

1. Monolithic Architecture
2. Microservice Architecture
3. Single thread, Multithread
4. Synchronous, Asynchronous behavior
5. Proxy, Reverse proxy

Thread is a component of OS that executes the portion of code which are in the single control flow.

single thread programming language w/ Asynchronous behavior can output real time accurate output (ex: Node.js, ex: JavaScript).

Proxy: In a software which hides client from server.

Frontend API provides methods to interact with server that stores browser's outgoing requests to server.

Frontend makes requests to server. Server returns data to front end.

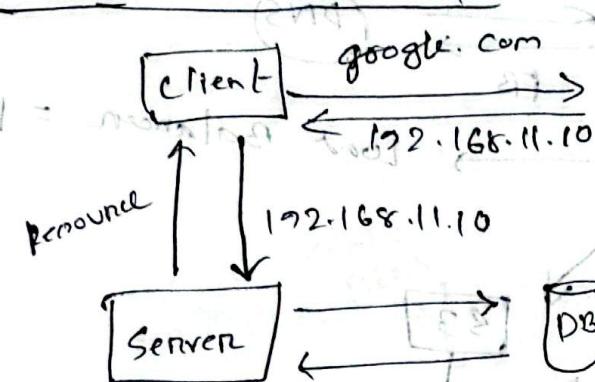
note that next section, Chatbot is not discussed.

Otherwise no connection exists to their API.

stub name

HTTP reverse proxy: Hides servers from clients

Single server setup:



There are two regions with problems with this setup.

→ ~~Central~~ single point of failure

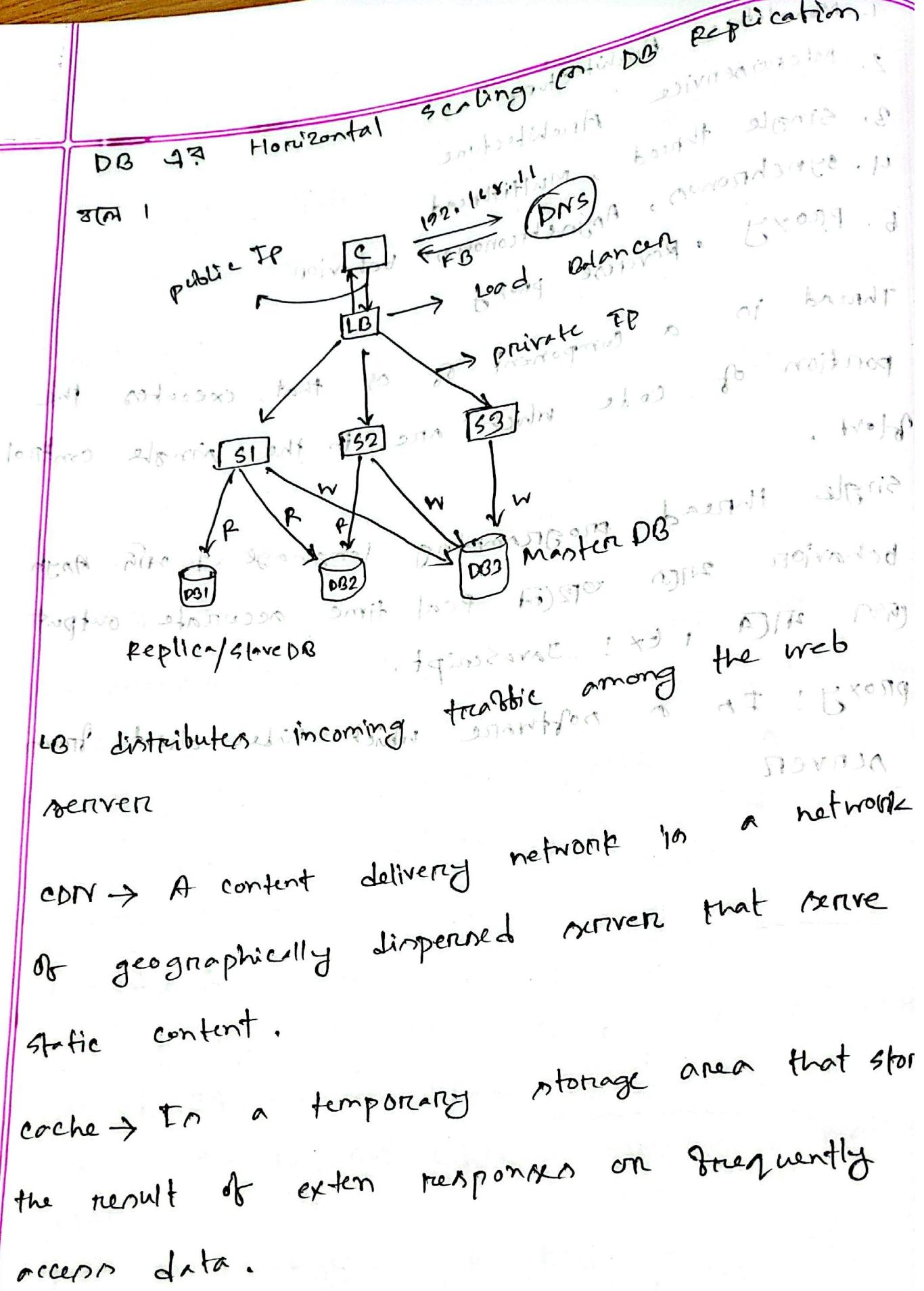
→ High latency

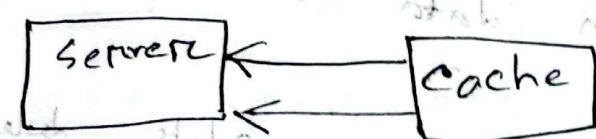
To solve this problem we need scaling:

Vertical scaling / Scale up: Means adding more power to our CPU, getting rid of single point of failure

