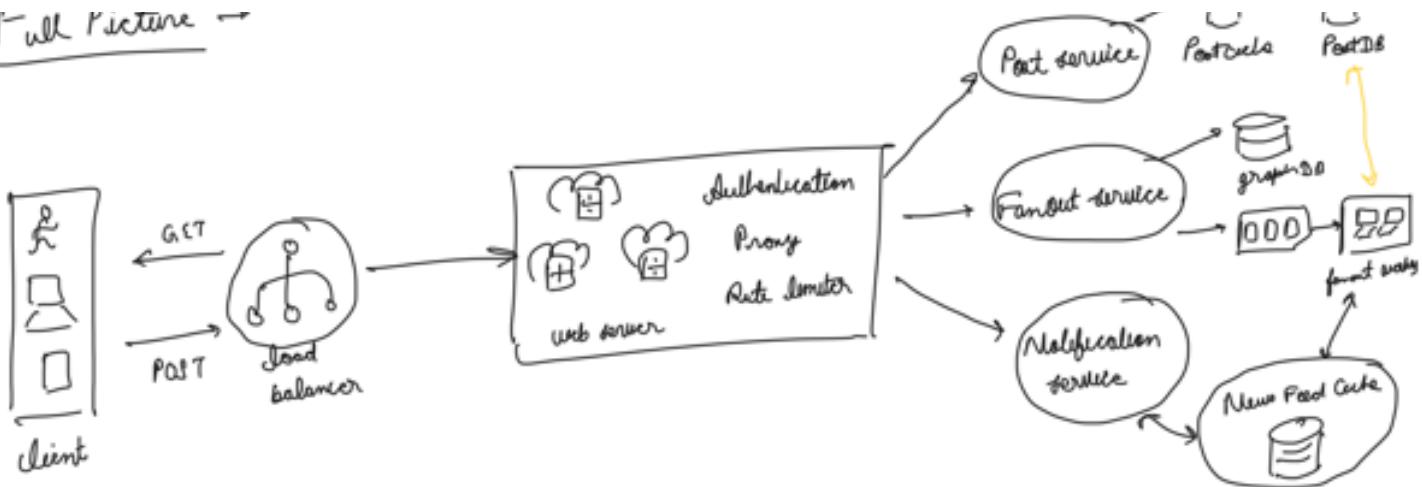
Fanout on Read (Pull Model)

- News Feed Cache is populated during read time
- for inactive users, fanout on read works better.
- data is not pushed to everyone, saves bandwidth.
- the process is slow
- User can have a bad experience.

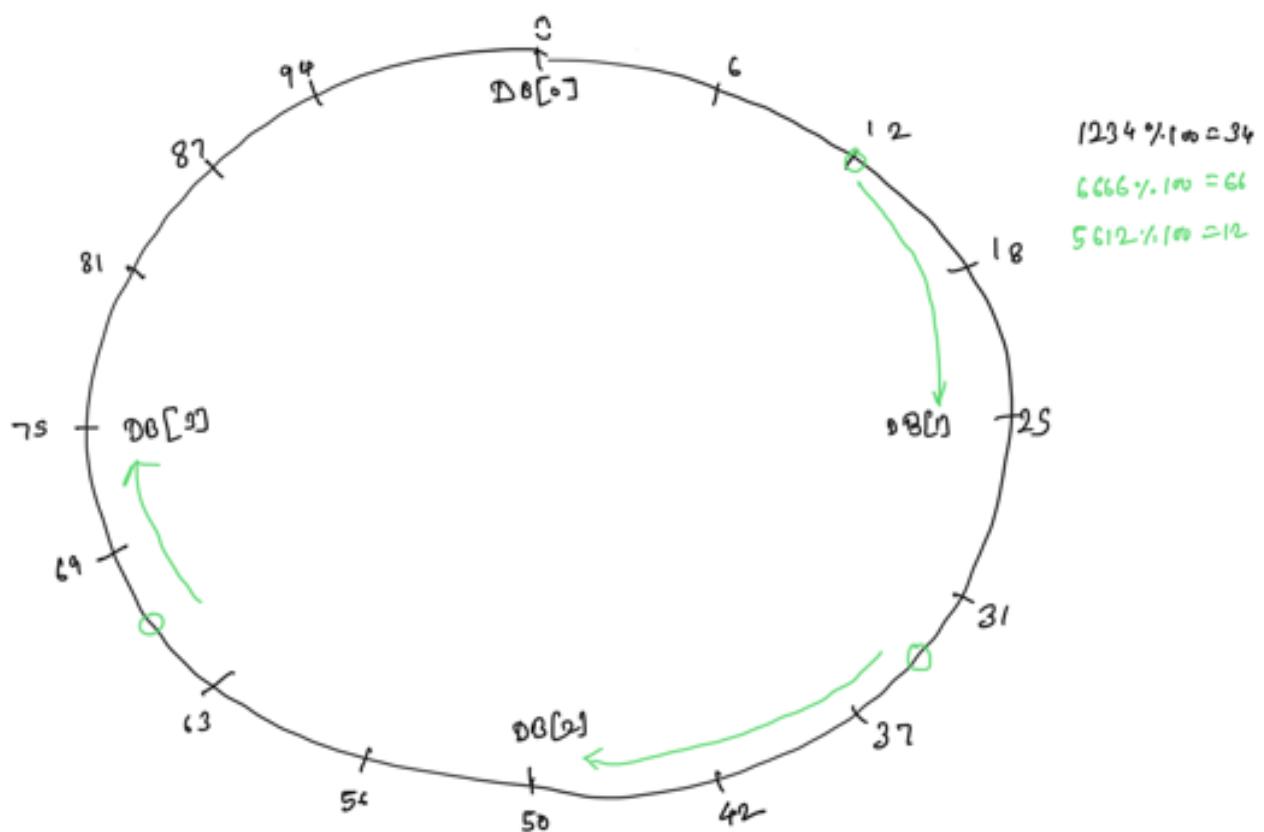
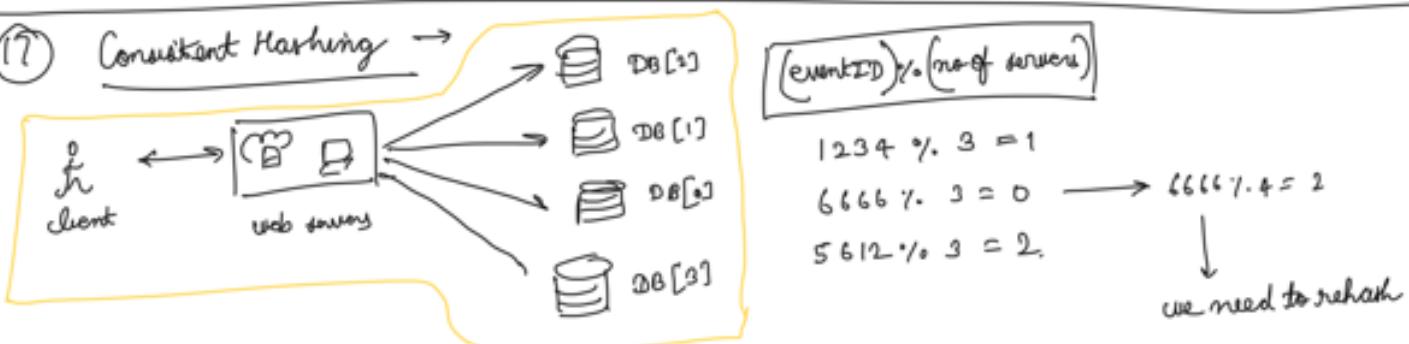
Design a Fanout Service →