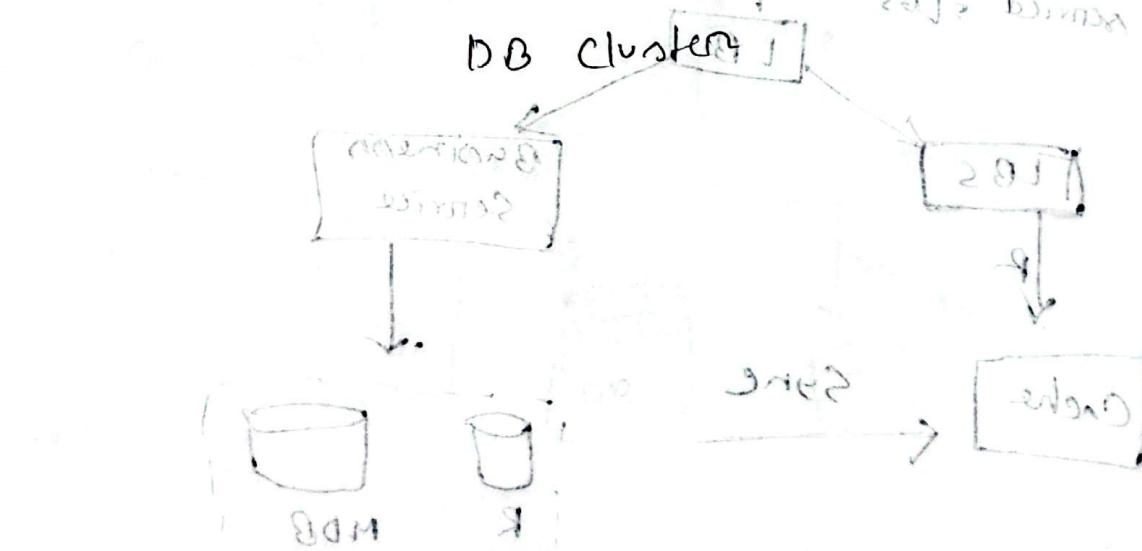
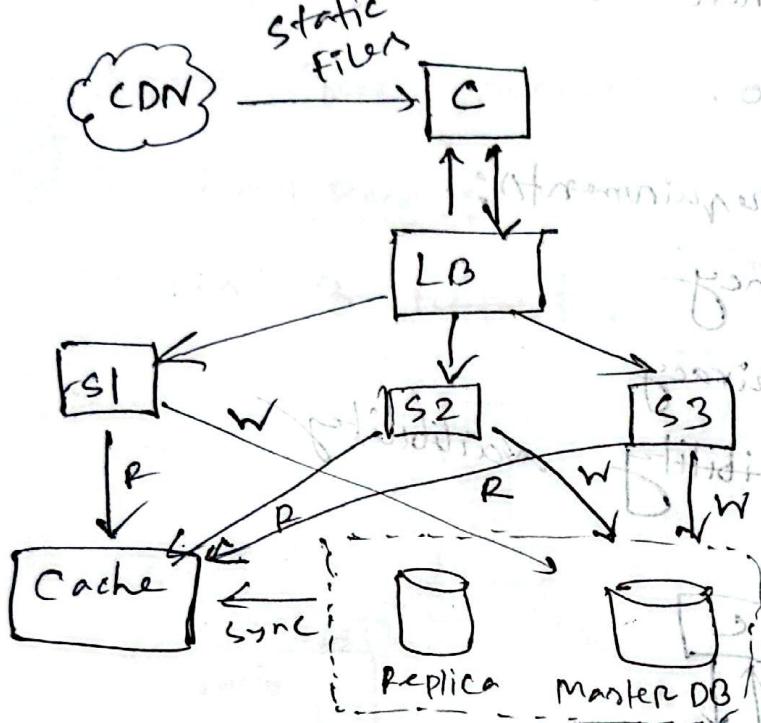
Horizontal scaling / Scale out: . that's what

means adding more resources to our pool of workers. Both problems have been solved in this





If data is not in Cache, save in Cache



proximity service design

functional requirements:

1. Return all business location.

2. Business owner can add, update, delete their business info.

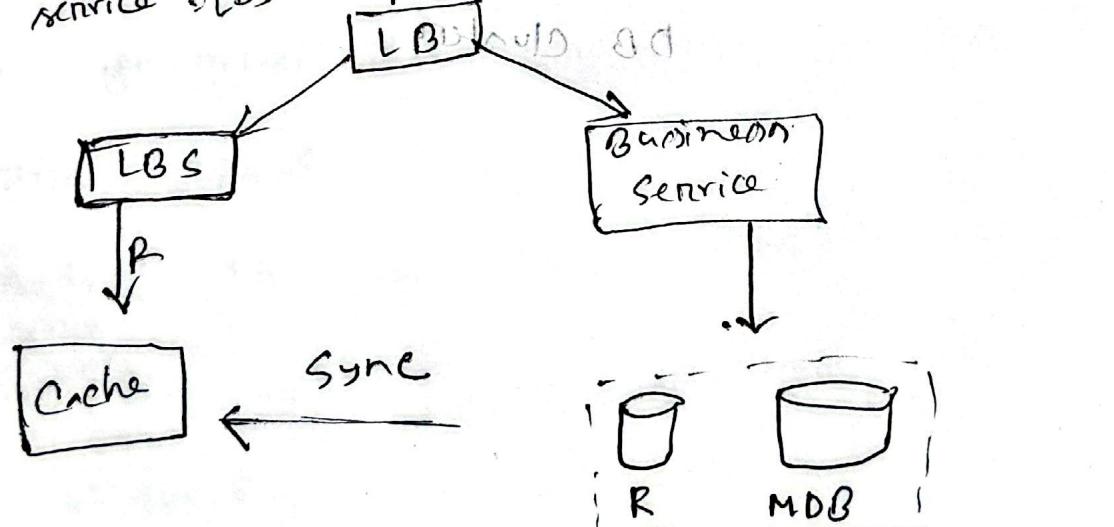
non-functional requirements:

1) low latency

2) data privacy

3) High availability

Location Based Service or LBS



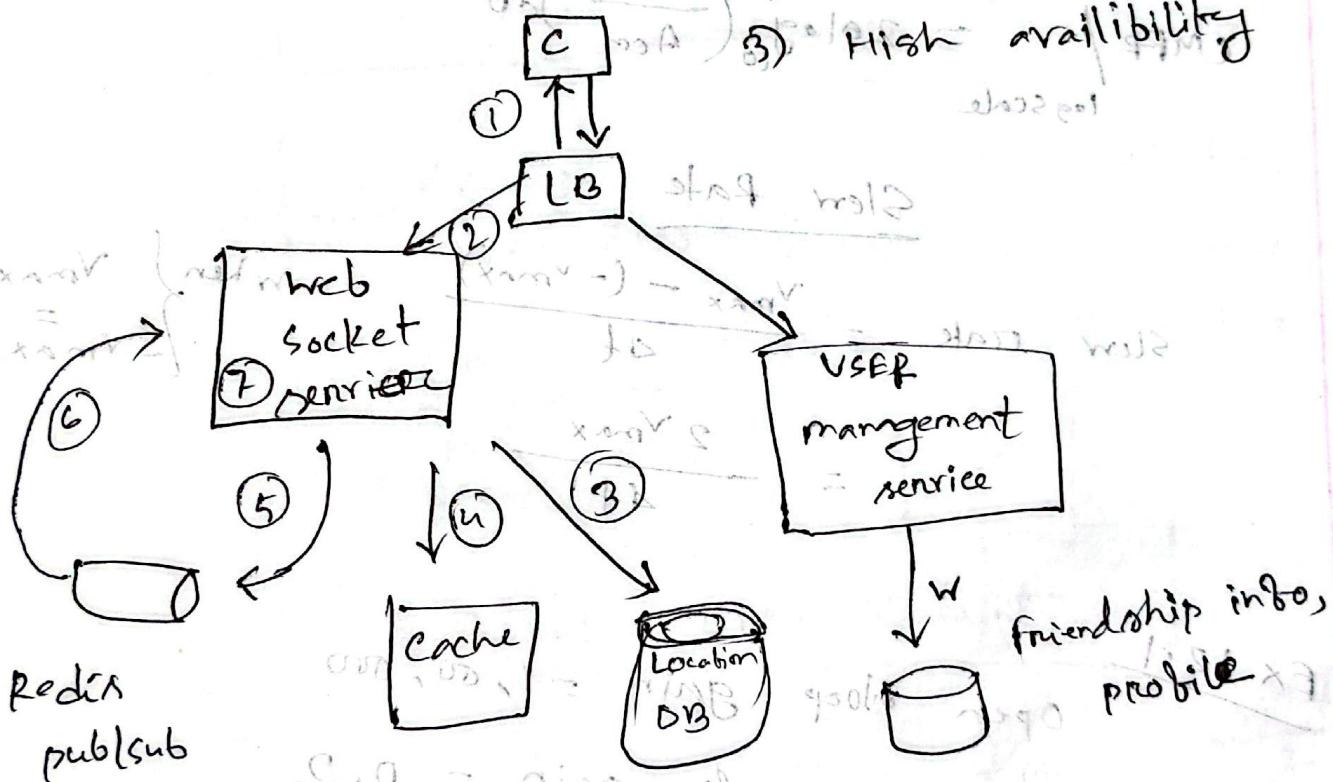
(with existing Nearby friends) subsystem design

Functional reqs.

1. User can see their nearby friends on the mobile app
2. Nearby friends list should be updated every few seconds

Non-functional Req.:

- previous 3 repeat.
- D) low latency
 - 2) data privacy
 - 3) High availability



software

1. The mobile client sends location update to the load balancer.
2. The load balancer forwards the location update to the persistent connection on the websocket server. From that client server, the websocket server saves the location data to the location DB.
3. Web socket server receives update from the cache with new location.
4. Web socket server publishes the new location to the wschannel channel. It is (other medium) using pub/sub receiver.
5. The wschannel receiver publishes the new location to the wschannel channel.

6. Redis pub-sub server broadcast the update to all the subscribers.

7. web-socket server (friends) computes the distance between the user and the subscriber. If the distance doesn't exceed the search radius, the location is sent to the subscriber client.

System design of google maps

Functional Requirements:

- i) User's location update time
- ii) Navigation service including estimated time of arrival.
- iii) Map rendering