Full Picture →

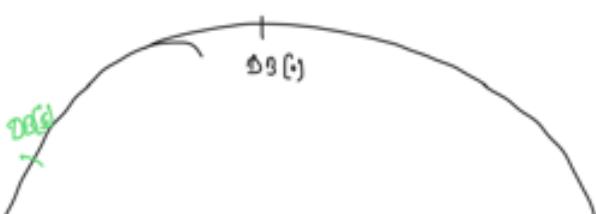


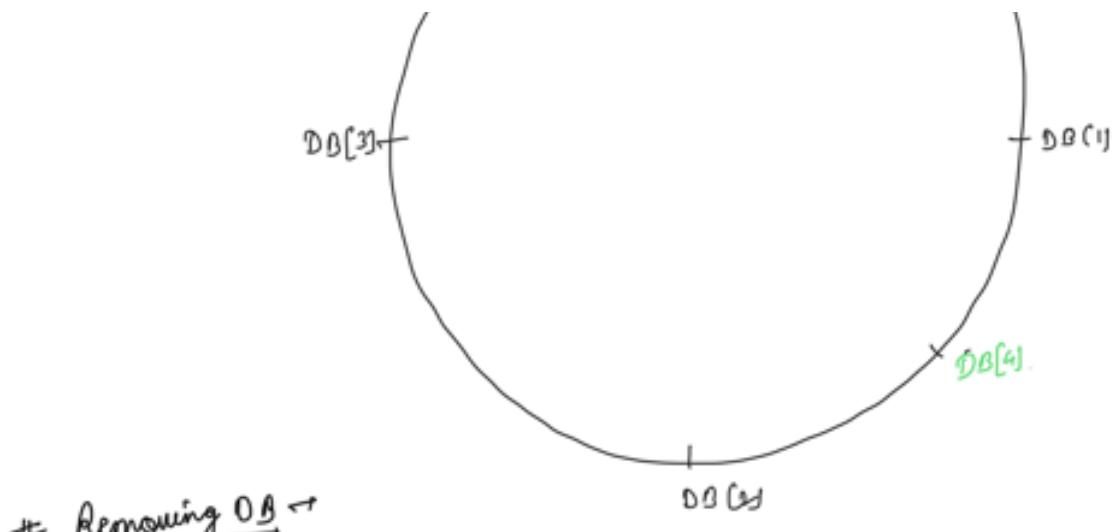
⑦ Consistent Hashing →



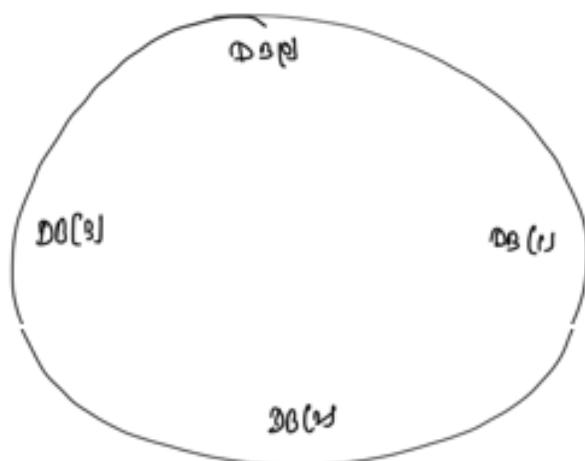
Let's say we need to add new DB[4]

Adding DB →

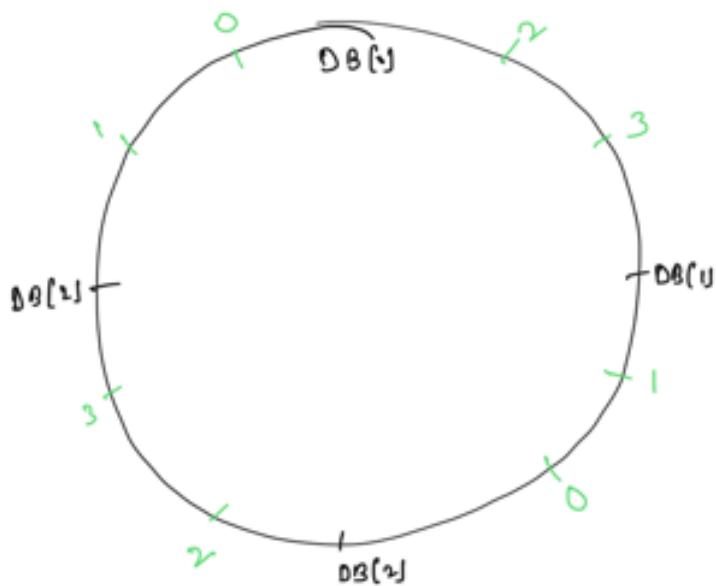




Removing DB →

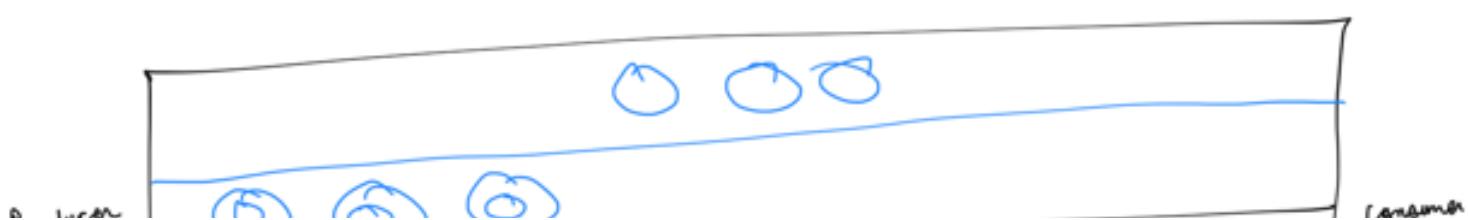
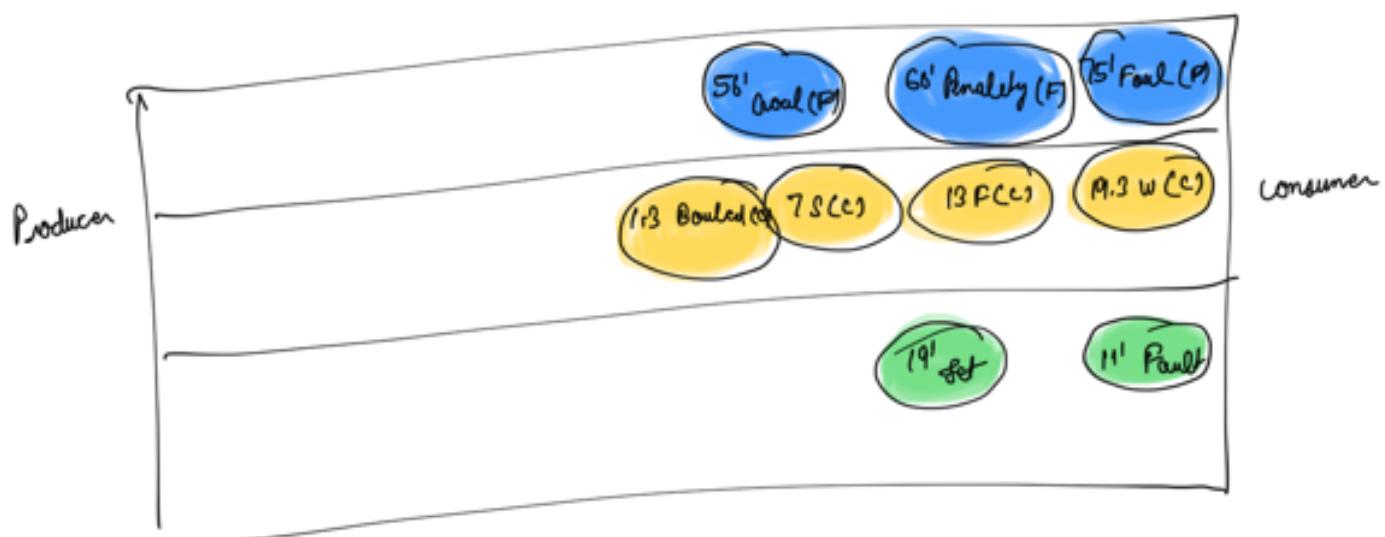
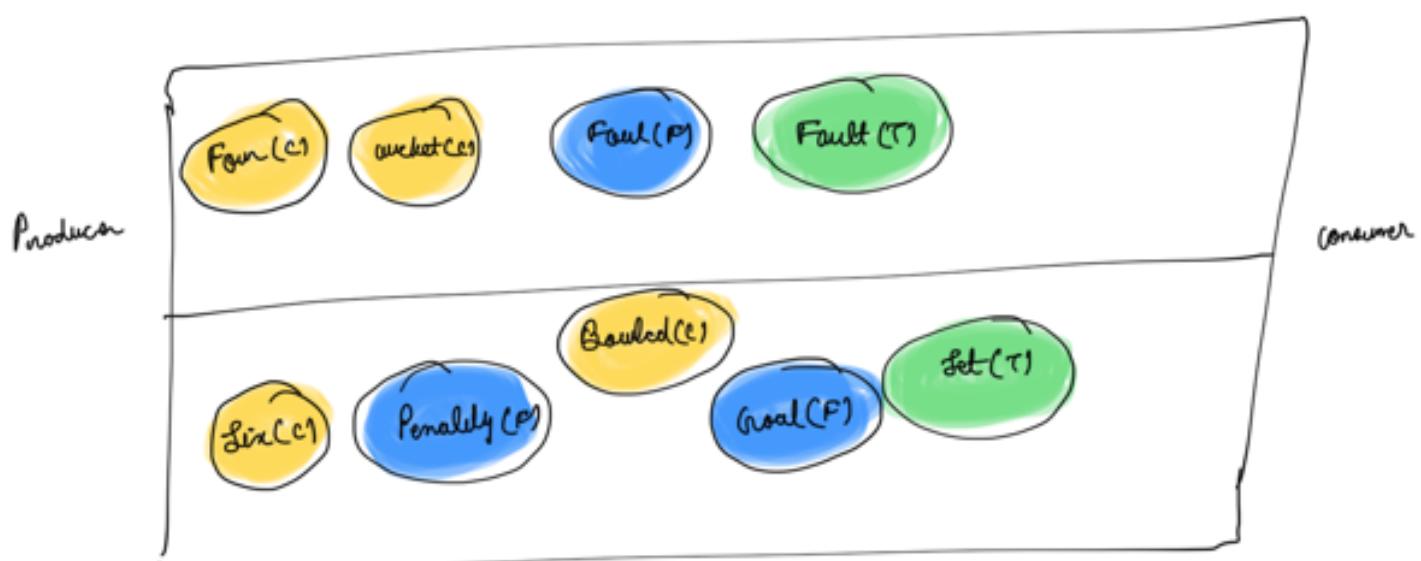
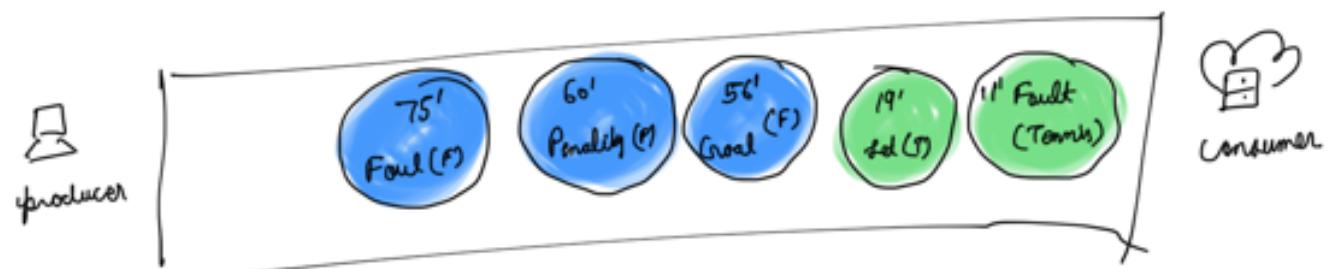
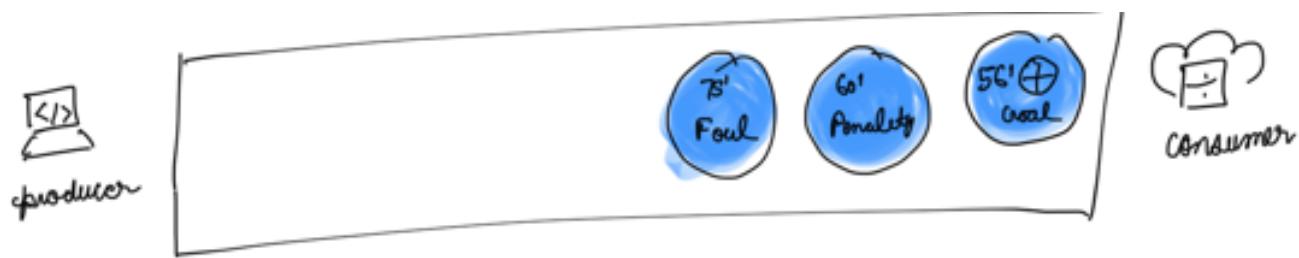