Non-functional requirements

Repeat all three

1) low latency

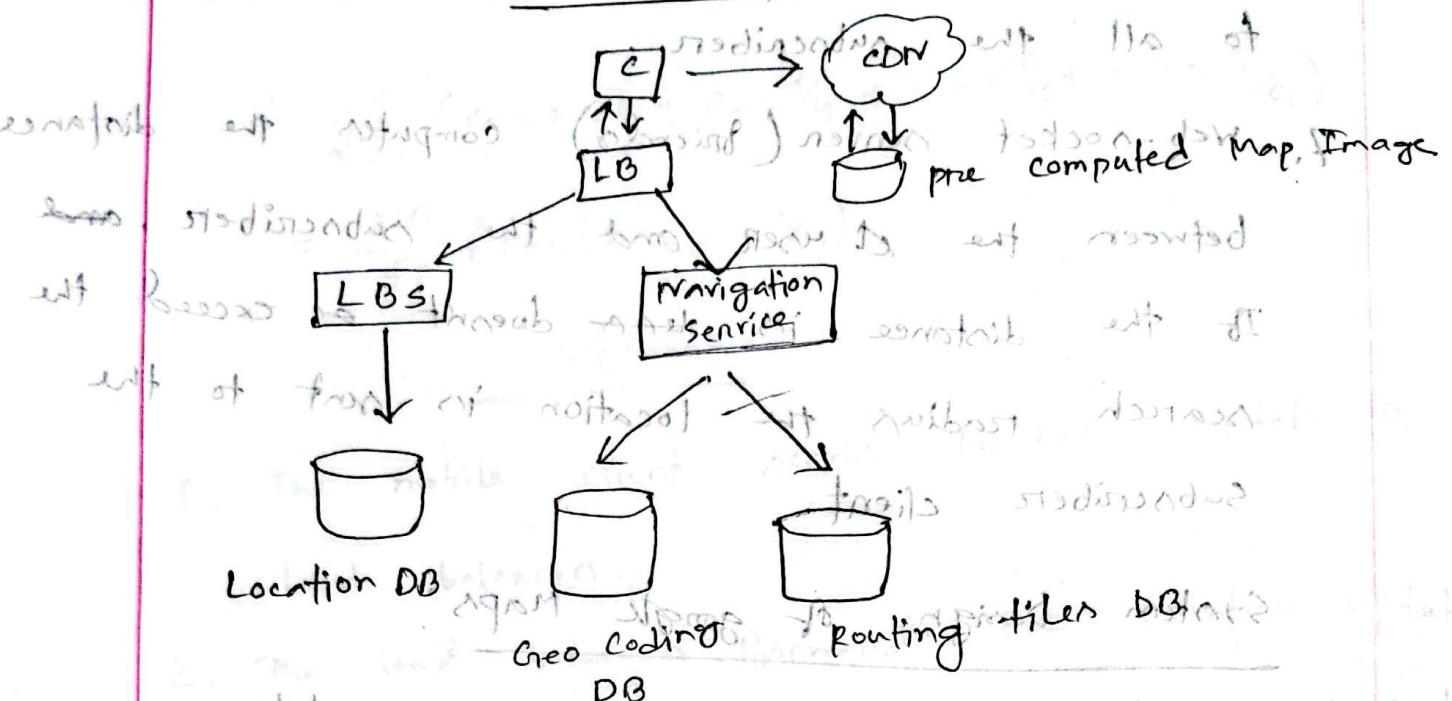
2) High availability

3) data privacy

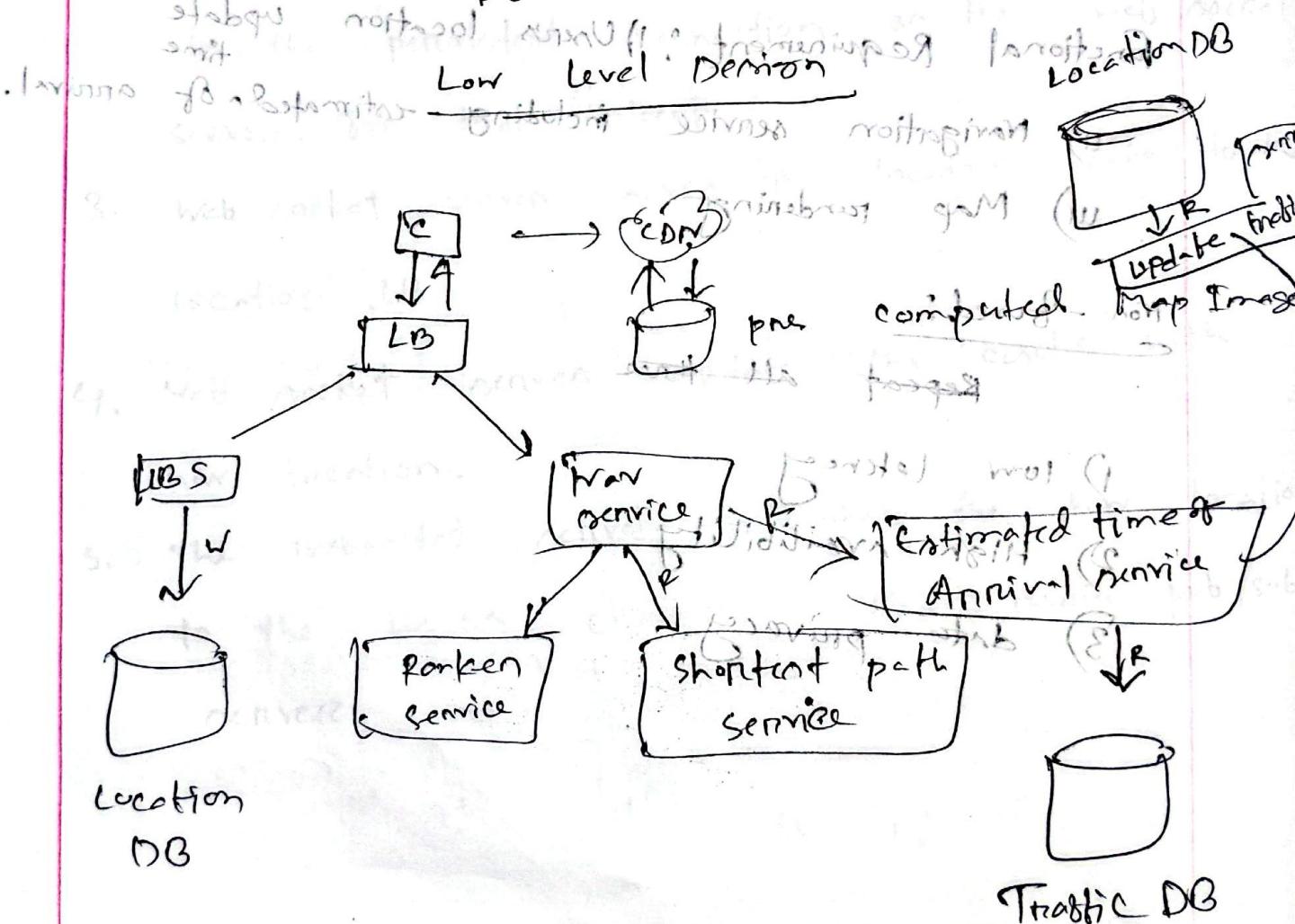


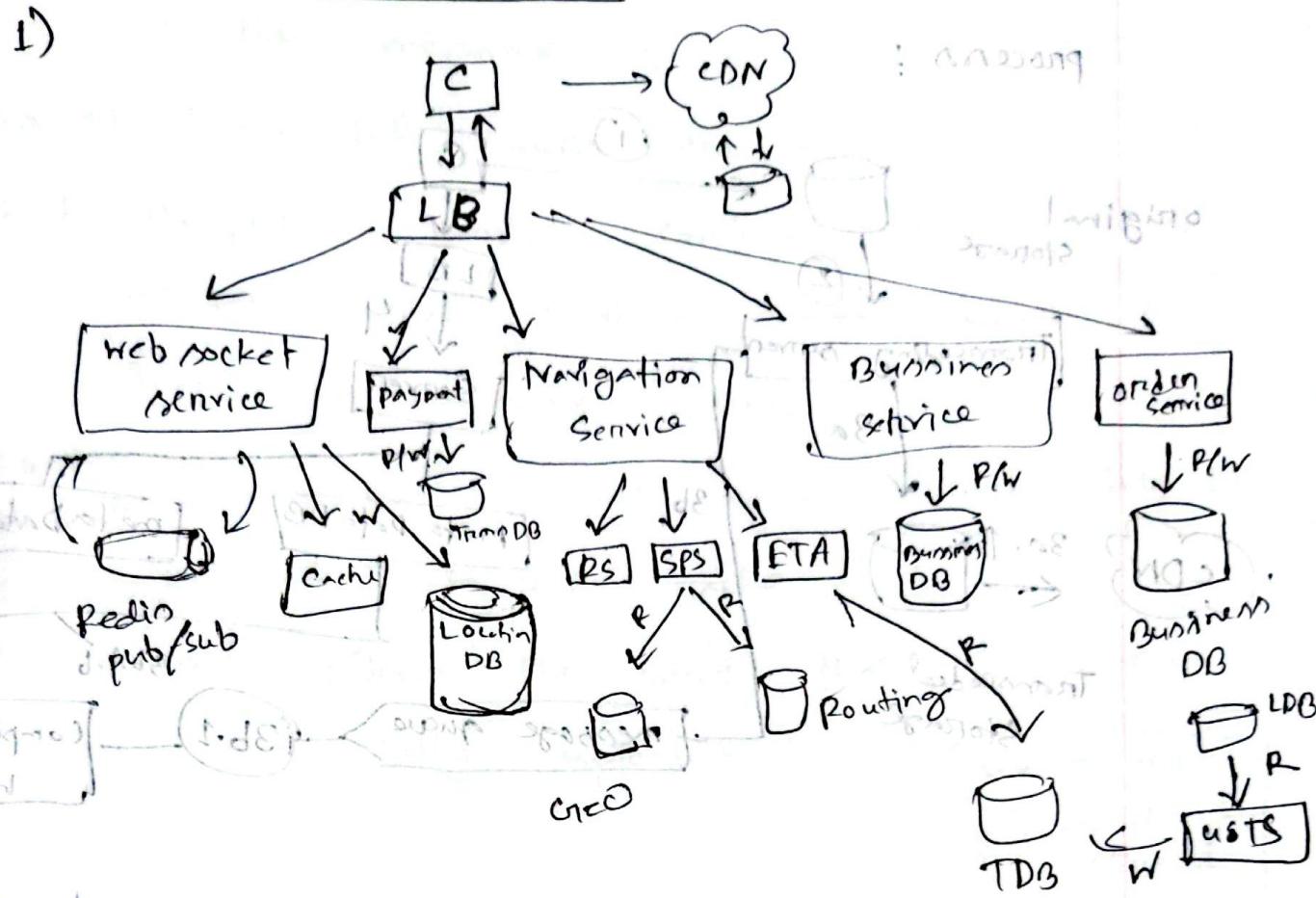
Different

High level Design



Low level Design

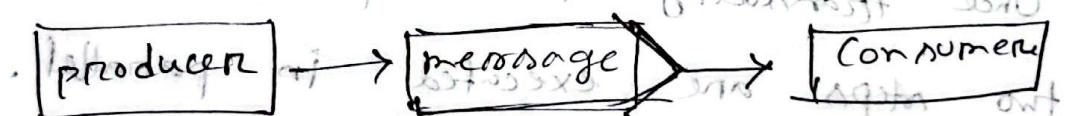




Message queues: A message queue is a durable store holding a set of messages.

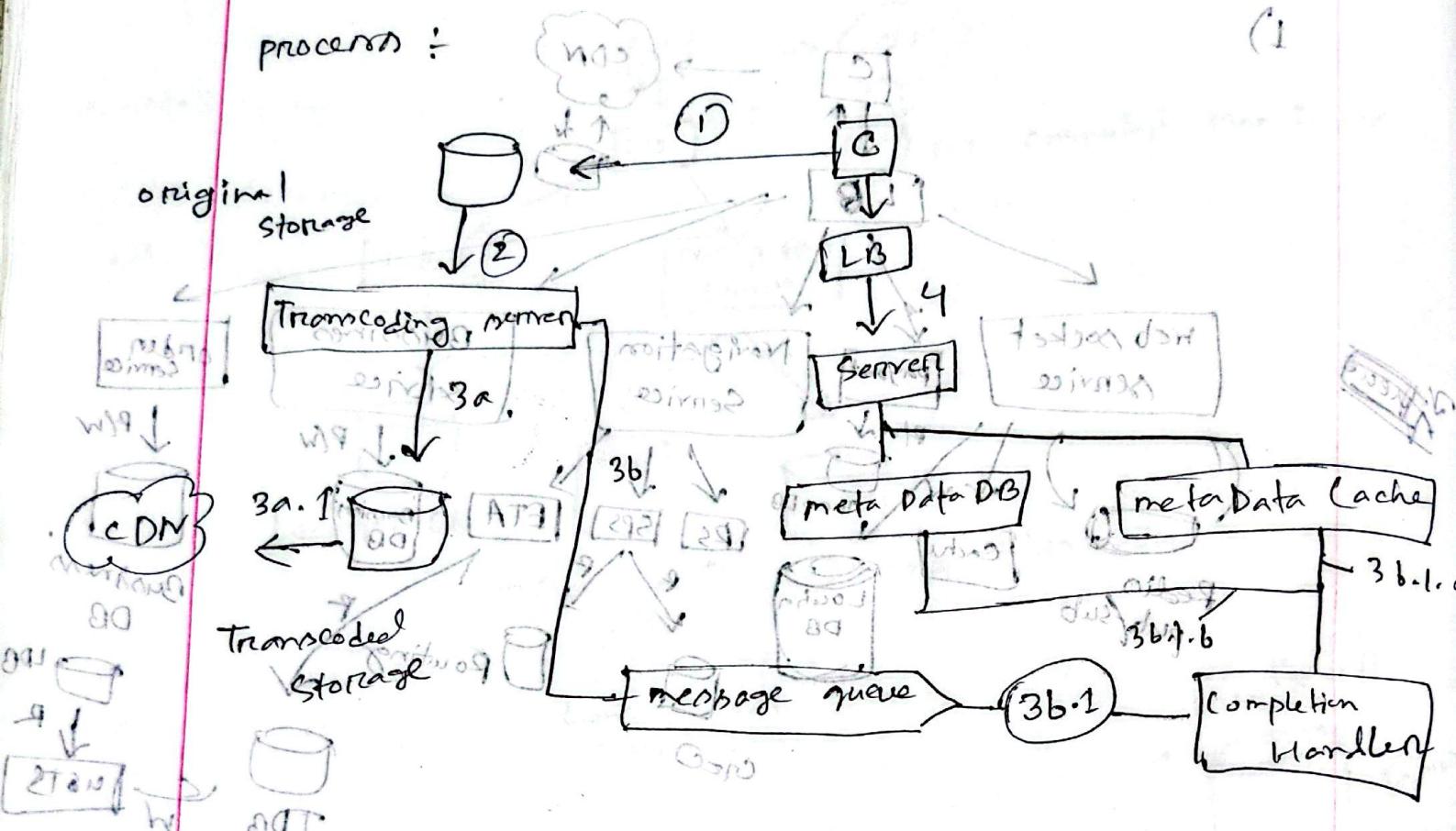
Component stored in a memory box that is used

to process any nonconform requests.



System ~~video~~ Design of youtube's video uploading process :

(1)



1. Videos are uploaded to the original storage.
2. Transcoding server's fetch video from original storage and starts transcoding.
3. Once transcoding is complete, the following two steps are executed in parallel:
 - 3a. Transcoded videos are sent to transcoded storage.
 - 3b. Two steps are executed in parallel:
 - 3b.1. Transcoded videos are stored in transcoded storage.
 - 3b.2. Transcoded videos are sent to CDN.
 - 3b.3. Transcoded videos are sent to meta Data Cache.
 - 3b.4. Transcoded videos are sent to Server.