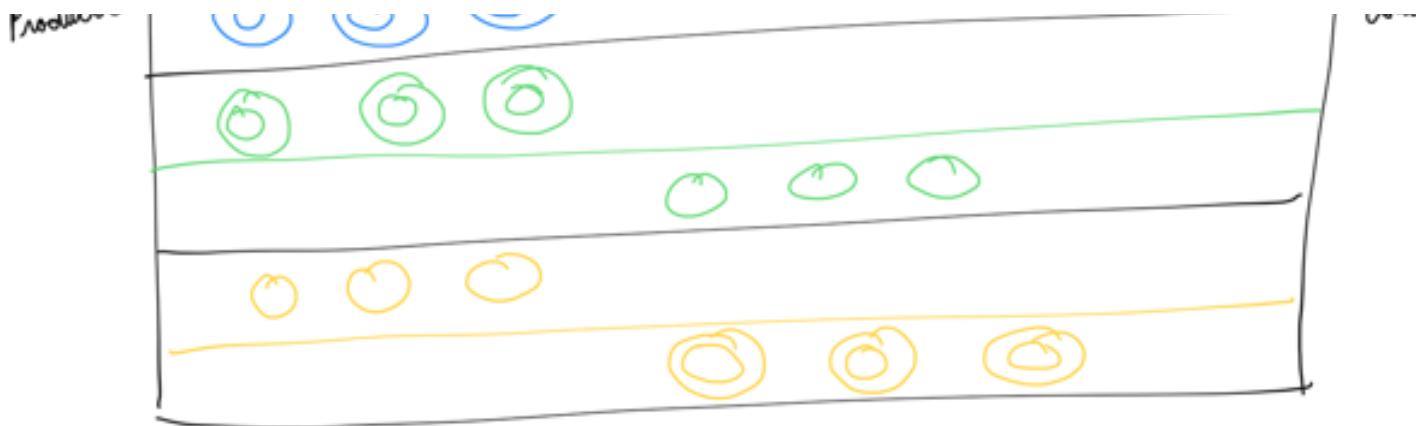
Virtual DB Nodes →



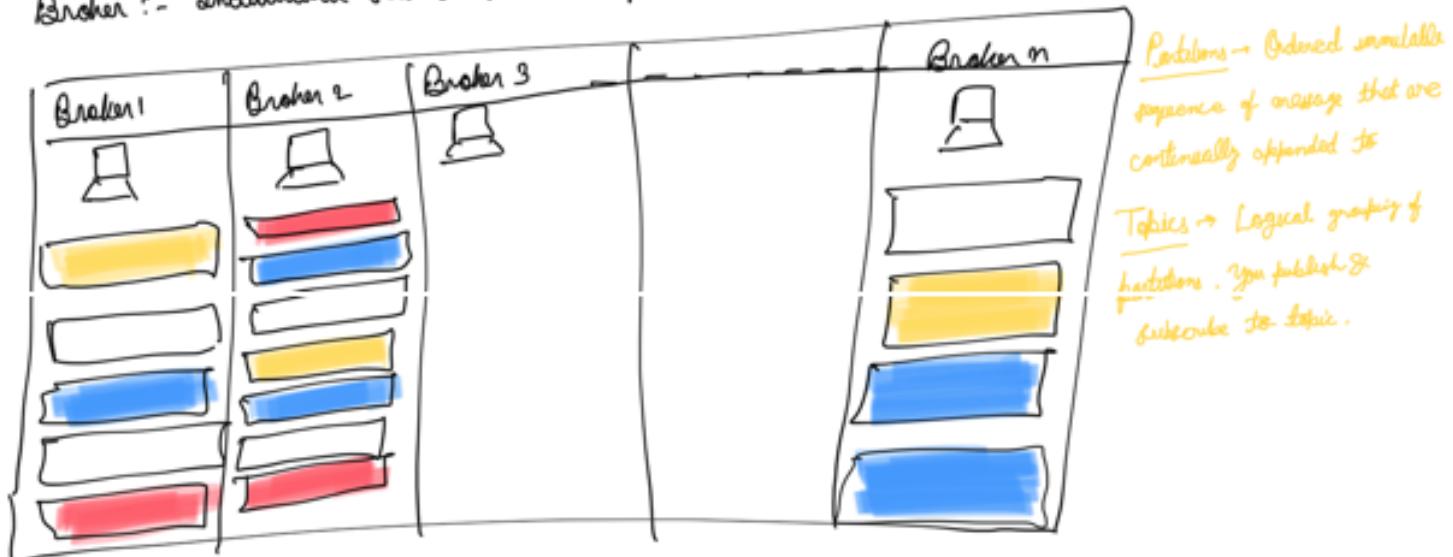
⑮ KAFKA →

→ A motivational Example



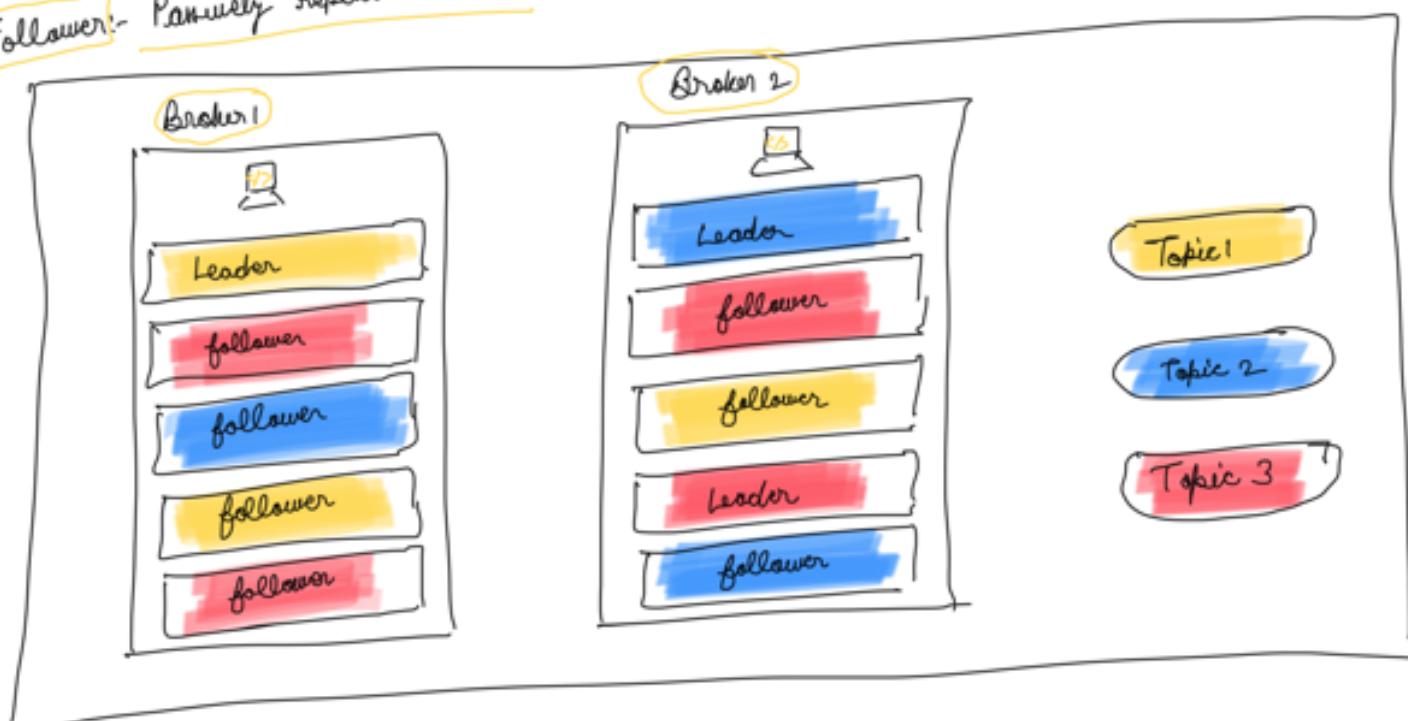


Broker :- Individual servers that make up the cluster and store information.



Leader: Handles all incoming data.

Follower: Passively replicate leader

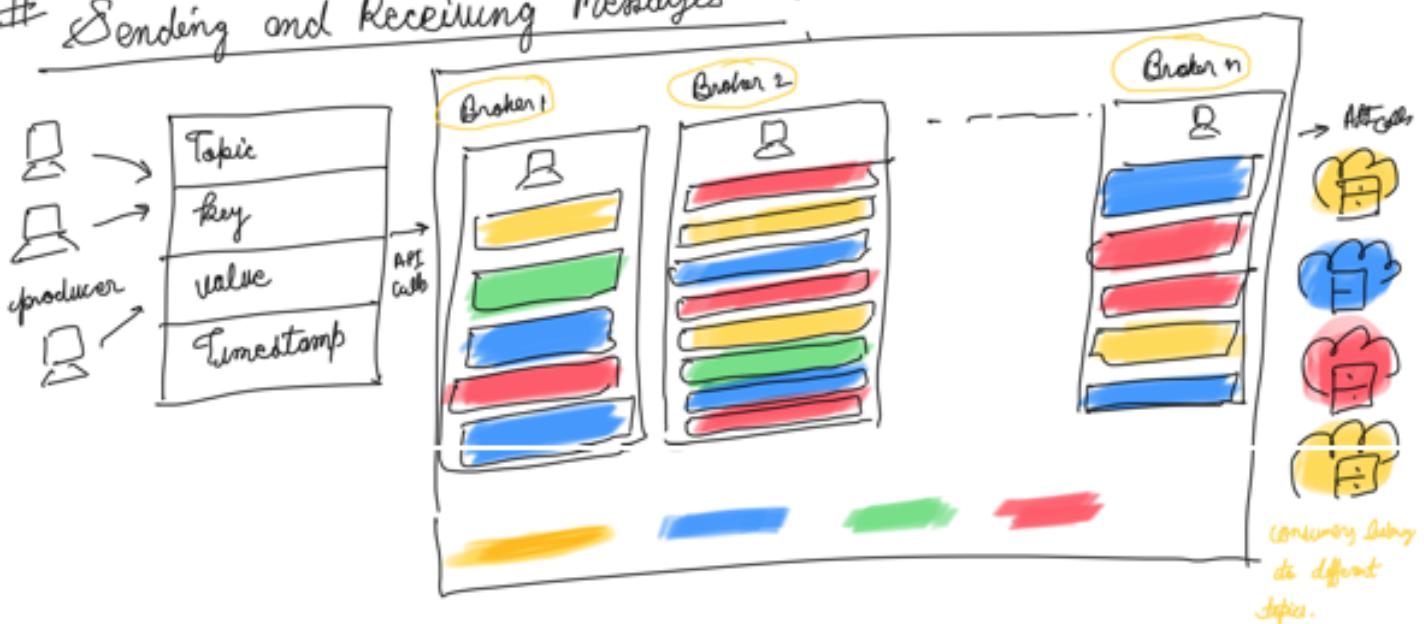


Topic 1 →
1 Leader
2 follower

Topic 2 → 1 Leader
2 follower

Topic 3 → 1 Leader
3 follower

Sending and Receiving Messages →

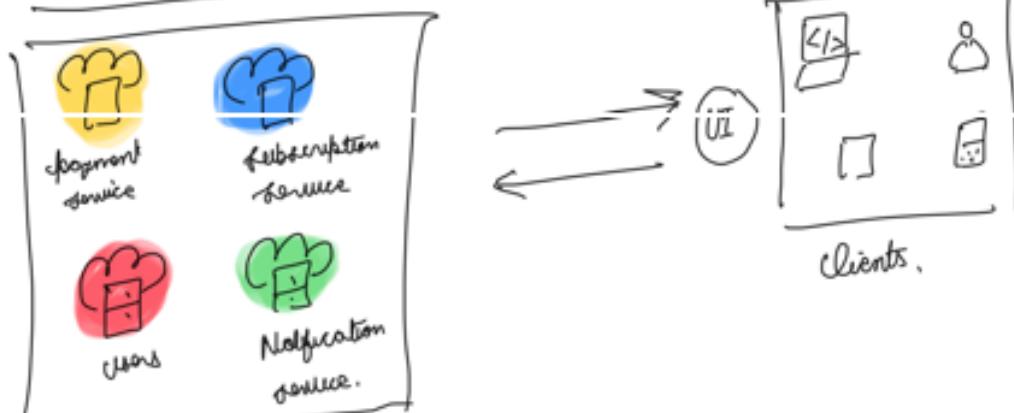


(19) API Gateway →

A Basic Example ↗



A Basic Example (Zoom In) →

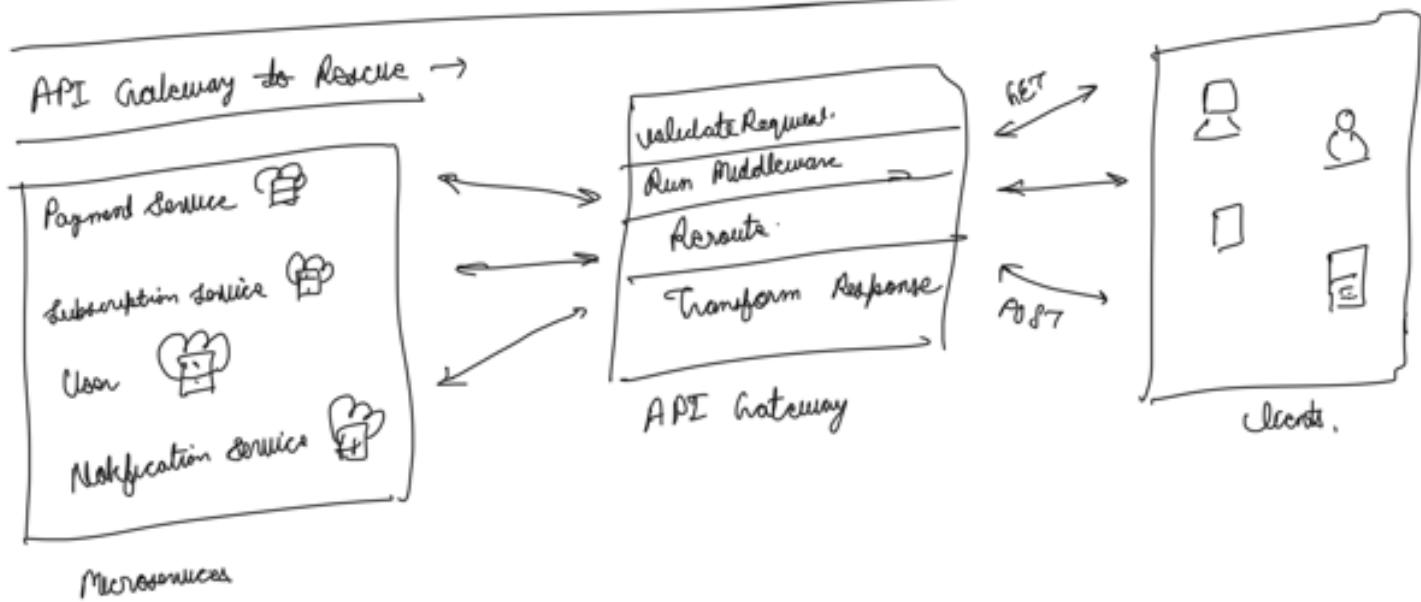


Microservice

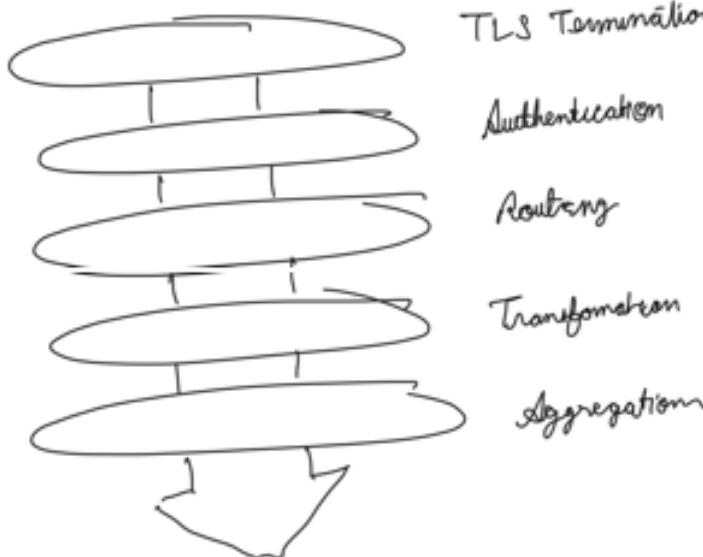
A chaos of Microservices →

- chaotic direct connection
- No unified Monitoring.
- scaling nightmares
- Security Vulnerabilities.
- inefficient data aggregation
- Redundant logic in service

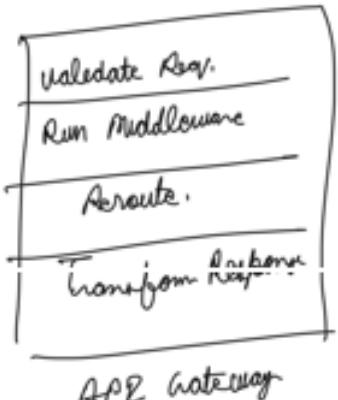
API Gateway to Rescue →



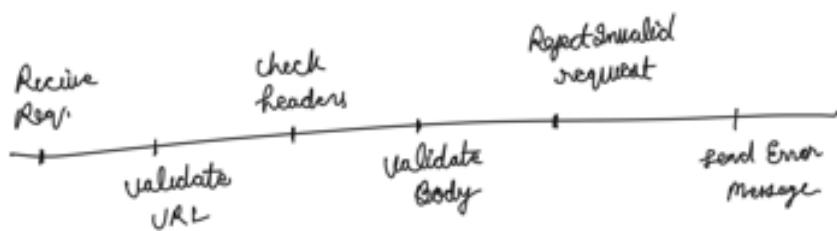
Filtering Visualized →



validate Req.
Run Middleware
Route
Transform Response



Request Validation process →



Middle Ware →

Monitoring
Compliance
Security
Performance.

- Authenticate req. using JWT Tokens
- Limit request rates to prevent abuse.
- Terminate SSL connections.
- Log and monitor traffic
- Compress responses
- Handle CORS headers
- Whitelist/Blacklist IPs
- validate request size
- Handle response timeouts
- Version APIs
- Throttle traffic
- Integrate with service discovery.

API Gateway

Pros	Cons
<ul style="list-style-type: none">→ centralized security→ Reduced Coupling→ Improved performance→ Scalability<ul style="list-style-type: none">- fail tolerance	<ul style="list-style-type: none">→ operational Complexity→ cost.→ Latency

→ Fault Tolerance

Which API Gateway should we choose?