 - 3b.5. Transcoded videos are sent to meta Data DB.
 - 3b.6. Transcoded videos are sent to Completion Handler.

Transcoding → Device dependent to independent

3b. Transcoding completion events are queued

in a list in other message queue.

3b.1) Transcoded video are distributed to CDIV.

3b.1) Completion handler contains a bunch of tasks of workers that continuously pull event data from the message queue.

3b.1.a) Completion handler updates meta data DB
3b.1.b) Completion handler pushes to meta data cache.

4. servers inform the client that the video is successfully uploaded.

5. notification system starts a background task

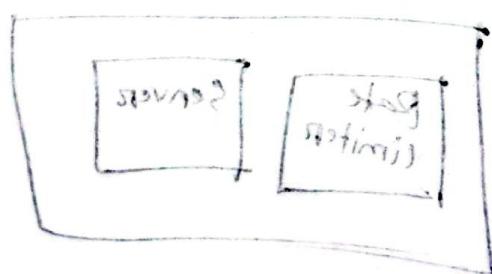
workers are kind of threads

servers set

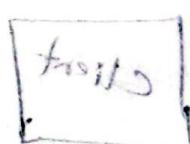
adding addition to skin traits



initial trait



; skin moves



et tragschicht direkt zu Problemlösung
fragegestellt;

System eignen sich für Rate Limiter:

In a network system, a rate limiter is used to control the rates of traffic and by a client/server. In the HTTP world, a rate limiter limits the number of client requests.

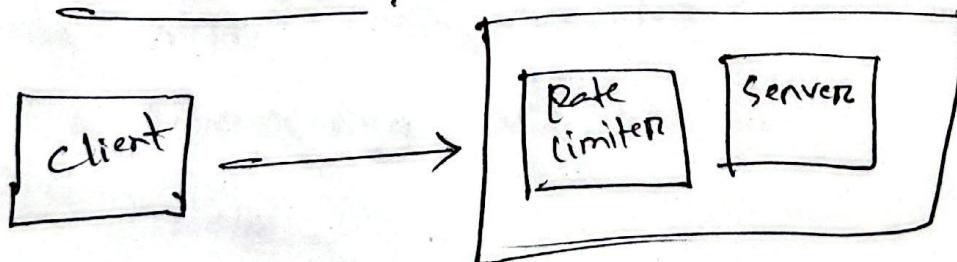
Benefits of using a rate limiter:

Denial of service prevents resource starvation caused by denial of service attack.

We can implement a rate limiter either a client side or server side

Client side: An unreliable place to enforce rate limiting.

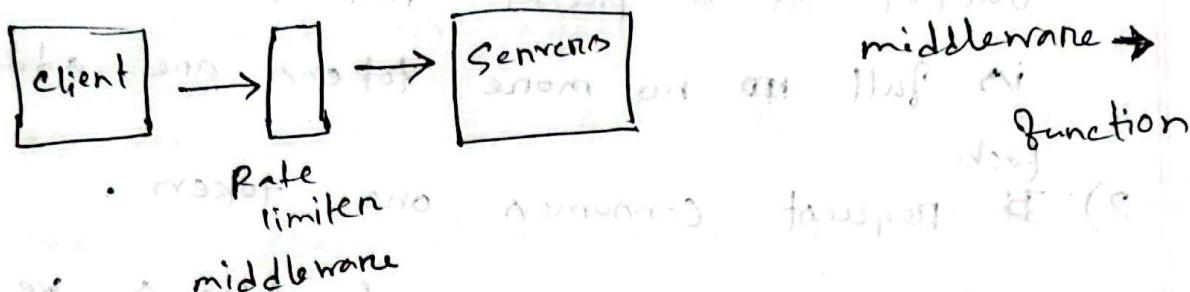
Server side:



2/3

Software

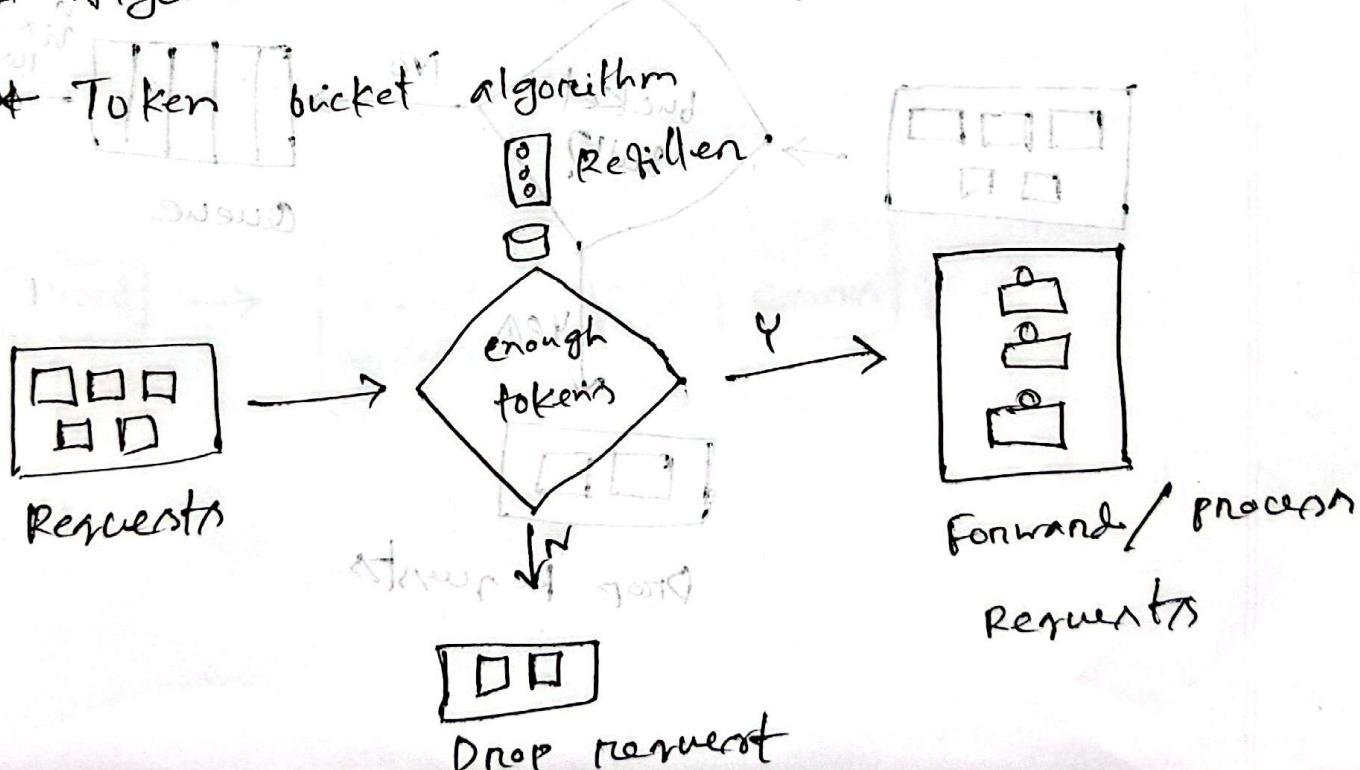
Instead of putting a rate limiter at the server side we can create a rate limiter which throttles request to our server.



Let's assume our API allows two requests per second. and a client sends 3 requests within a records. The first two requests are routed to servers, however the rate limiter throttles the third request and returning HTTP status code 429.

~~Algorithm for rate limiting~~

* Token bucket algorithm



Alg 0.

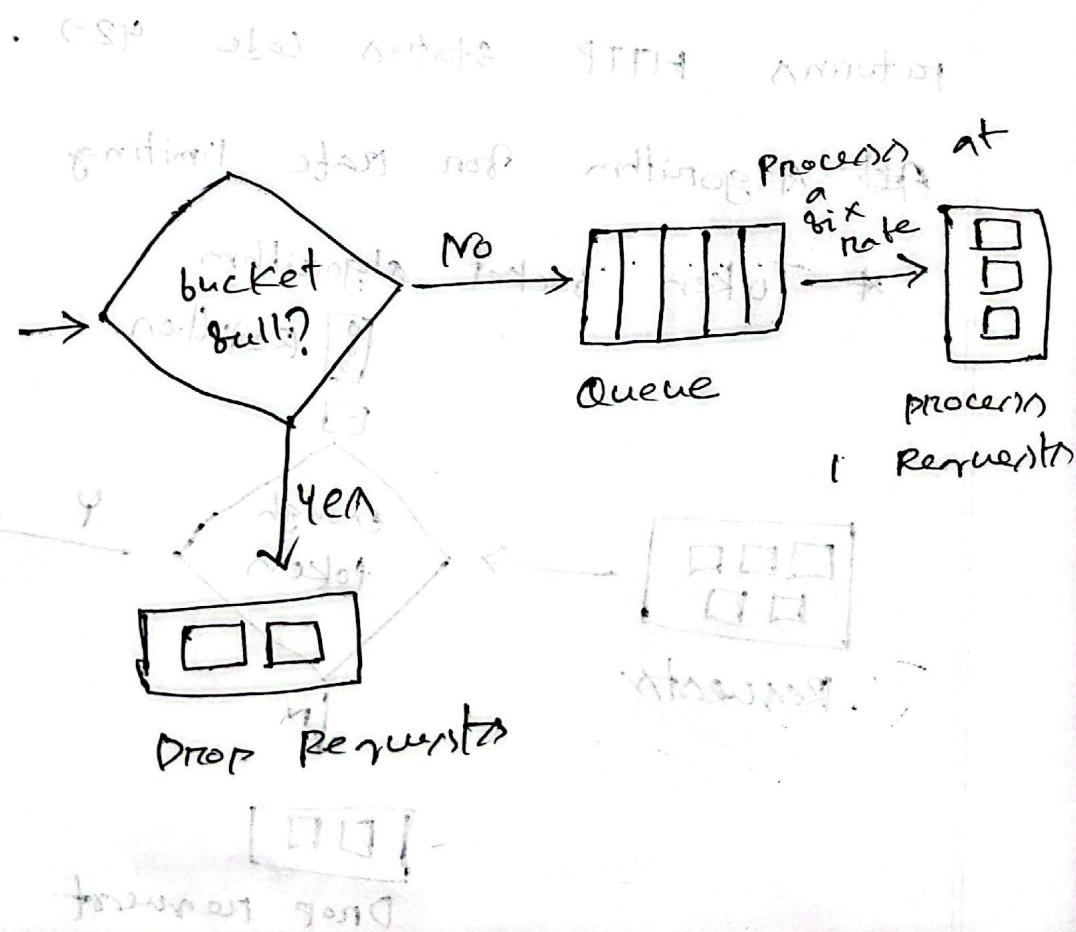
The token bucket algorithm follows:

- 1) A token bucket is container that has predefined capacity. Tokens are put in the bucket at a preset rates. Once the bucket is full ~~no more~~ tokens are added.

- 2) Request consumer one token.

- 2.1) If there are enough tokens we take one token out for each request.
 - 2.2) No token request drop.

Leaking Bucket Algorithm:

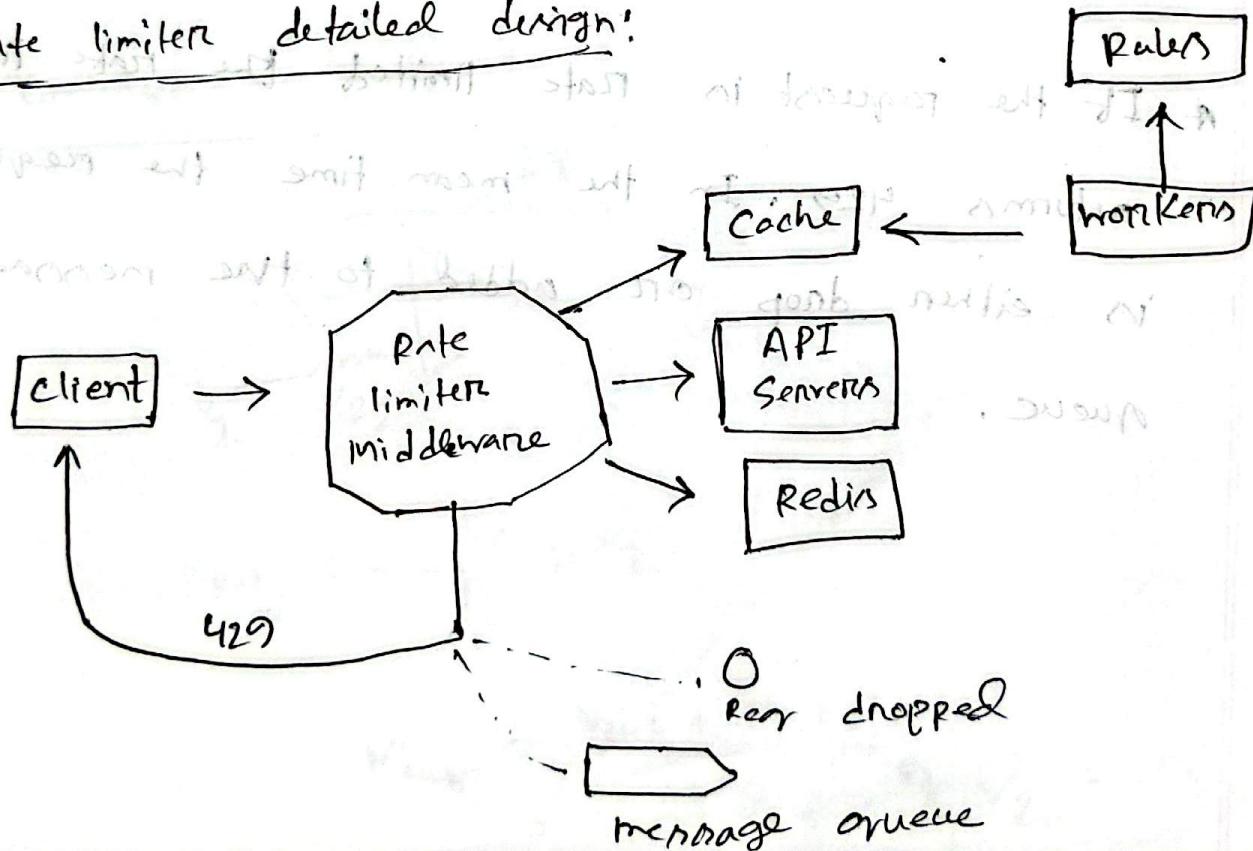


The algo works as follows:

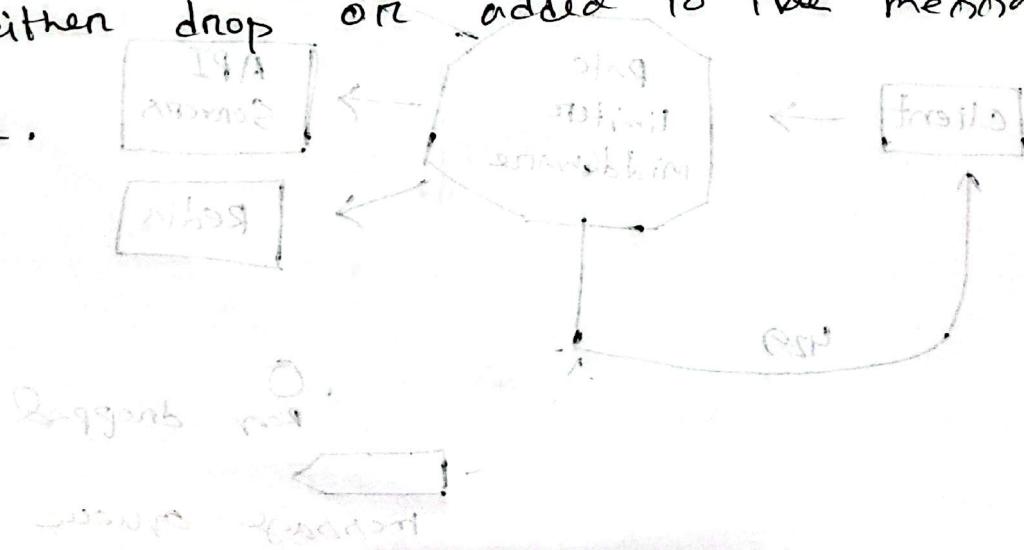
- 1) When a request arrives the system checks if a the queues is full. If it is not full the request is added to queue.
- 2) If the queue is full the request is drop.
- 3) Requests are pulled from the queue and process at a regular intervals.

[N.B If a user is allowed to make one facebook friend per second, add 150 a per day, like 50 per second. Three buckets are required for each user. Annotate LIA of behaviour of if]

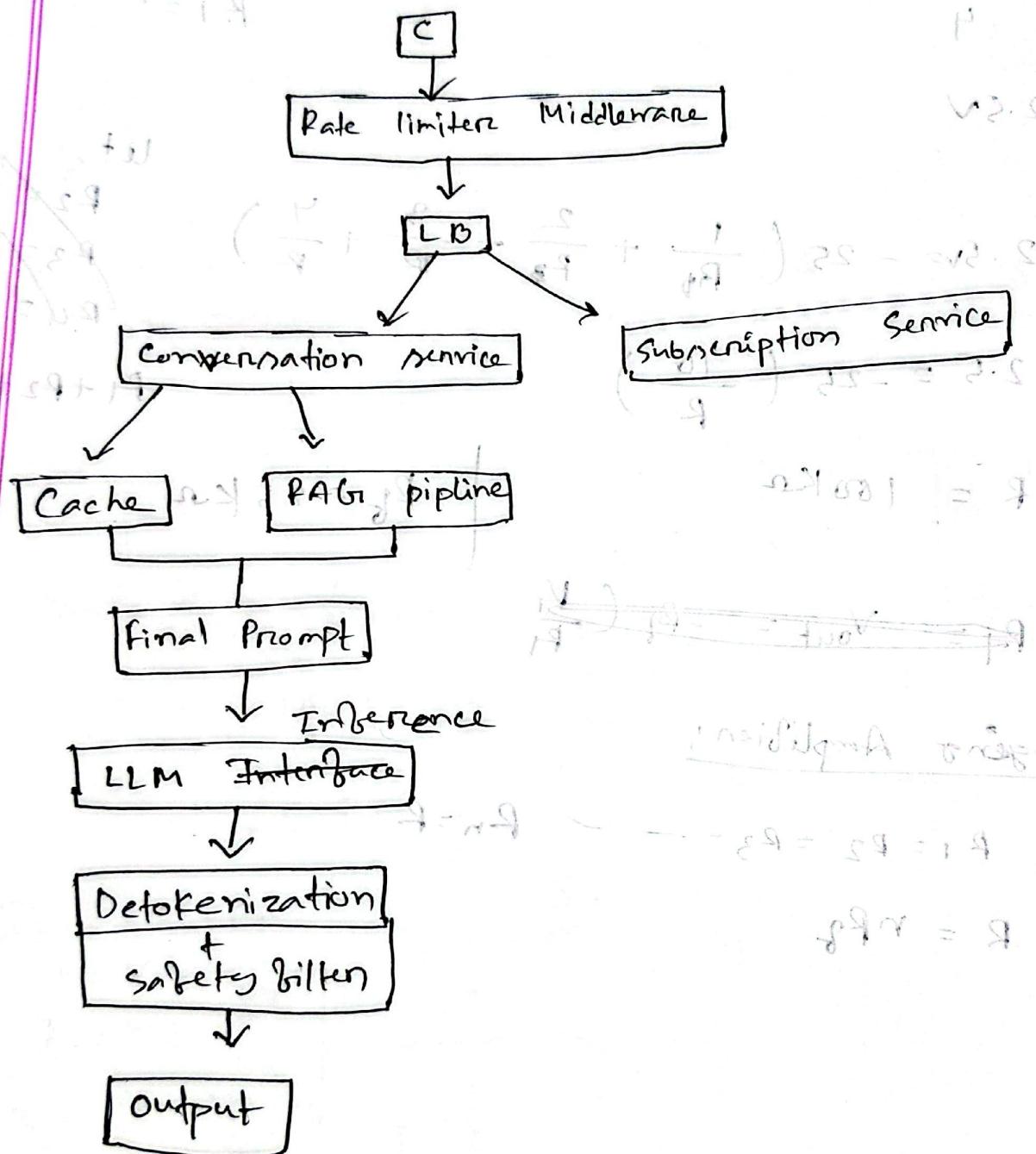
Rate limiter detailed design:



- * Rules are stored on the disk. workers frequently pull rules from the disk and stored them in the cache.
 - * When a client sends a request to the server, the request is sent to the rate limiter middleware first.
 - * Rate limiter middleware loads rules from the cache. It fetches the value of counter variable and last request time-stamp, from redis.
 - * If the user request is not rate limited it is forwarded to API servers.
- Rate Limiter
- If the request is rate limited the rate limiter returns 429. In the mean time the request is either drop or added to the message queue.



System Design of chatgpt



How LLM works!