Managed Service → offers seamless integration and ease of use but is more expensive. (Amazon API Gateway)

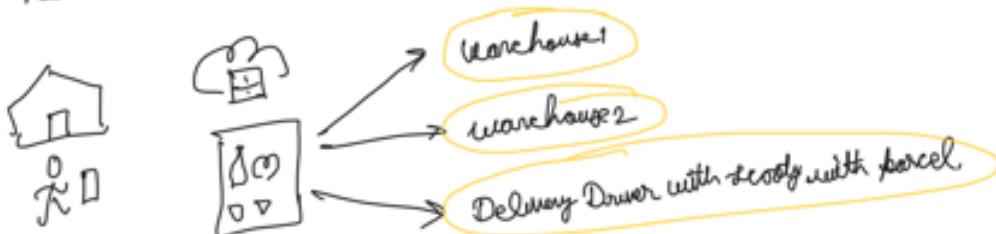
(Azure APG)

Open source sol'n → Provides more control and flexibility but requires more setup.

TYK & RONIN

Q) Design a Quick-Commerce Service →

The 3-minute magie.



System Requirements →

→ Availability connections → Connect to any data centre within one hour

→ Ordering Consistency → Ordering must be strongly consistent across all system

→ Order Volume → System must be handle one million orders per day

→ Order Items → Ability to order items through the system

→ Availability Speed → Fast availability with 100ms response time

→ System Size → Ten Thousand data centers, one hundred thousand catalog items

→ Query Availability → Query availability by location within one hour

Inventory Management Components →

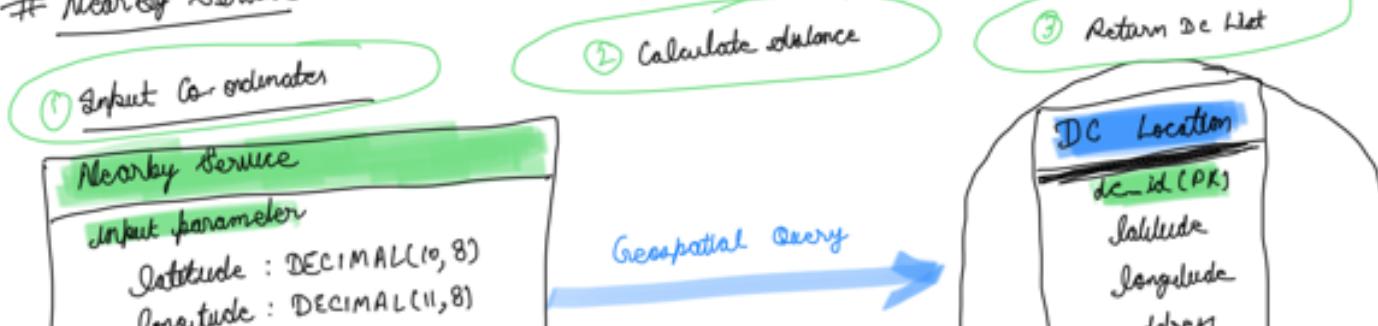
(i) Inventory → Physical Item instance at DC

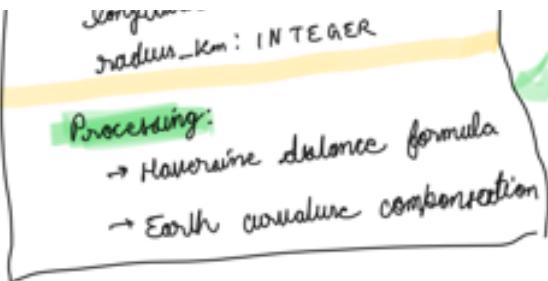
(ii) Item → Type of Item e.g. Cheetos.

(iii) Distribution Center → Physical location where item stored

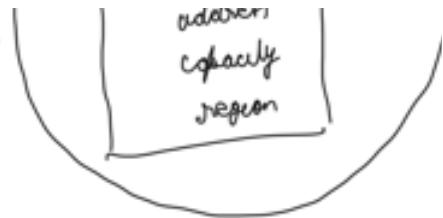
(iv) Order → Collection of Inventory ordered by user.

Nearby Service →

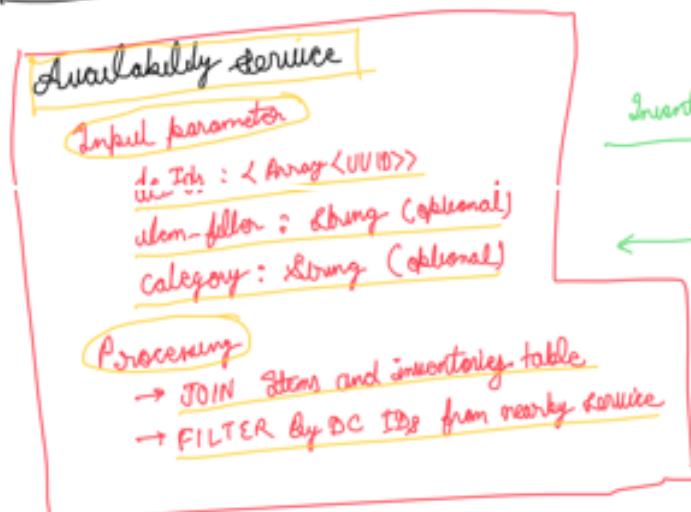




DC List Response



Availability Service →



Inventory query [DC-0]

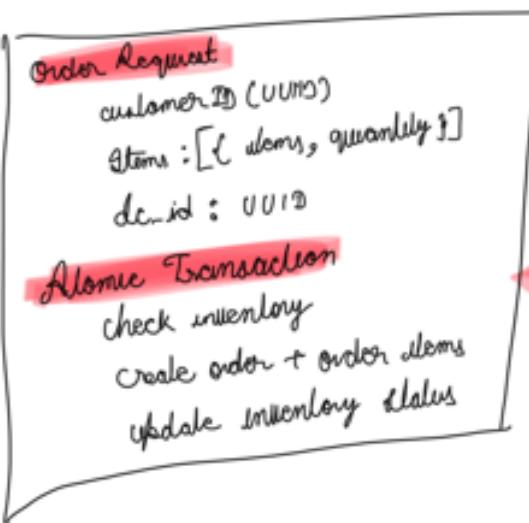
Available items [ITEM: QTY]

Distribution Center Database

Items	Inventory
id	id
name	item_id
description	quantity
category	reserved_qty
price	available_qty

Join on id

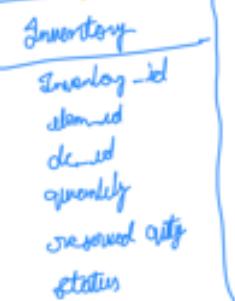
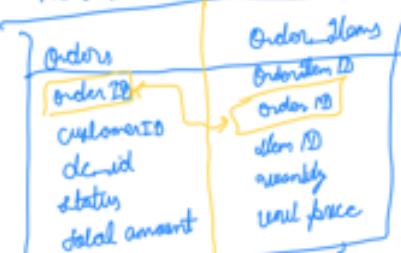
Order Service



Atomic transaction

order confirmation

PostgreSQL Database
ACID Transaction * Row-level Locking



Tying Together →



dark stock

find DC

place order

Availability Service
(Real time inventory check)

Location Service
(find nearest DCs)

Order Service
process order and payments

Query inventory

Location Service

Order Transaction

Glass	Inventory
order	order_glass
order_glass	order

Distribution Center

21 Chat System \Rightarrow

The Basic Need \rightarrow



Design Requirement

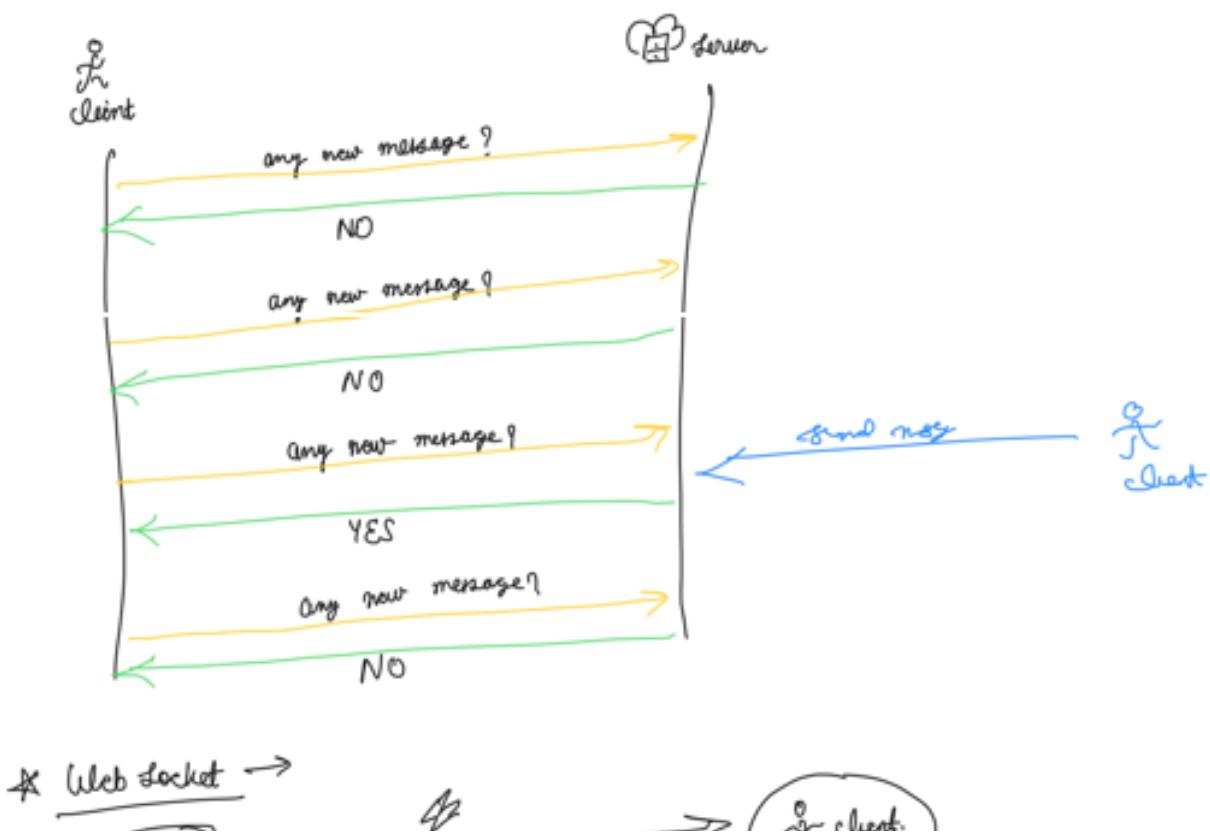
- ① One on One Chat \rightarrow focuses on personal communication with low latency.
- ② Group Chat \rightarrow support interaction among up to 100 members
- ③ Mobile App \rightarrow Accessible on smartphone and tablets
- ④ Web App \rightarrow Accessible through web browsers
- ⑤ Large Scale \rightarrow Designed to support 50 million daily active users.

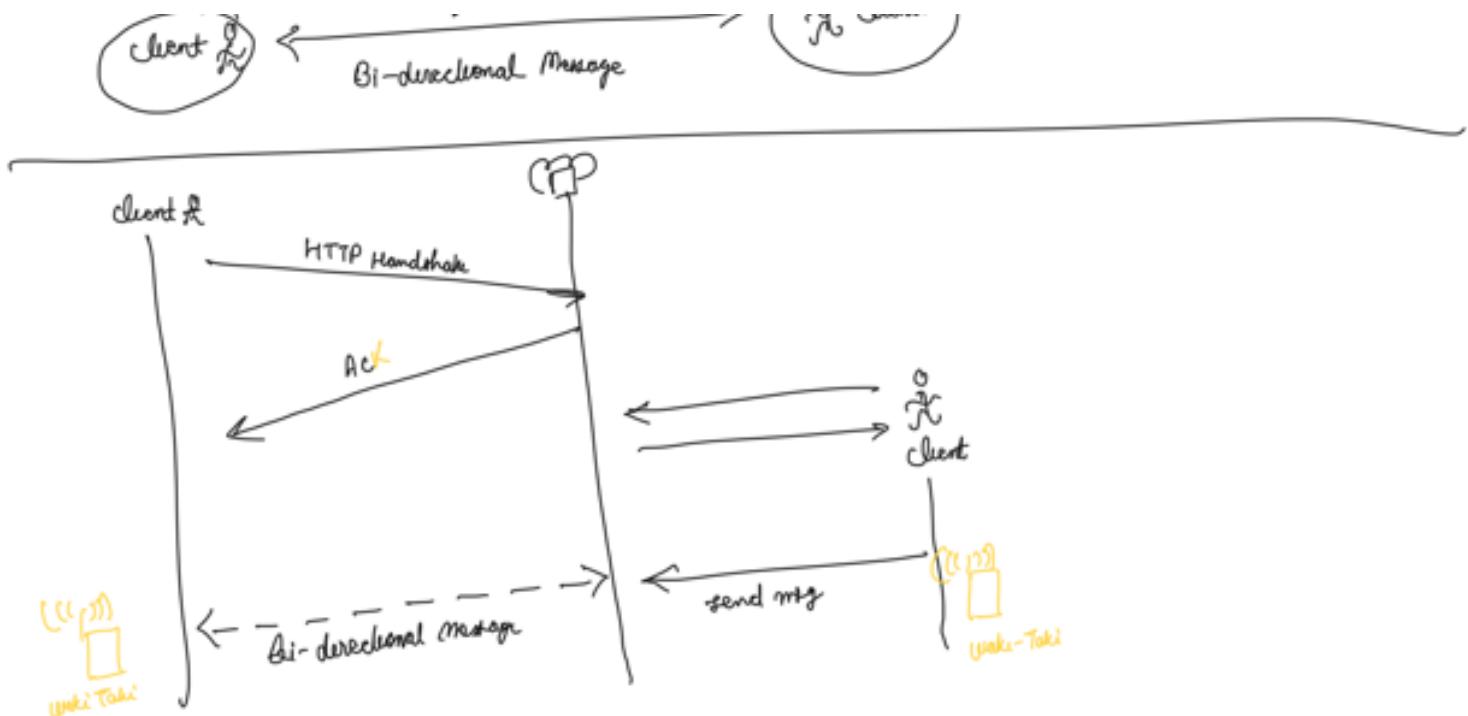
Communication Technique \rightarrow

Web Socket \rightarrow A persistent connection for real time communication

Polling \rightarrow Clients periodically check for new message from the server.

Long polling \rightarrow clients wait for a server response before checking again.

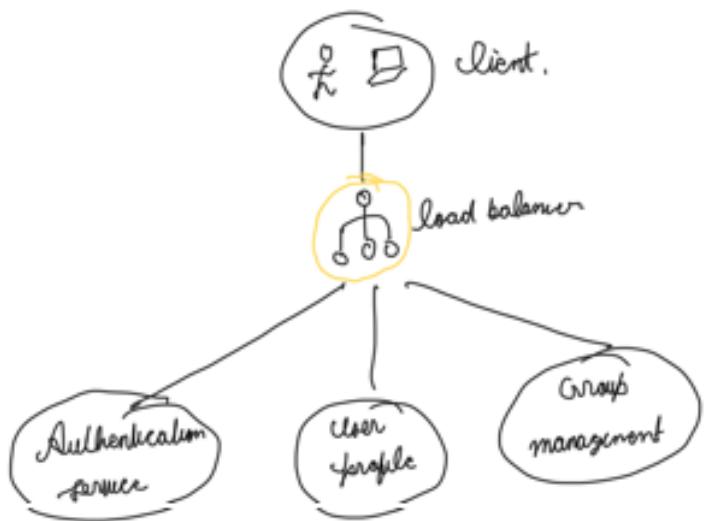




Stateless System →

Each request from the user is independent. Server does not remember any previous request.

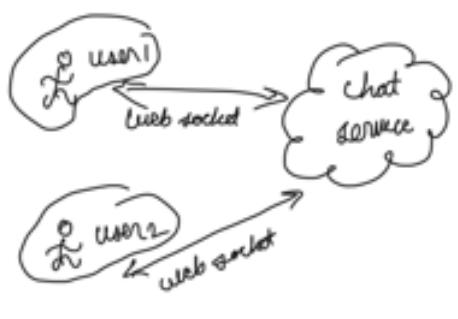
- All information must be sent with every request.
- No session information
- Easier to scale and load balance



Stateful System →

Server remembers the connection OR user session across multiple interaction

- Maintain ongoing context
- Usually needed for real-time communication
- Harder to scale.

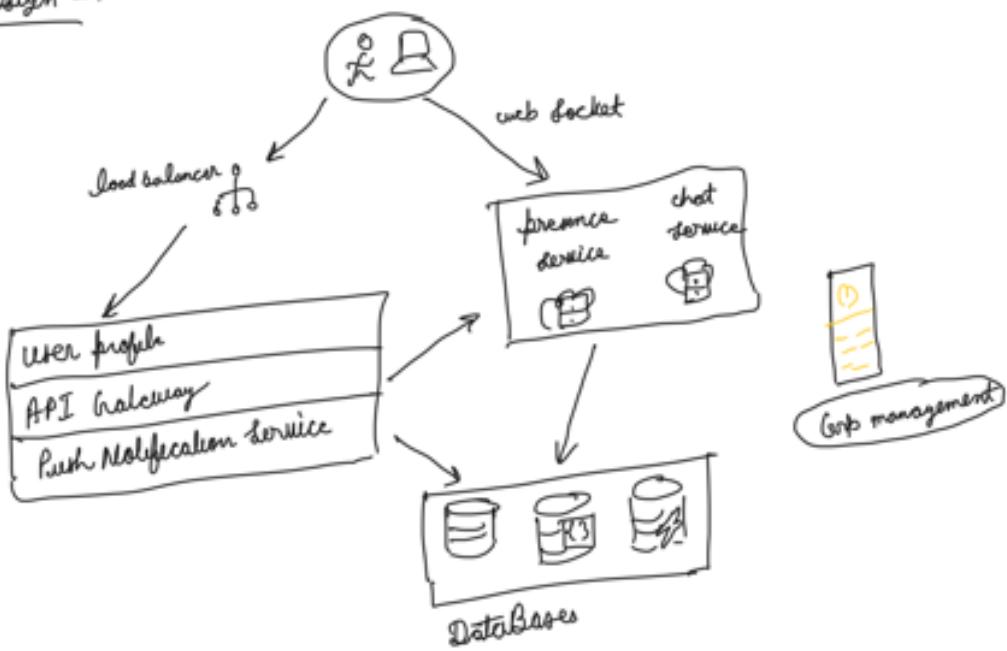


High Level Design →





* High Level Design →



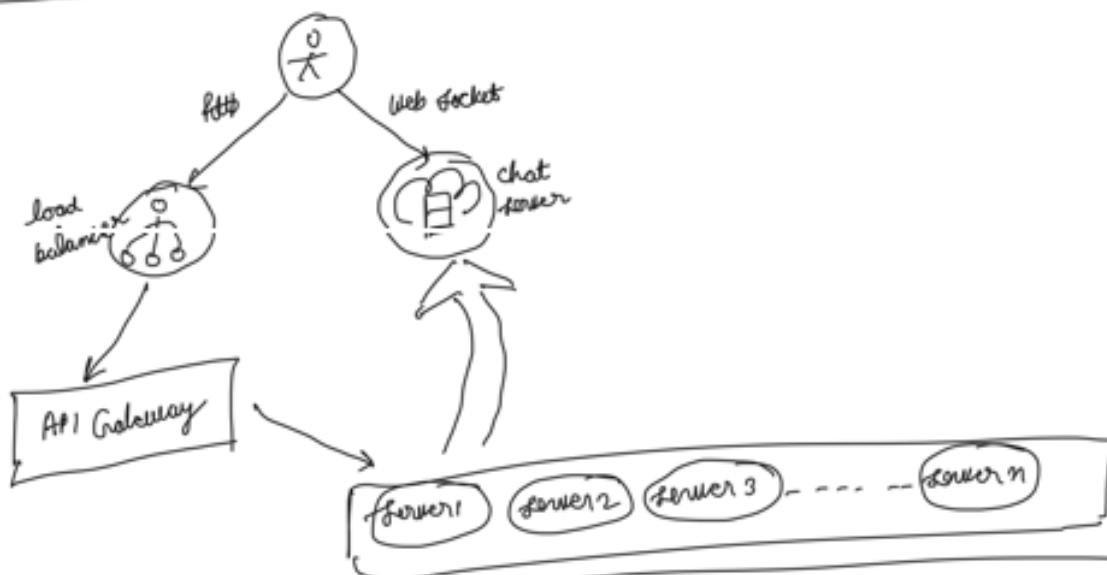
(2) Chat System (Part 2) →

* sending a Message →

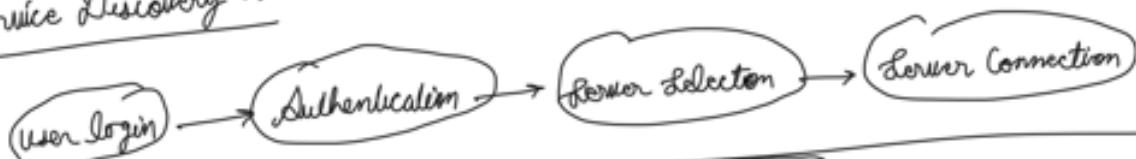
Message	
MessageID	BigInt
MessageFrom	BigInt
MessageTo	BigInt
Content	text
Created At	timestamp

- Auto Increment → Suitable for SQL databases but not for NoSQL.
- Global Sequence Generator → Ensures uniqueness and consistency across systems.
- Local Sequence Generator → Simpler implementation but limited access.

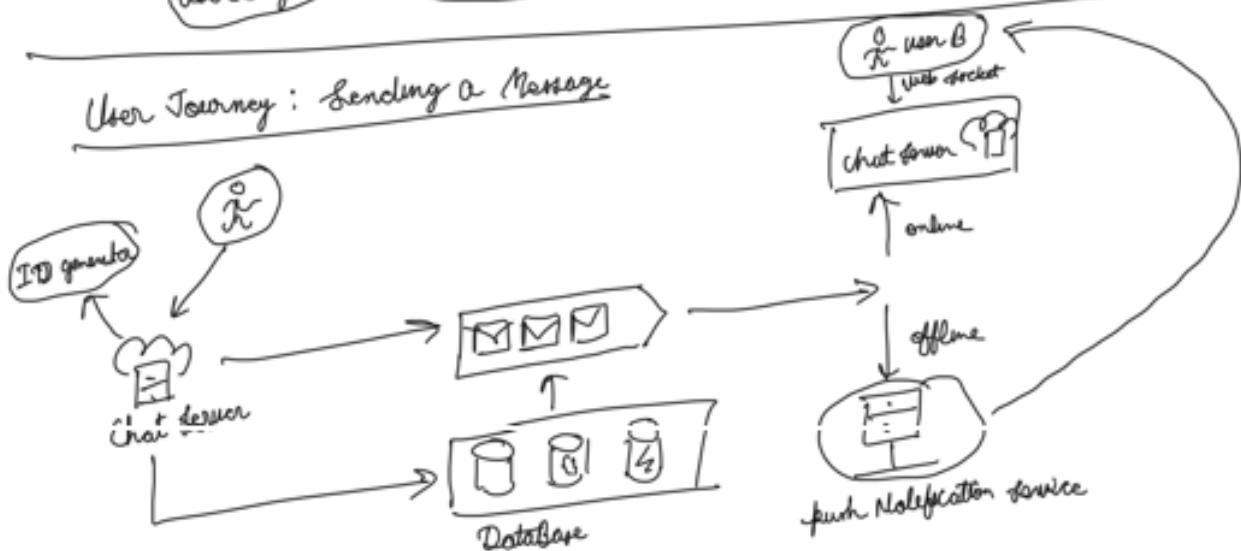
* User Journey →



* Service Discovery →



User Journey : Sending a Message



Group Message →





Synchronizing Device →

max message-id = 765

user's
module



web socket



user laptop

max message-id = 690

1 2 3 4 5 6 7 8 9 10 11 12

1 2 3 4 5 6 7 8 9 10 11 12

... 6 7 8 9 10 11 12

: