

① Back of the Envelope Estimation →

(Estimates) + (Thought Experiment) + (Common performance Numbers)

→ which design will meet your requirements.

→ Good sense of scalability.

→ Following concepts should be aware →

↳ power of 2

↳ Latency Numbers

↳ Availability Numbers

<u>Power</u>	<u>Approximate Value</u>	<u>Full Name</u>	<u>Short Name</u>
10	1 Thousand	1 Kilobyte	1 KB
20	1 Million	1 Megabyte	1 MB
30	1 Billion	1 Gigabyte	1 GB
40	1 Trillion	1 Terabyte	1 TB
50	1 Quadrillion	1 Petabyte	1 PB

→ Latency Numbers →

<u>Operation Name</u>	<u>Time</u>
L ₁ Cache Reference (Level 1 cache Reference)	0.5 ns
L ₂ Cache Reference (Level 2 cache Reference)	7 ns
Main Memory Reference	100 ns
Compress 1 K bytes with Zipper	1000 ns = 10 μs
Send 2 K bytes over 1 Gbps Network	2000 ns = 20 μs
Read 1MB sequentially from memory	250 000 ns = 250 μs
Round trip within the same datacenter	500 000 ns = 500 μs
Disk Seek	10 000 000 ns = 10 ms
Read 1 MB sequentially from the N/W	10 000 000 ns = 10 ms
Read 1 MB sequentially from disk.	20 000 000 ns = 20 ms
Send packet CA (California) → Netherlands → CA	150 000 000 ns = 150 ms

After seeing latency numbers →

→ Memory is faster than disk. (RAM fast होता है डिस्क से)

→ Avoid disk seek if possible (Cache का use करें) (Memory (Ram) जावा होता है)

→ simple compression algorithm are fast (compress कर दो)

→ compress data before sending it if possible

→ Data centers are far located (Time taking)

Availability Number →

- High availability = system continuously operational for long.
- Measured in % ⇒
 $100\% \text{ availability} = 0 \text{ down time}$
 Mostly service = 99% to 100%

SLA (Service Level Agreement) →

- Agreement b/w service provider and consumer/ customer.
- ~~for all~~ available $\geq 99\%$ (uptime)

e.g. → Amazon Google Microsoft

Cloud service = $99.9\% \uparrow$

Availability	Downtime per day	Downtime per year
99 %	14.40 minutes	3.65 days
99.9 %	1.44 minute	8.77 hours
99.99 %	8.64 seconds	52.60 minutes
99.999 %	864.00 millisecond	5.26 minutes
99.9999 %	86.40 millisecond	31.56 seconds

② Twitter QPS and Storage Estimation →

Q → Estimate Twitter's QPS (Query per second) and storage requirement

Note → Following numbers are taken e.g. only

~~Assumptions~~ →

V.V. Jmp

Monthly active user = 300 million
Daily users = 50 %
Tweets per day on avg per user = 2
Tweets containing media = 10 %
Tweets stored for = 5 years

Estimations →

Query per second (QPS) estimate →

→ Daily Active Users (DAU) = 50% of [300 million] = 150 million.

→ Tweets QPS = $150 \text{ million} * 2 \text{ Tweets} / 24 \text{ hours} / 3600 \text{ seconds}$
 $\approx 3500 \text{ Tweets/sec}$

→ Peak QPS (100% user) = $(2 * \text{QPS}) \approx 7000 \text{ Tweets/sec}$

Storage Requirement Estimation →

Note → we will only estimate media storage here

→ Average tweet size →

→ tweet_id = 64 bytes

→ text = 140 bytes

→ media = 1MB

→ Media storage = (Avg day active users) * (no of tweets per user) * (10% contain media) * (size of media)

$$= (150 \text{ million}) * (2) * (10\%) * (1 \text{ MB}) = (30 \times 10^6) = (30 \text{ TB}) \text{ per day.}$$

$$\therefore 1 \text{ TB} = 10^6 \text{ MB}$$

$$5 \text{ year media storage} = (30 \text{ TB}) * (365) * (5) = (54750 \text{ TB}) \approx 55 \text{ PB}$$

$$\therefore 1 \text{ PB} = 1000 \text{ TB}$$

→ Interview Tips →

→ Rounding & Approximation → $99989 / 9.2 \times \cancel{\textcircled{X}}$
 $\approx 10000 / 10 \checkmark$

→ Always write down your Assumption

→ Label your units → size 5 ~~☒~~
size 5 MB

→ Commonly asked question → QPS, Peak QPS, storage, Cache, number of servers etc.

(3) Red Flags and Myths →

Dream Company → selected for interview

- ↳ DS
- ↳ Algo
- ↳ Space & time Complexity
- ↳ System Design.

→ Intimidating after

→ Design a "well known product "ABC"

→ How? 1 hour? → Man, it was built by 100-1000s of engineers.

Good News →

No one expects you to design a perfect real world system in 1 hour

e.g. Real world system design is complicated

Google Search :- Deceptively simple but amount of technology under it is astonishing

If no one expects us to design in an hour, then what is the purpose of System design Interview??

- simulated real life problem solving
- Two co-workers (you & interviewer) collaborate on an ambiguous problem and come up with sol?
- open ended problem
- No perfect answer
- Process and effort > final result.
- Demonstrate your design skill, Defend your design choices, Respond to feedbacks.
- They want to accurately assess your ability.

Myth →

System Design Interview is all about your technical design skills " 

It gives strong signals about your ability to →

→ collaborate → you can collaborate in a very good manner

→ work under pressure

→ resolve ambiguity

→ Ask good questions.

Red flags →

↳ over Engineering is a disease

↳ Too much design purity

↳ Ignoring tradeoffs

↳ Ignoring compounding costs of other engineered systems.

↳ others: Narrow mindedness stubbornness etc.

④ Design News Feed System (step 1 of 4) →

4 step process for effective system Design Interviews →

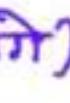
Step ① Understand the problem & Establish Design Scope

→ Ask good question (most imp skill)

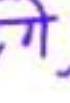
→ Don't jump to answers without clear understanding of requirements.

→ Clarify requirements and assumptions

→ What kind of question to ask ??

↳ What specific features to build? (not specific features )

↳ How many users do the product have?

↳ How fast company anticipates to scale up? (not user )

↳ what is company technology stack?

* News Feed System →

Q → Is this the mobile app? web app? or both?

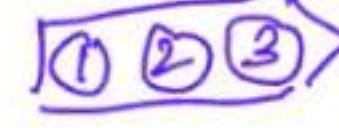
A → Both (Mobile App + Web App)

Q → What are the most important features of the product?

A → Ability to post

→ Ability to see friend's news feed.

Q → News feed sorted in any particular order? first closefriends's feed then group's feed?

A → Reverse chronological order (to keep things simple) 

Q → How many friends can a user have?

A → 5000

Q → What is the traffic volume?

A → 10 million daily active user (DAU)

Q → Can feed contain image/video OR just text?

A → It can contain media file (Both images & videos)

प्राप्त करो above questions, Twitter estimation फिर हमें देखा था

so similar pattern के question you can ask.

⑤ Design News Feed system (Step 2 of 4) →

Step 2 Propose High Level Design and Get Agreement →

→ Make high level design Reach an agreement with interviewer.

→ Blue Print → Ask

→ Box Diagrams with Key Components.

↳ Clients (web/mobile)

↳ API

↳ Web servers

↳ DataBase Cache

↳ CDN

↳ Message Queue etc.

→ Ask if Back of the envelope is required.

→ Ask if Back of the envelope is required.

↳ If yes → Then do the calculations to evaluate if your blueprint fits the scale constraints.

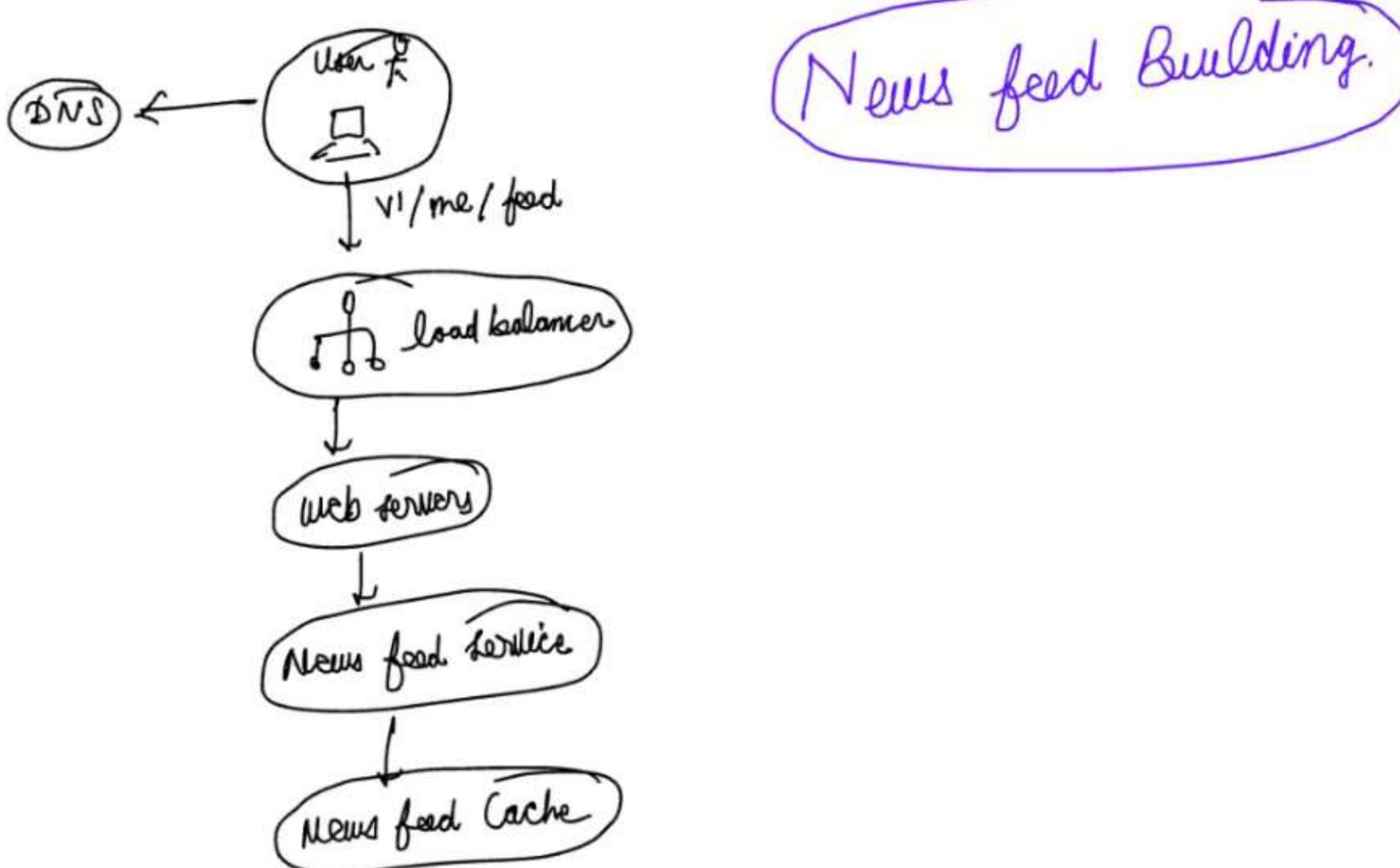
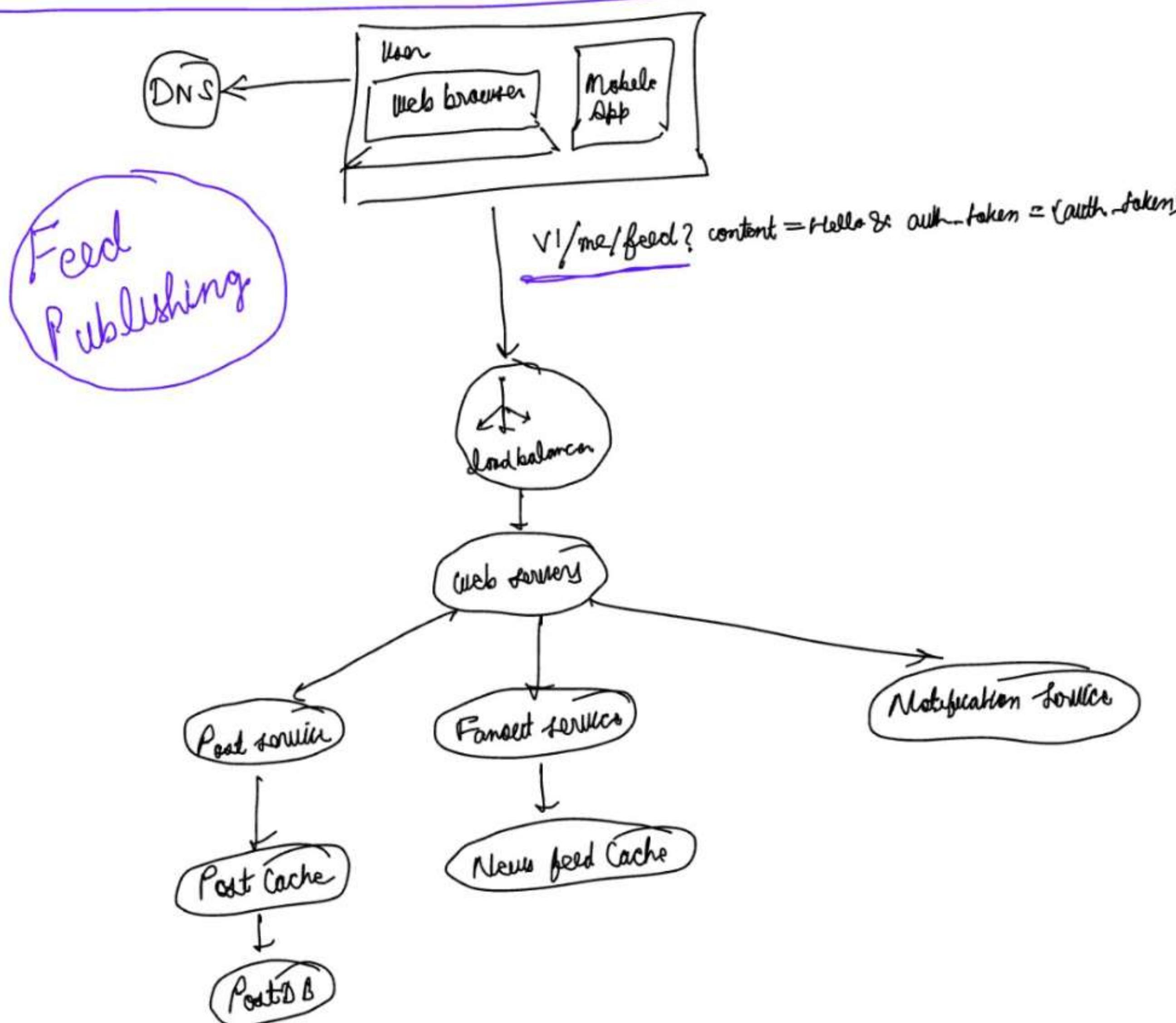
* News Feed System (High Level Design) →

Note → Few Components in design, I will cover later bcz they are long topic and need to be cleared properly in separate lecture.

High Level 'News Feed System' will have 2 flows

↳ Feed Publishing → Published post → DB/cache
→ Friend's News feed.

↳ News Feed Building : Aggregating friend's post in "reverse chronological order".



⑥ Design News Feed System : (step 3 of 4) →

(step 3) Design Deep Dive →

By this time following should have already been achieved.

↳ overall goals and scope - Agreed.

↳ High level blue print.

↳ Feedback on High level design.

↳ what area to focus on - Based on Interview feedback

→ Identify & prioritize components → Interview items

→ May focus on → High level design (HLD)

→ Back of the envelope estimation

→ Any particular component

e.g. ④ URL shortener - Hashing is important

⑥ Chat system - Reduce latency support (online/offline) states etc.

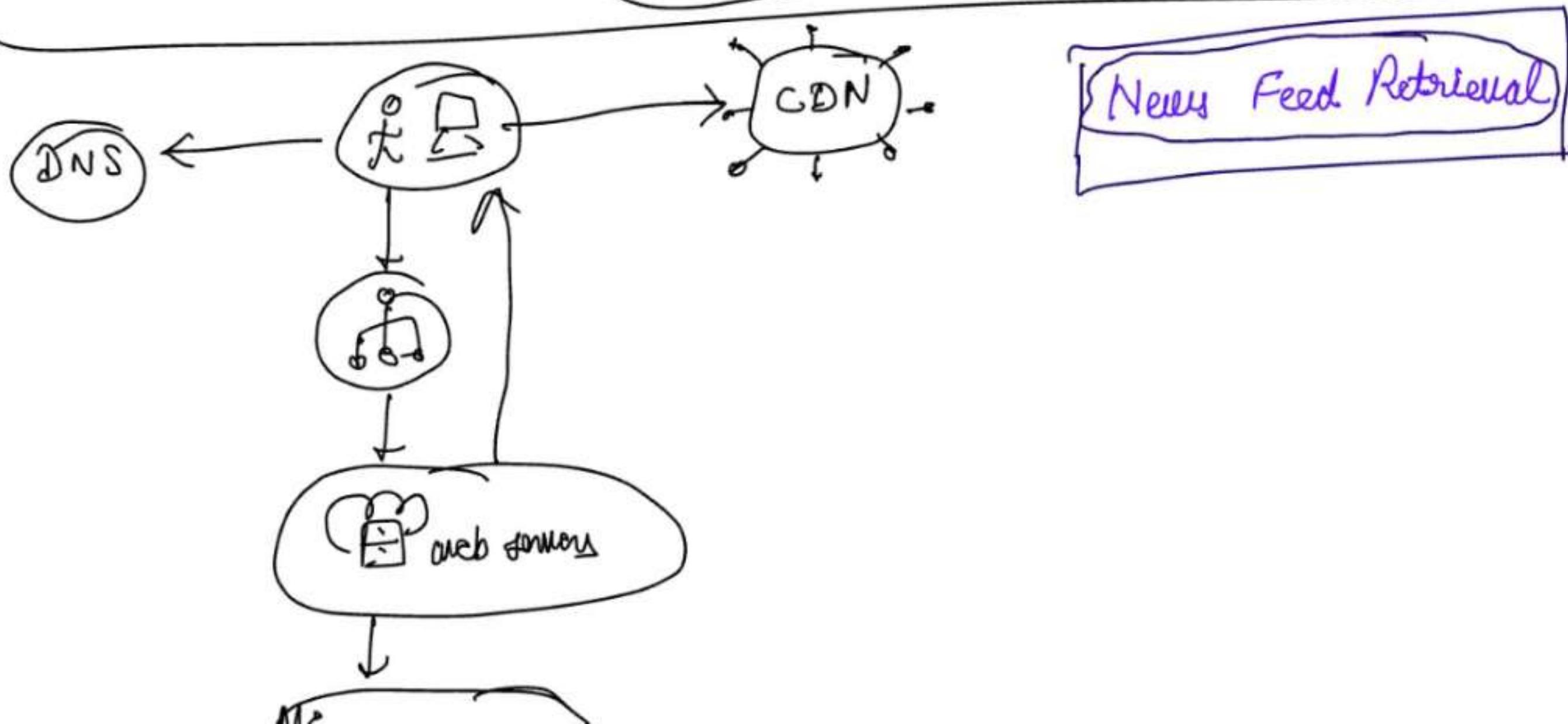
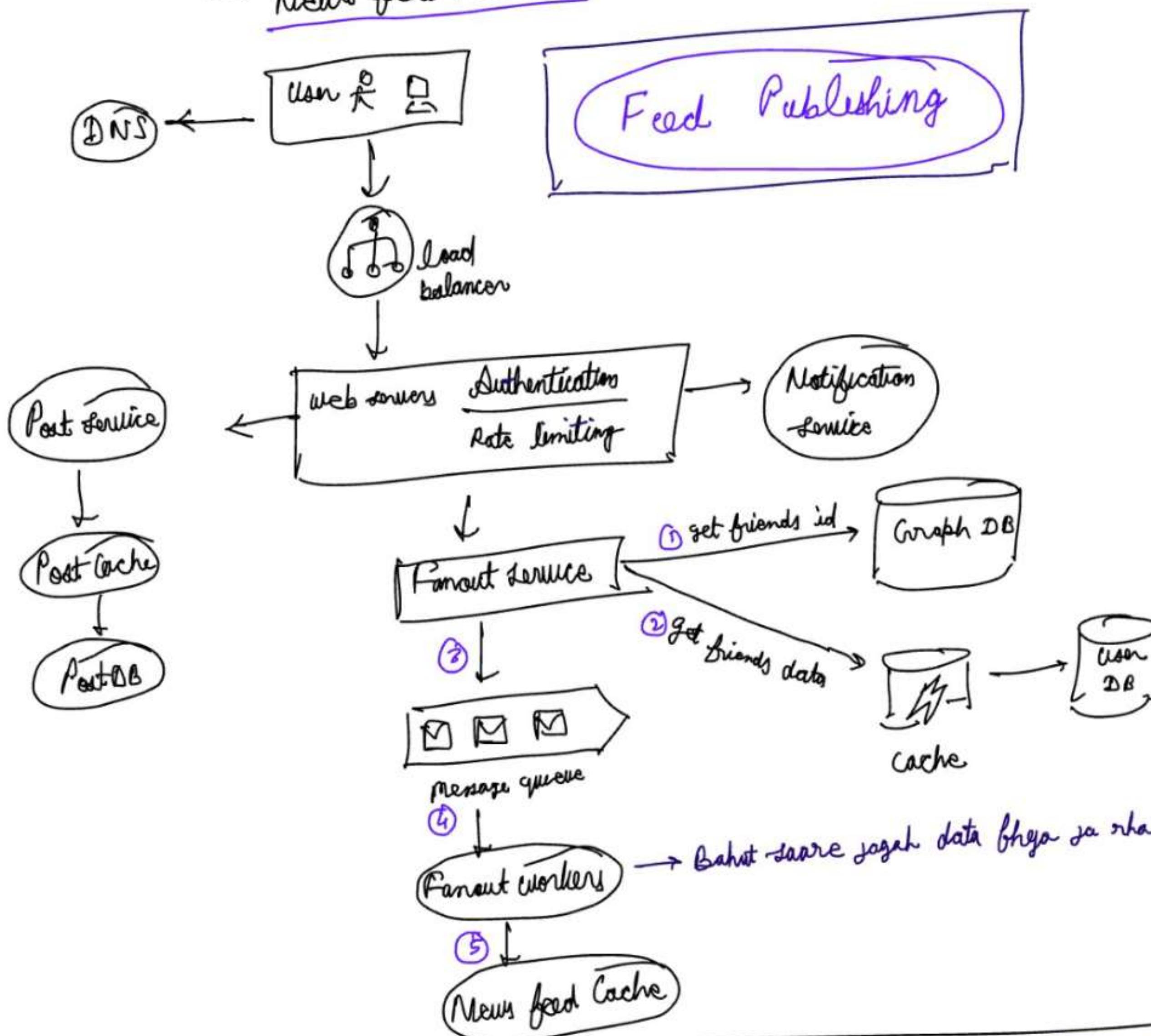
→ Don't waste time in un-necessary component details

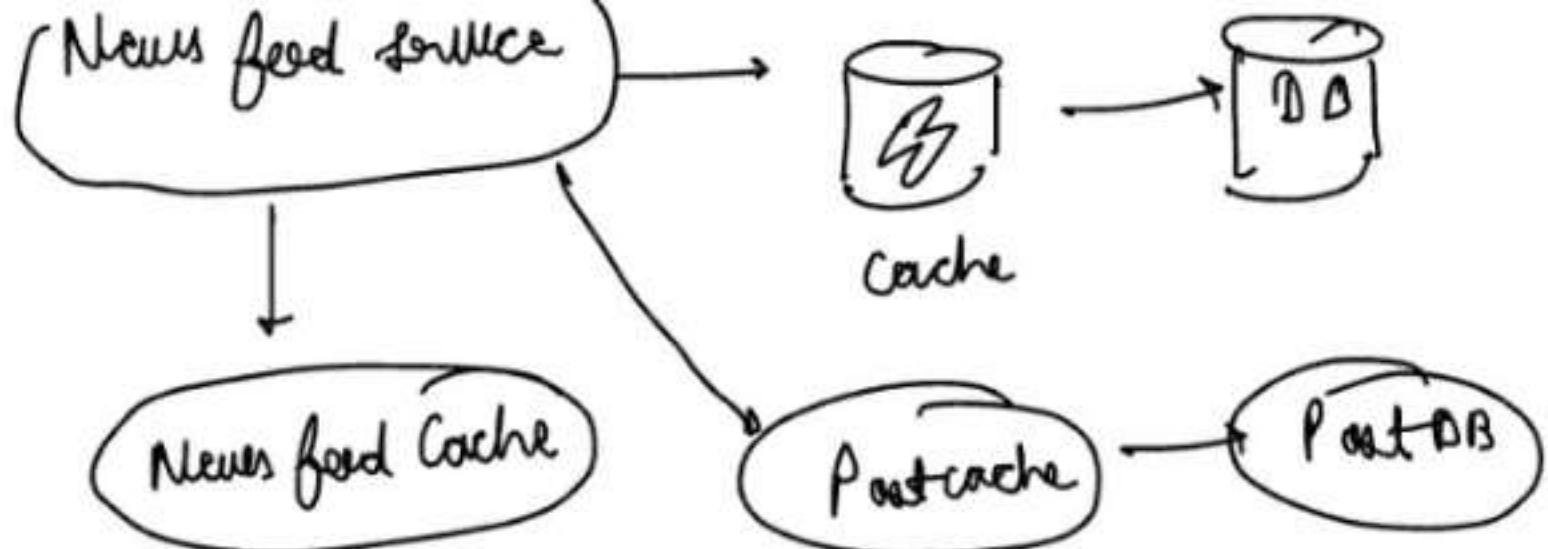
* News feed system (Design Deep Dive) →

2 imp use case

↳ Feed Publishing

↳ News feed Retrieval





⑦ Design News Feed System (step 4 of 4) →

Step 4 Wrapup.

→ Interviewer might ask Follow-Up question.

- ↳ Improvements / system bottlenecks
 - ↳ Never say your design is perfect.
- ↳ Recap of your design.
- ↳ Error Cases (Network loss / server failure) etc.
- ↳ How do you monitor metrics / logs?
- ↳ Next Scale Curve . 1 million to 10 millions users

Dos

- Ask for clarification
- Understand requirements
- Communicate - talk out your thoughts
- Multiple approaches (if possible) (else not)
- Agreement → Detail each component → most critical one first.

Dont

- Don't basic Q's don't answer
- Don't directly jump to solⁿ
- Too much detail on un-necessary
- Not asking for hints.
- Not asking feedbacks early and often.

* Time Allocation on Each step →

- ① Understand problem + Establish Design scope (3-10 min)
- ② Propose high level design + Agreement (10-15 min)
- ③ Design Deep Dive (10-25)
- ④ Wrap (3-5 min)

⑧ Design Rate Limiter →

In Network Systems →

Rate Limiter → Used to control the rate of traffic sent by client or service.

In HTTP world →

Rate Limiter → Limits the no. of client requests allowed to be sent over a specified period.

```
if count(API request) > Threshold  
{  
    "Access Call Block"  
}
```

Examples →

→ No more than 2 post / second

→ Maximum 10 accounts / day from the same IP address

Before designing, see the benefits of Rate Limiter

→ Prevent DoS Attack →

→ Every org. uses rate limiting. (But API call API)

→ Example → Twitter → no. of tweets = (300 per 3 hours).

→ Blocking excess intentional or unintentional requests.

→ Reduce Cost →

→ Limiting excess request = (Few Services) + (More resources for high priority APIs)

→ Important if we use 3rd party APIs → ऐसे लगाते हैं।

→ Prevent Worker Overloading →

→ Filter out excess requests

→ Caused by Bot or user misbehaviour

⑨ Design Rate Limiter (Step 1 of 4) →

Step ① Understand the problem and Establish Design Scope →

Note →

Rate Limiter

Interact with interviewer.

Get to know what kind to implement

Different Algorithms

Pros & Cons

Q → What kind of Rate Limiter?

↳ client side rate limiter

↳ server side rate limiter → mostly पैर उतारते हैं

A → Server side API rate limiter

Q → Limit API requests based on

↳ IP? (same IP)

↳ userID? (same user ID)

↳ other properties?

A → Flexible enough to support different sets of limiting rules.

Q → What is scale of system?

↳ start up

↳ big company (large user base)

A → Able to handle large no of requests.

Q → Will system work in distributed systems?

A → Yes → बड़ी company distributed environment में काम करते हैं।

A → Yes → Should it implement in application code?

Q → Rate limiter separate service होगी ??? पर क्या यह should it implement in application code?

A → It's a design decision, (up to you).

Q → Do we need to inform users who are throttled (limited)? (जिनके लिए limit किया है उनको पता देना चाहिए)

A → YES

* Requirements →

→ Limit excess requests

→ Large User Base

↳ should not slow down HTTP response time. (slow nahi hona chahiye)

↳ use as little memory as possible. (Rate limiter के लिए कम से सोशी उतारना चाहिए)

↳ High fault tolerance. (Back up jaisa such hona chahiye)

↳ If any issue (Cache server goes offline, server down etc).

↳ Must not affect entire system.

→ Exception Handling → show clear exception to user (user को inform करना है कि उसकी request block हो गई)

→ Distributed rate limiting → shared across multiple servers. (shareable across servers)

⑩ Design Rate Limiter (Step 2/4) → Part 1 → where to put Rate Limiter →

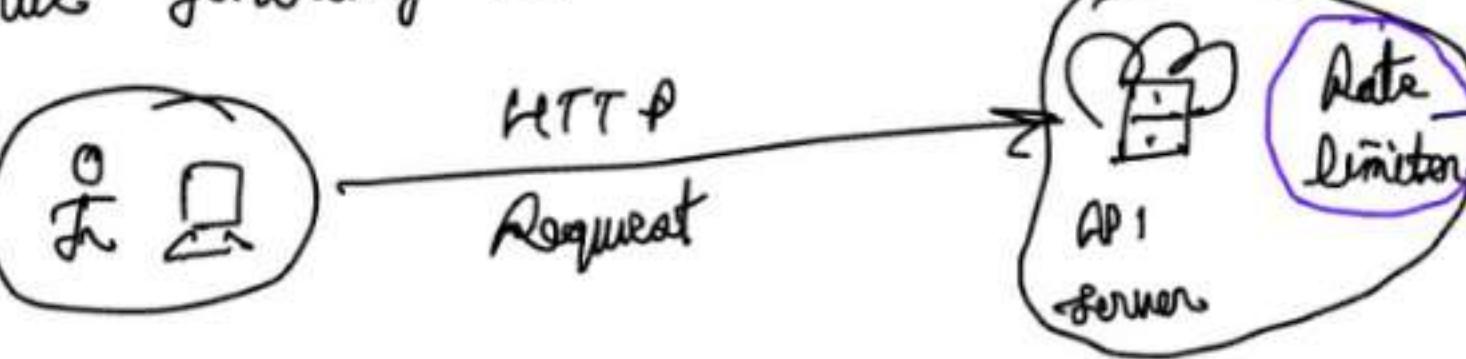
Step ② Propose high level design of Great Guy in →

Note → Where to put Rate Limiter?

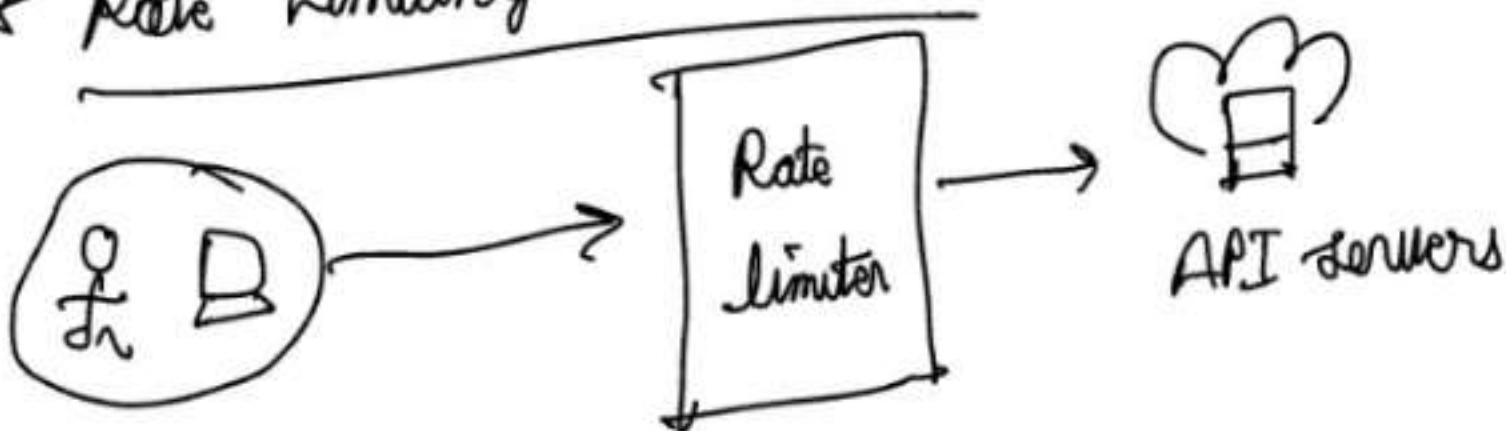
→ Rate limiter algorithm

→ High level Architecture

Where to put Rate Limiter?

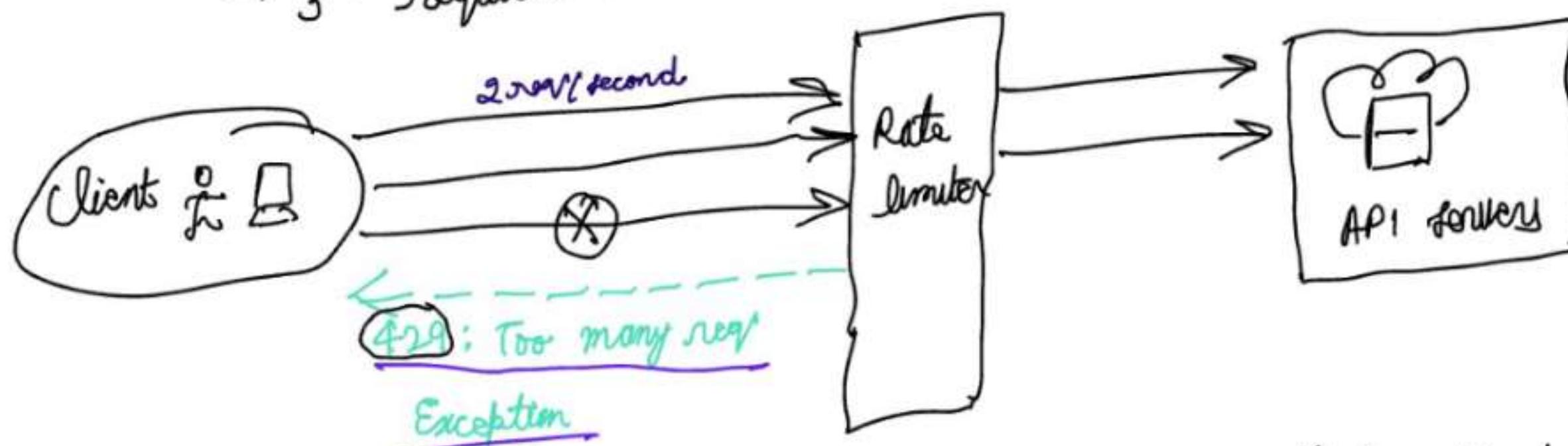
- ↳ client side → we generally have no control on client side. (client side mein nahi rkha jao)
- ↳ server-side → 

* Rate Limiting Middleware →



Example →

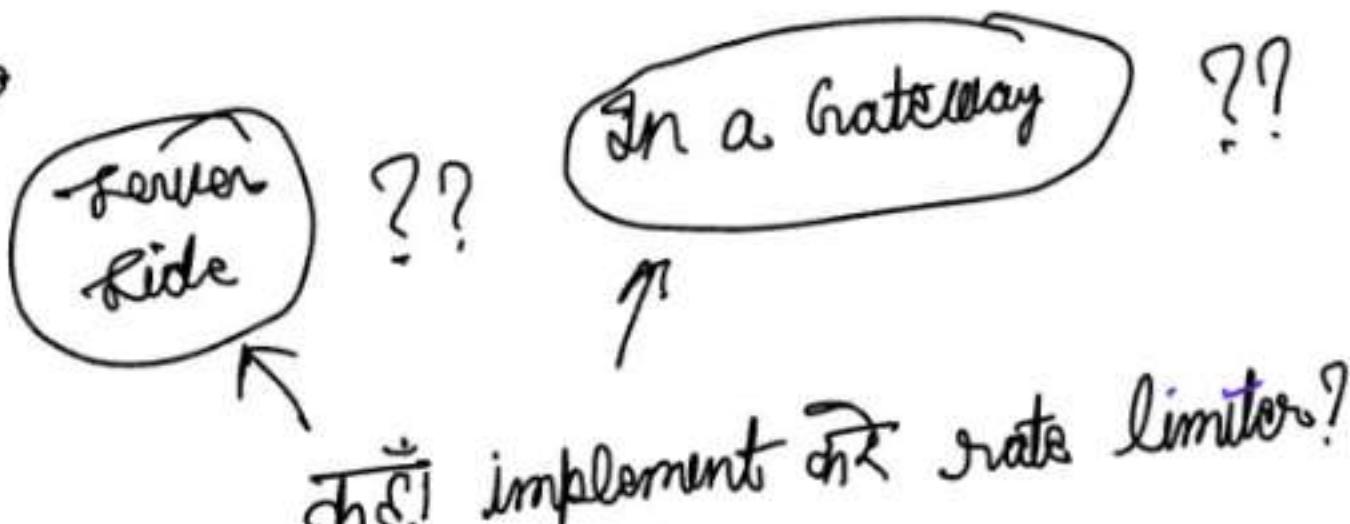
- Assume → API allows 2 requests / second.
- Client sends 3 requests within a second
 - 3rd request throttled with error: HTTP status code 429.



* API Gateway (Middleware) → Fully managed service that supports →

- ↳ Rate limiting
- ↳ IP whitelisting
- ↳ Authentication
- ↳ static Content etc.

So,



Chill!!! It depends Bro →

- ↳ current tech stack in your service.
- ↳ Programming language is efficient?
- ↳ Rate limiting algorithm that suits your business. (server side / or / gateway side?)
- ↳ Already API gateway use ~~SSL~~ ~~HTTP~~ then implement in it.
- ↳ Use Commercial API Gateway to save time & resource?

11 Design Rate limiter (2/4) part 2 → Token Bucket Algorithm →

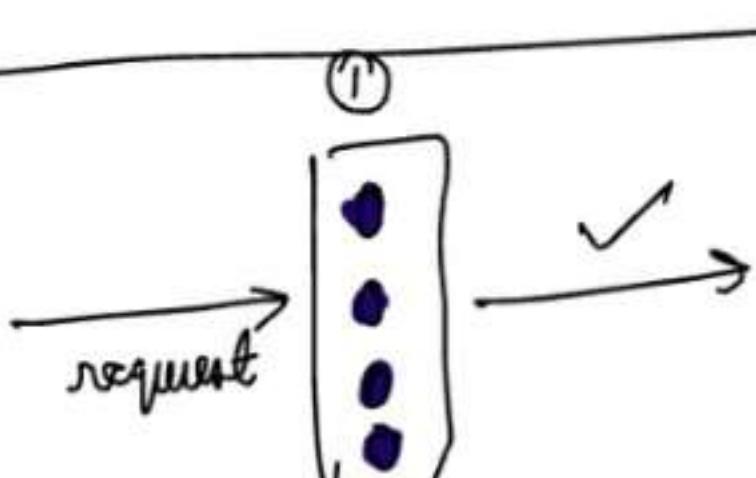
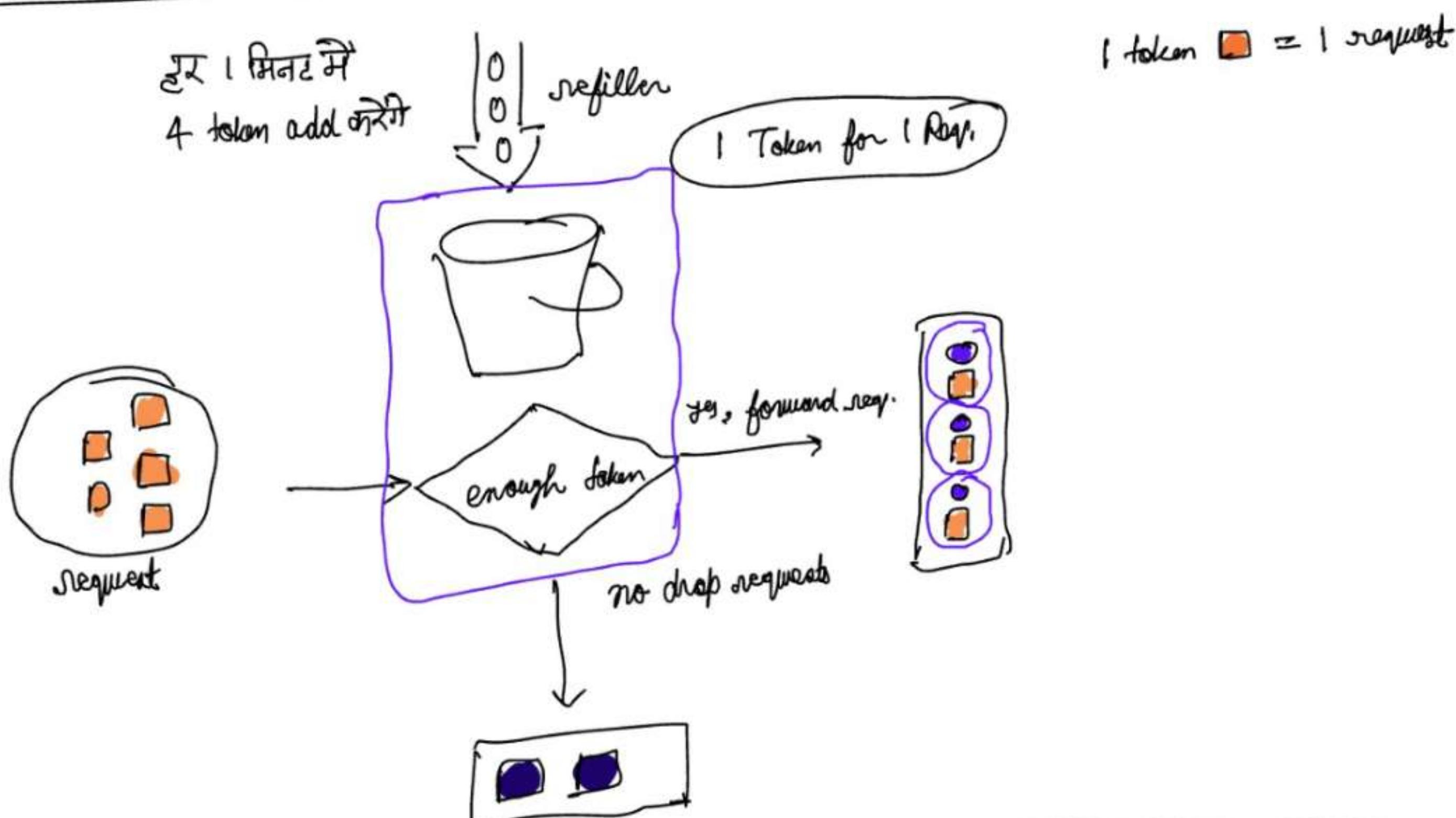
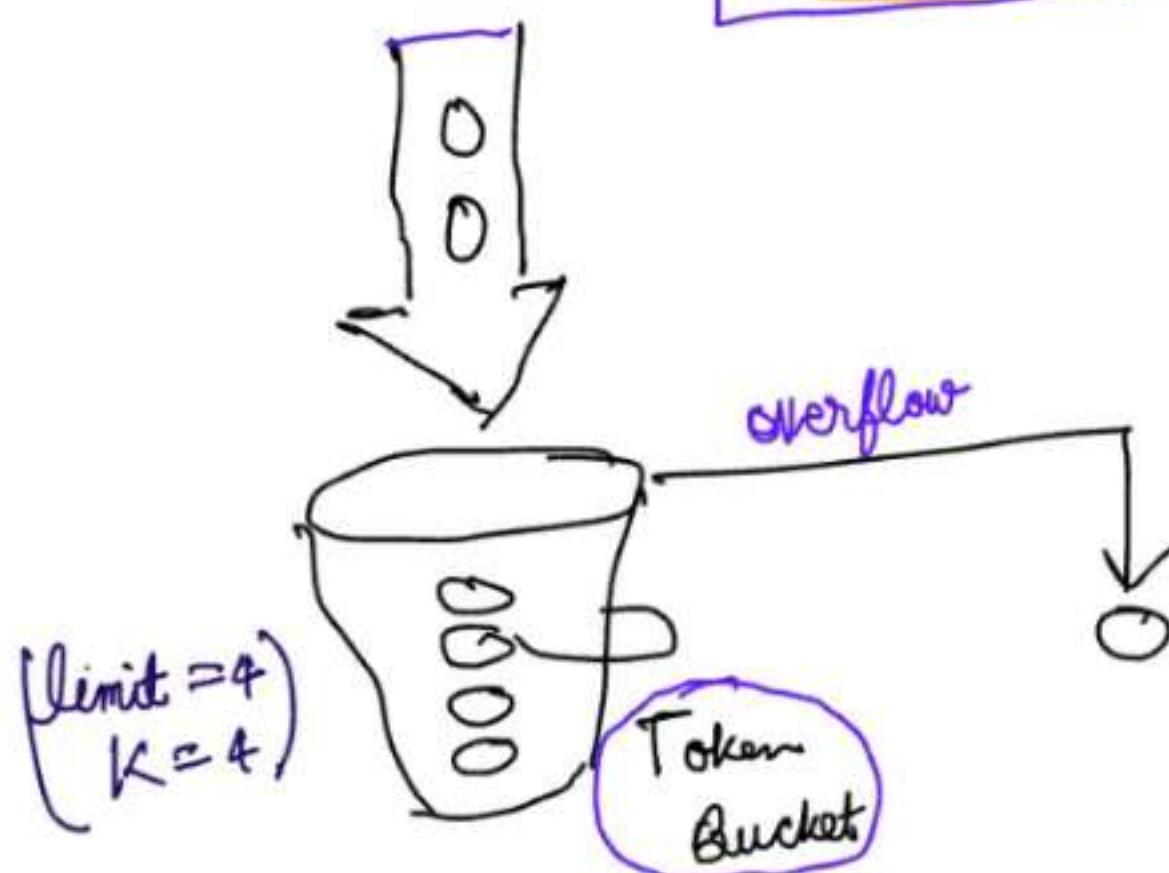
step 2 : part 2 → Proposed HLD & get Buy in →

Note →

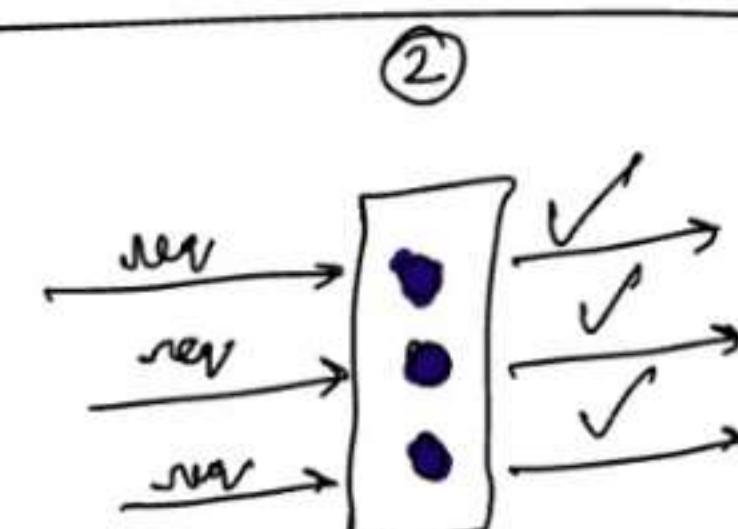
- where to put Rate limiter?
- Rate limiting Algorithms?
- High level Architecture?

① Token Bucket Algorithm →

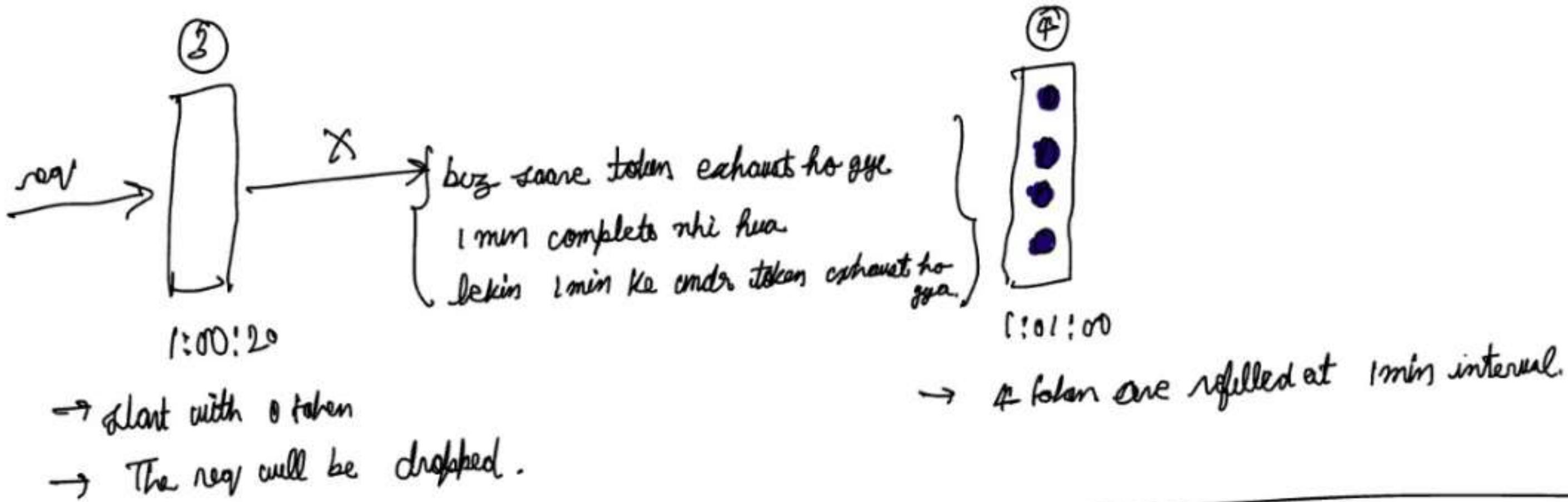
- Widely used & simple
- Used by companies like Amazon, Stripe etc.
- Takes two params
 - "bucket size"
 - "Refill rate"



- start with 4 token
- The request will go through
- 1 token is consumed



- start with 3 token
- All three requests will go through
- 3 tokens are consumed



* How many Buckets?

→ Depends on rate-limiting rules

Example → Different API ends points : Different Buckets
If you have 3 APIs, use 3 Buckets.

→ Throttle Based on IP addresses :- 1 Bucket / IP address.

→ System allows max 10000 reqs / second → A 'Global Bucket' shared by all requests.

Pros :-

→ Simple

→ A req. can go as long as we have tokens.

Cons :-

→ Bucket size and Token Full rate timing is challenging → Take 2 parameters (Bucket size, Token Full)

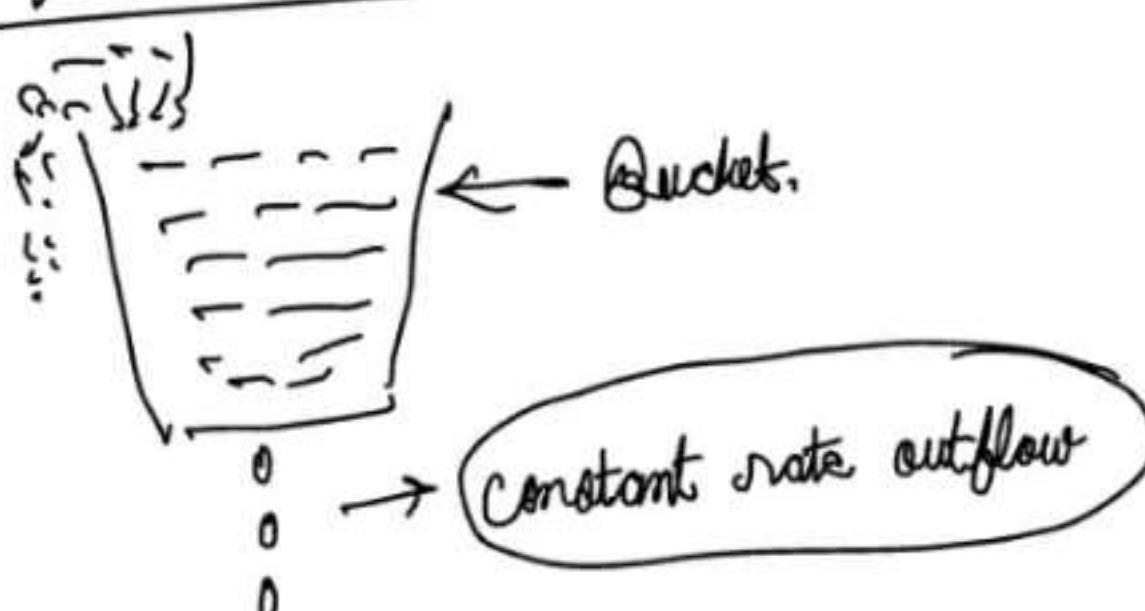
→ Bucket size and Token Full rate timing is challenging → Take 2 parameters (Bucket size, Token Full)

→ sends packet at an average rate. (Rate और भारतीय दूरी)

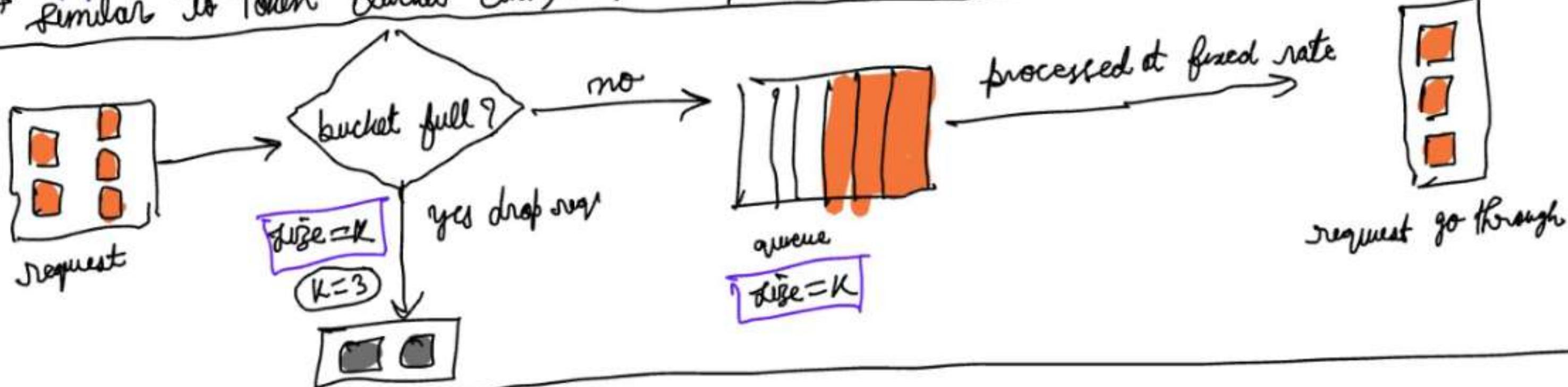
12 Design Rate Limiter (2/4) Part(3) →

② Leaking Bucket Algorithm →

overflows when full



→ similar to Token Bucket but, requests processed at fixed rate.



→ Using FIFO queue (First in First out)

Takes 2 parameter →

- Bucket size → Equal to Queue size
- Outflow rate → No of request at fixed rate.

* Shopify (e-commerce) → uses leaky bucket for rate limiting

Pros :-

- Memory efficient . Limited Queue size.
- Stable outflow rate.

Cons :-

- Takes 2 parameter . Tuning is challenging. (Bucket size , outflow rate)

⑬ Design Rate Limiter (2/4) part 4 → Fixed Window Counter Algorithm →

③ Fixed Window Counter Algorithm →

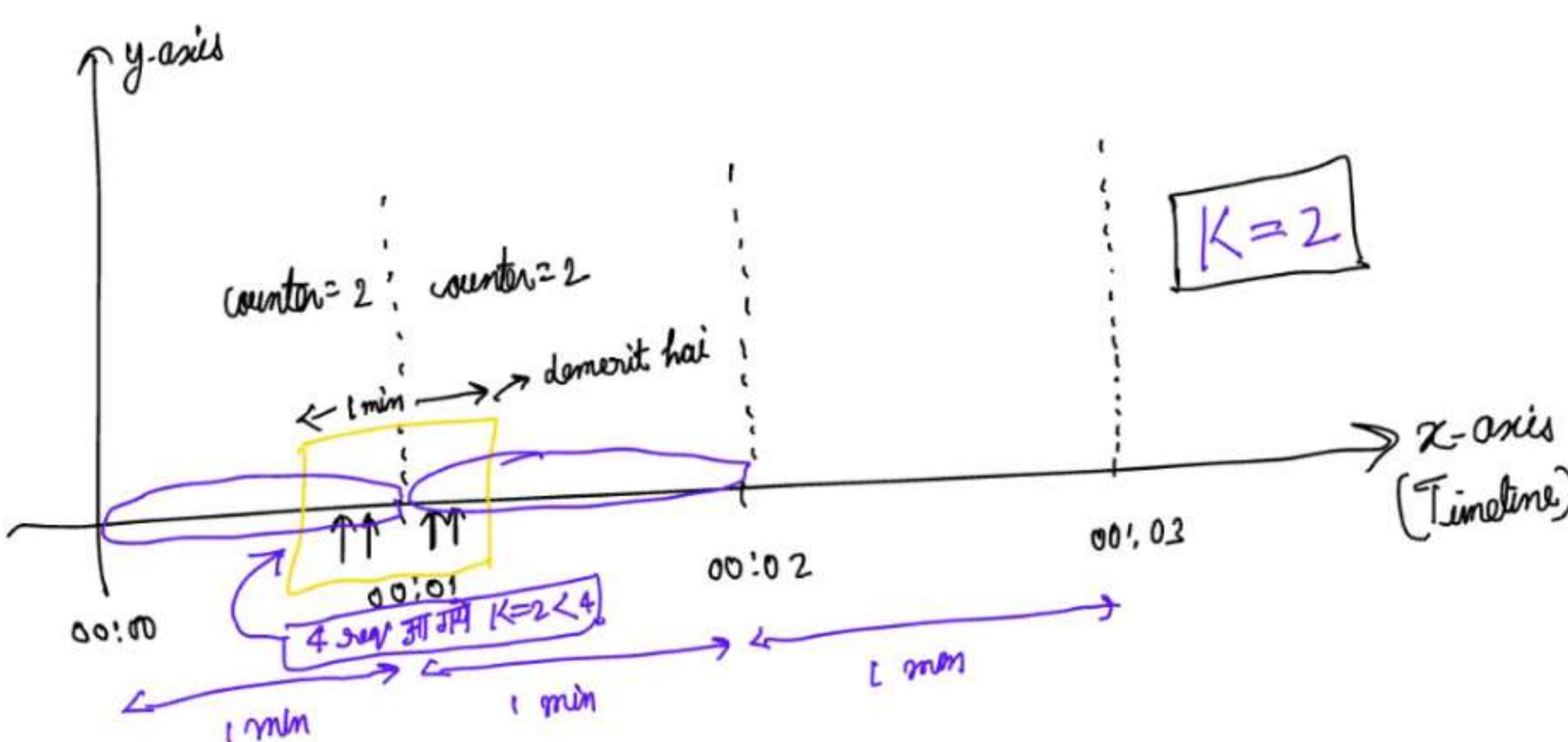
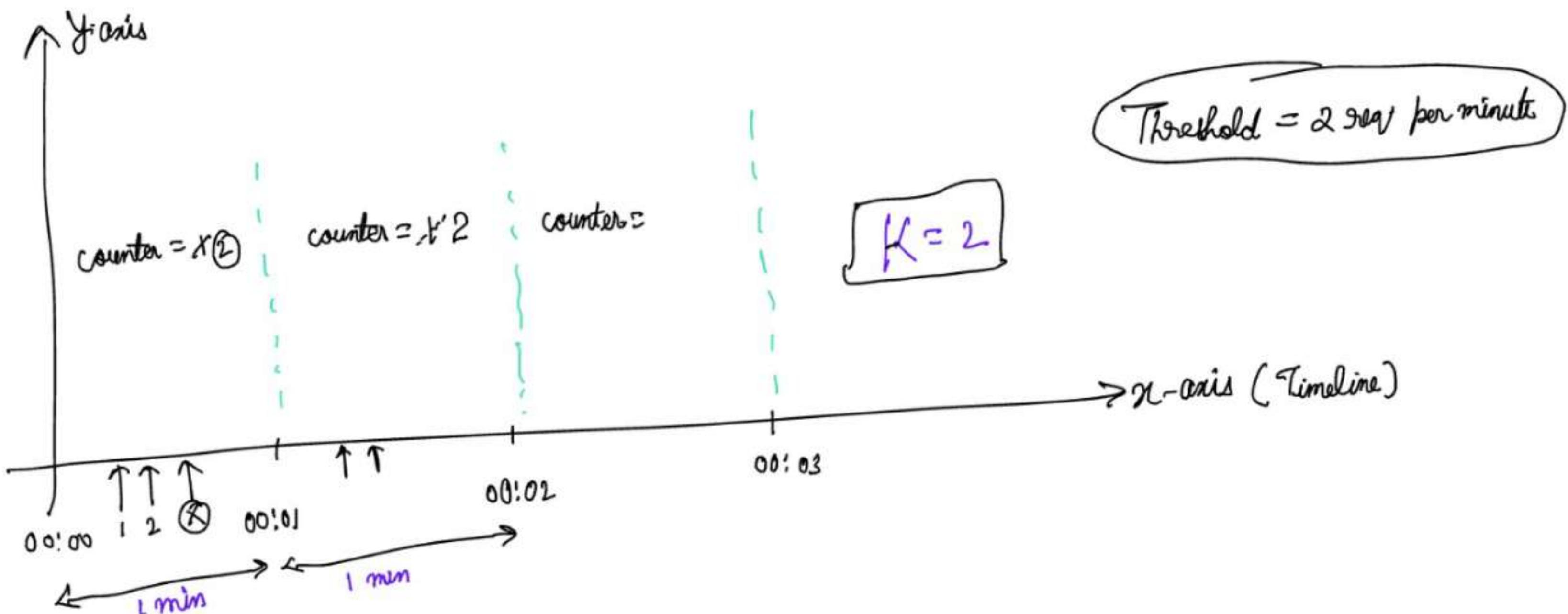
→ Timeline - Fixed time windows (Timeline की fixed windows में गाँठ देते हैं)

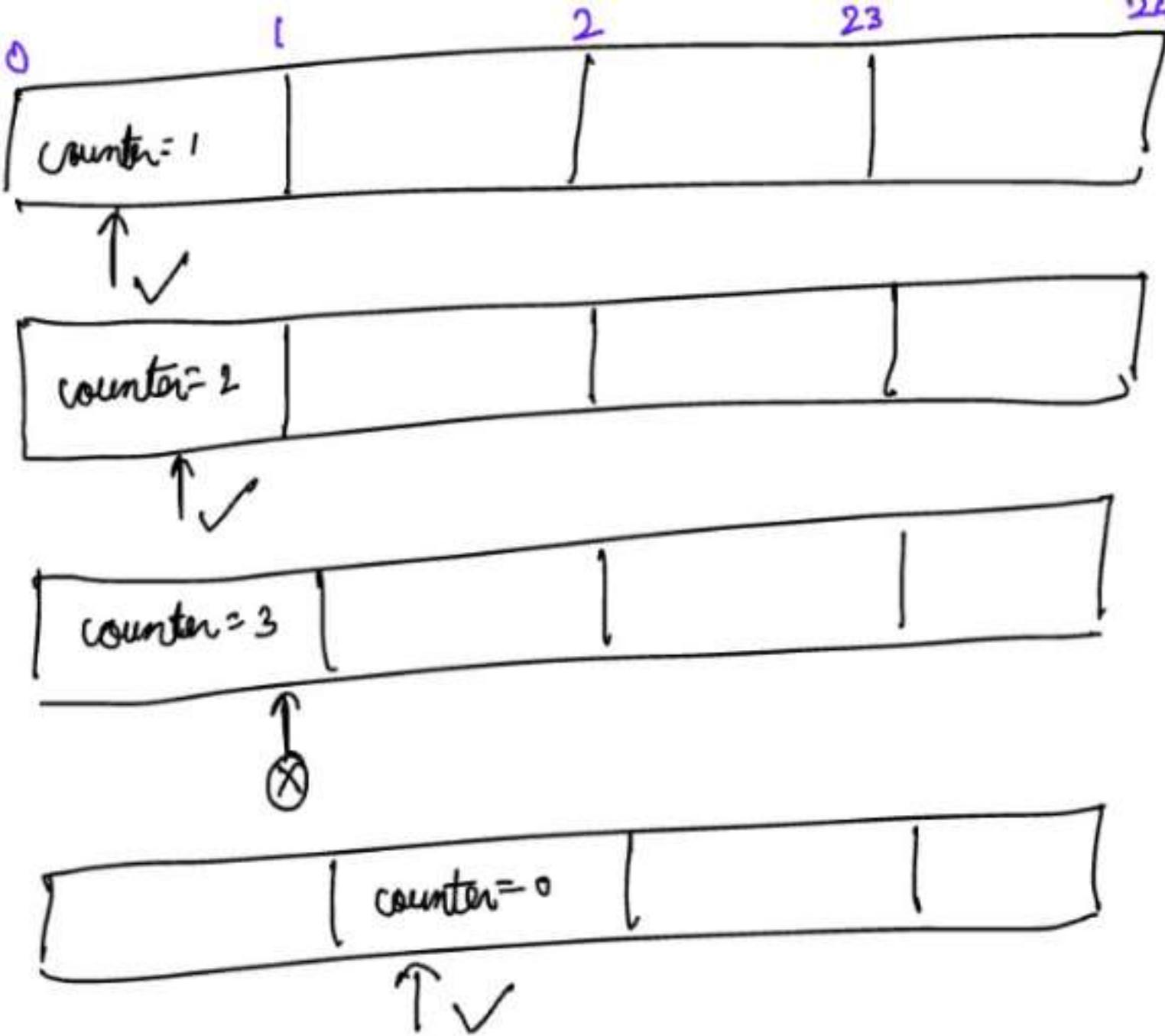
→ Each window has a counter

→ Counter increment +1 per request.

→ Counter reaches threshold - new requests dropped - Until new window starts.

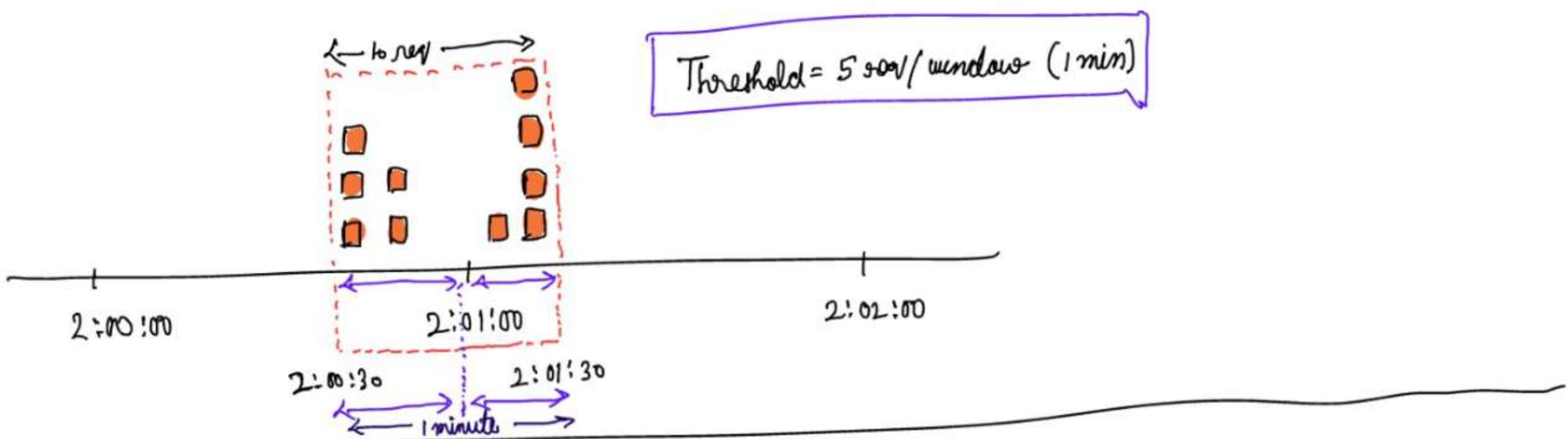
→ Counter reaches threshold - new requests dropped - Until new window starts.





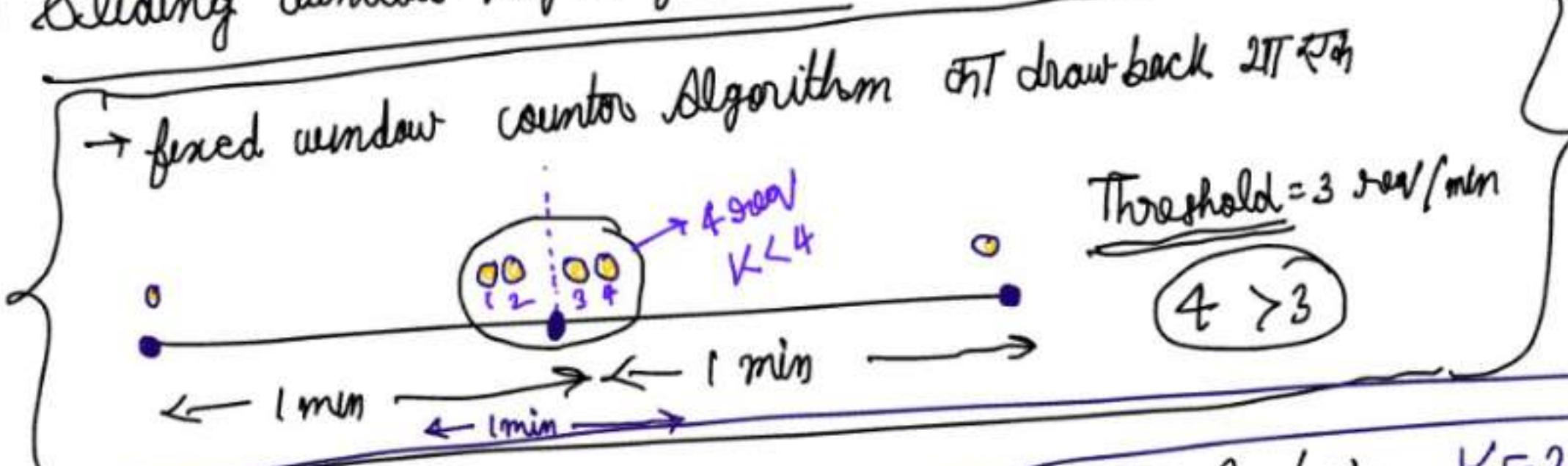
Pros :- Easy to understand & implement.

Cons :- Traffic Burst at edges of a window cause more request than the threshold.

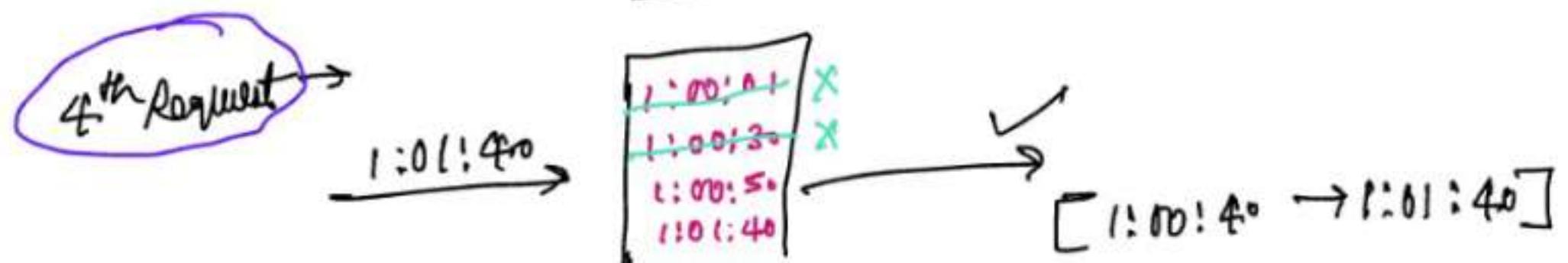
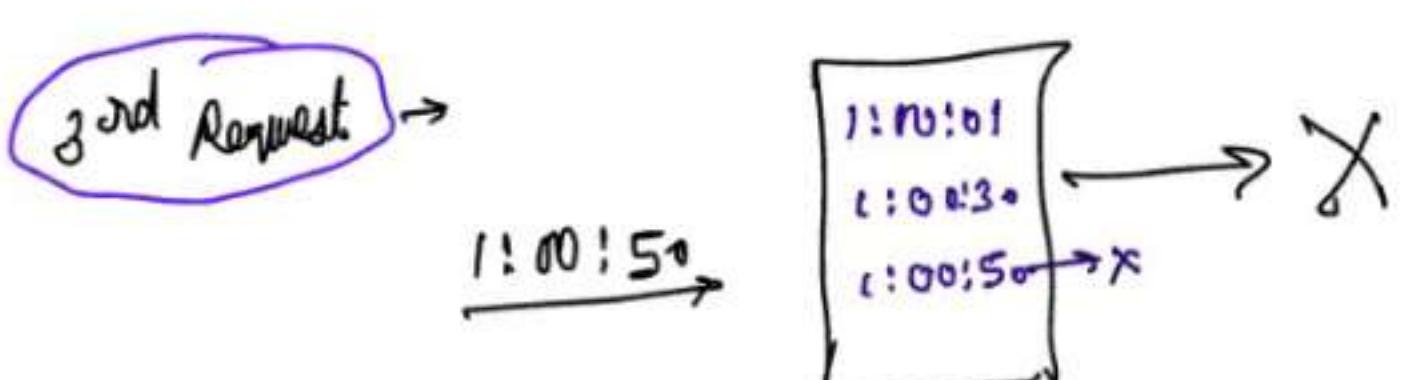
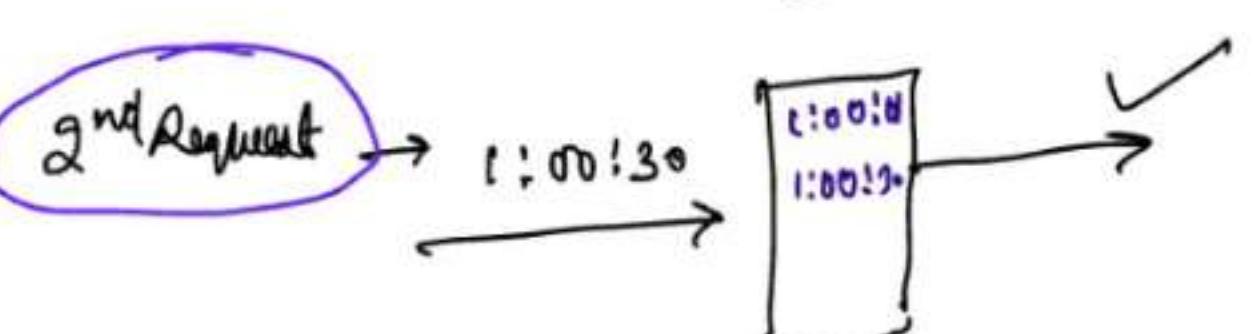
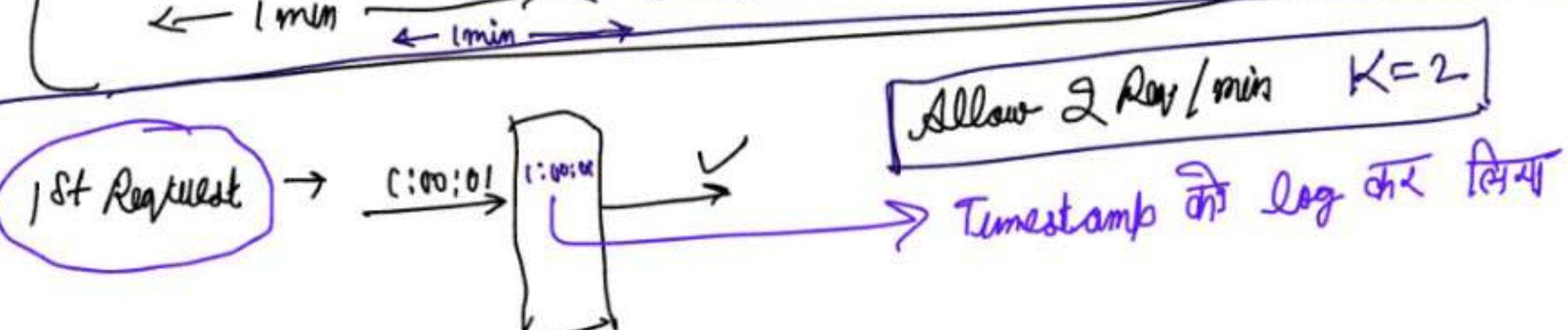


⑭ Design Rate Limiter (2/4) Part II →

④ Sliding Window Log Algorithm →



→ To fix this we use sliding window log algorithm.



- log every req's timestamp (हर request का timestamp सोके रख लेंगे)
- New Request — remove outdated timestamps (पुराना log हटाएंगे)
 - New Request का timestamp add कर दी
- Timestamps का count \geq threshold \rightarrow Rejected

Pros →

- Accurate
- Will not exceed threshold

Cons →

- Even if Req. Rejected, Logged in.
- Cache. (Memory usage high)  (Memory का उपयोग)

(15) Design Rate Limiter (2/4) part ⑥ →

High level Architecture →

- Basic Idea → High level में कैसे होते हैं
- Counter → To keep track of no of requests.
 - Same user, same IP etc.
 - counter > Threshold : Discard Request.

Interviewer → Where shall we store counters?

Candidate → रजिस्ट्रेशन से जावा बदल दिये तभी बदल दिया

Interviewer: Where to store counter?

Candidate: Using a database

Interviewer: Don't you think it will be slow due to slowness of disk access.

Candidate: In-Memory Cache

↳ Fast I/O

↳ Time based expiration strategy support.

Interviewer: Can you give an example?

Candidate: Redis Cache is a popular option.

Interviewer: Can you tell me more about redis?

Candidate: → In-Memory fast store

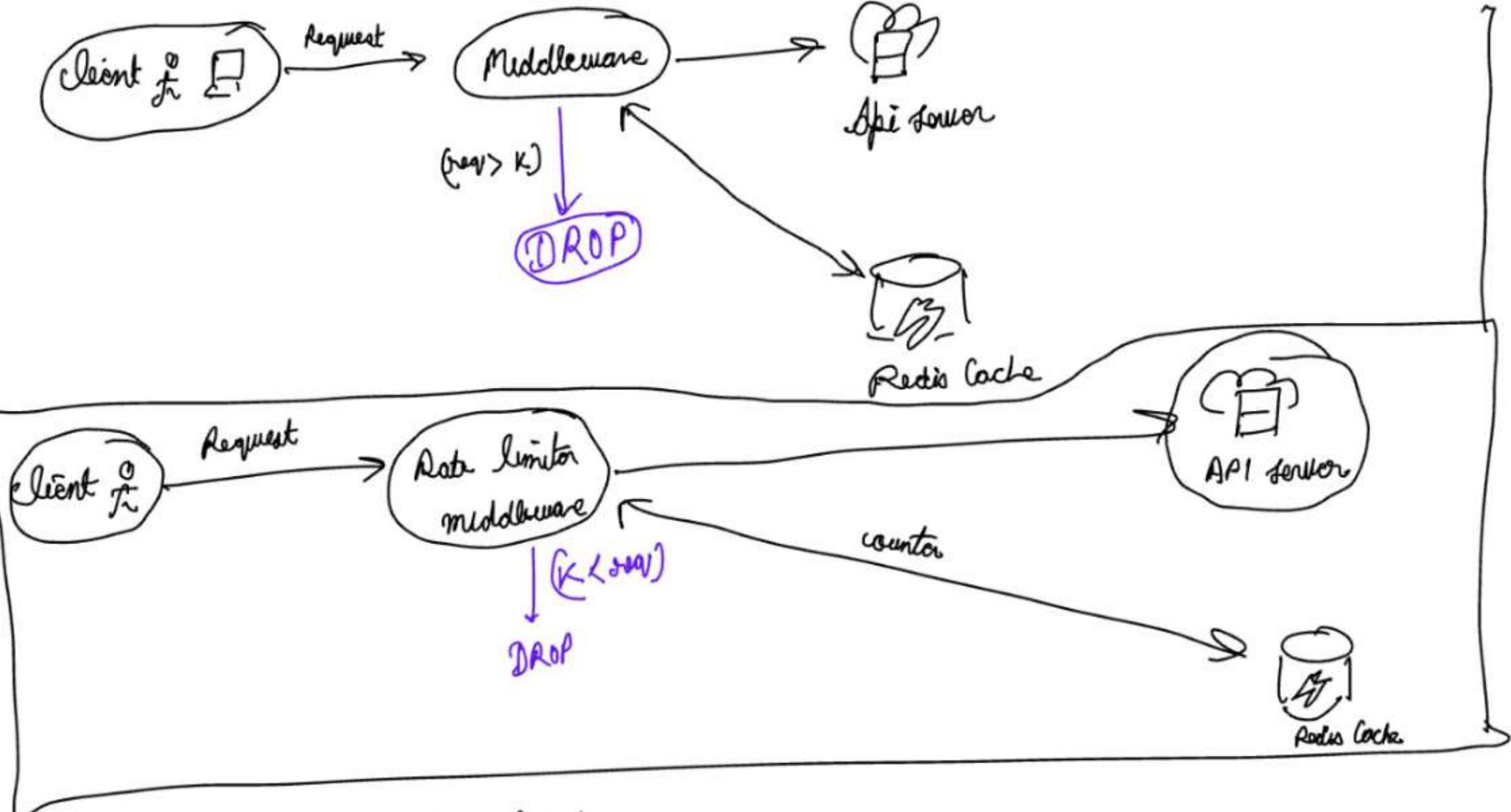
→ offers 2 commands

↳ INCR : counter++

↳ EXPIRE : Timeout for counter

High Level Diagram →

Rate Limiter



$\text{if } (\text{counter} > \text{Threshold}) : \text{Reject}$

⑯ Design Rate Limiter (3/4) Part ① →

part ① → Design Data Structure →

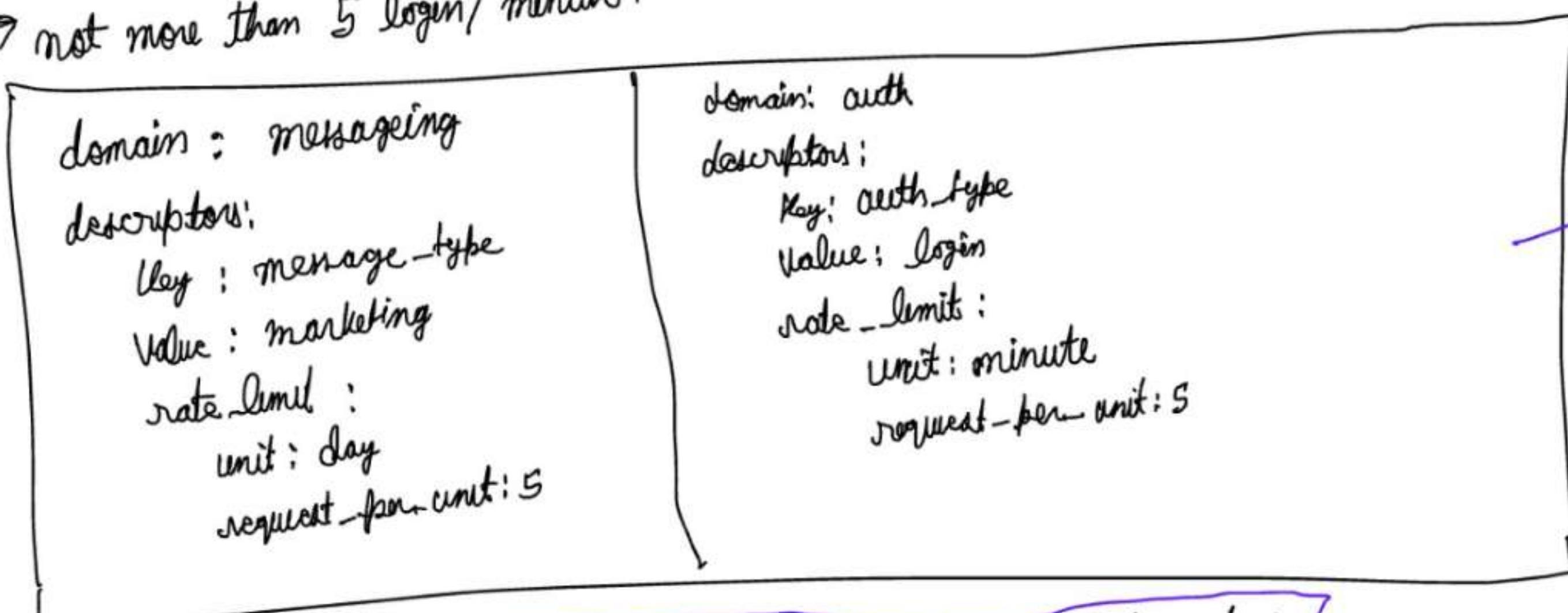
Interviewer → How are rate limiting rules created?
कैसे स्टोर करेंगे रुल्स की?
जो request limit तक तकाला करेंगे?

Candidate → First, lets talk about limiting rules

Rate Limiting Rules →

Example :-  (UBER, OLA)

- max 5 marketing message / day
- not more than 5 login/ minute.



Rules →
#

Rules are written in configuration files and saved on the disk.

Candidate: Coming to the next question

* Exceeding The Rate Limit →

जब कोई request limit हो जाए, then what.....

→ send HTTP response code **429** to client

→ **429**: Too many requests

→ Depends on use case what to do next:-

↳ Enqueue → processed later (Amazon की order queue हो जाए)

Example → Some of your orders got rate limited.

Interviewer:- How will the client know?

↳ Which Request rate limit हो गया?

↳ Which पहले कितने Request हो जाए before being rate limited?

Candidate → Via "HTTP Response Headers" sent by rate limiter.

Headers	Meaning
X - RateLimit Remaining	Ki तक Req लें सकते हैं।
X - RateLimit - Limit	Max kitne req you can send
X - RateLimit - Retry - After	Kitne time बाद you can send req.

So if the client has sent too many requests,

{ → 429 (Too many Req)
→ X - RateLimit - Retry after

will be sent to client

① Step ③ Design Deep Dive (Part 2)

Detailed Design →

→ Rules को **disk** में रखना है (Remember ???)

→ Disk slow हो तो **cache** भील में रख लें।

→ Req तो आता रहेगा - Use **workers** → time to time disk se Rate limiting rules nikal ke cache me daalte raha time to time

→ Client होगा

→ Server होगा (जिसको client req भेजेगा)

→ भील में **Rate limiter**

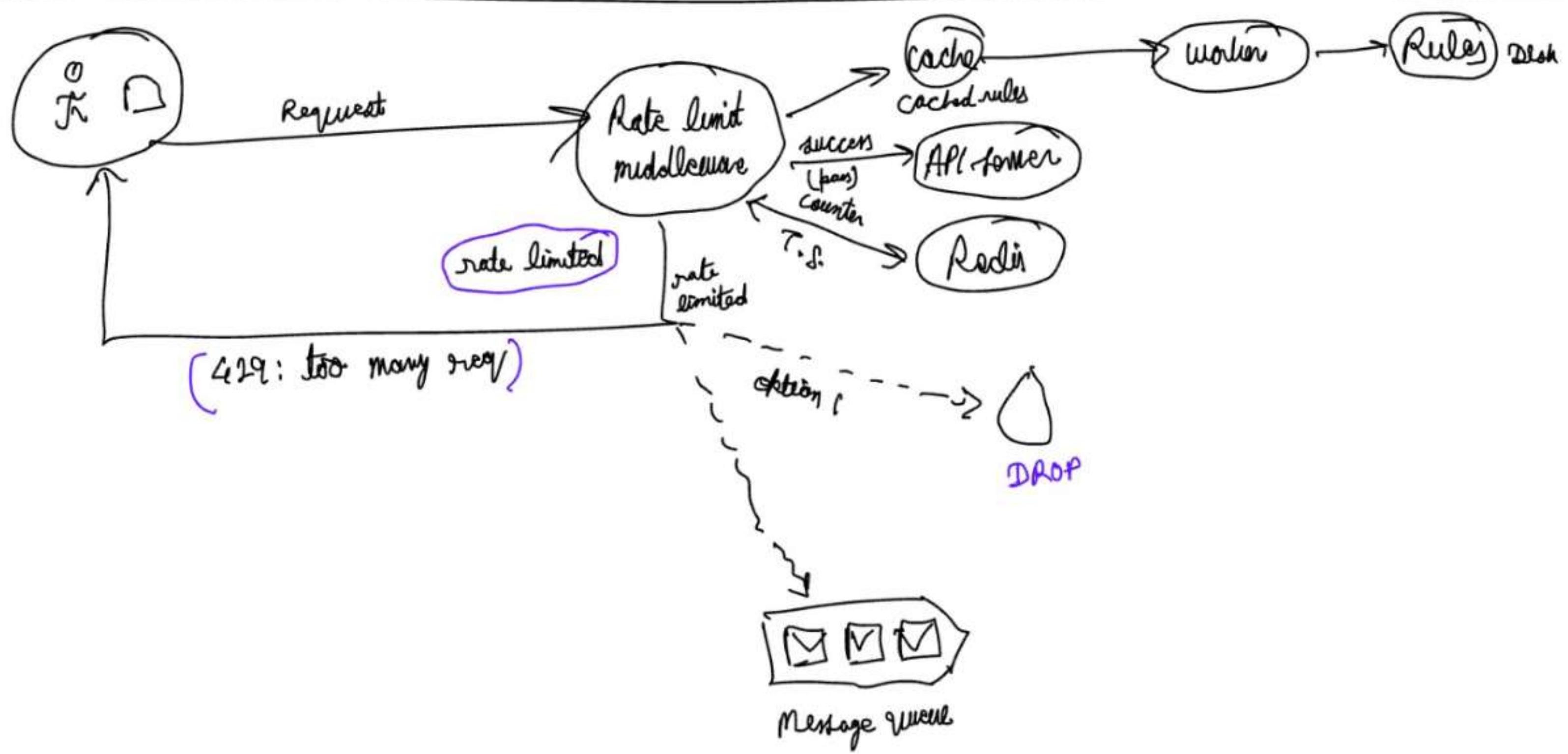
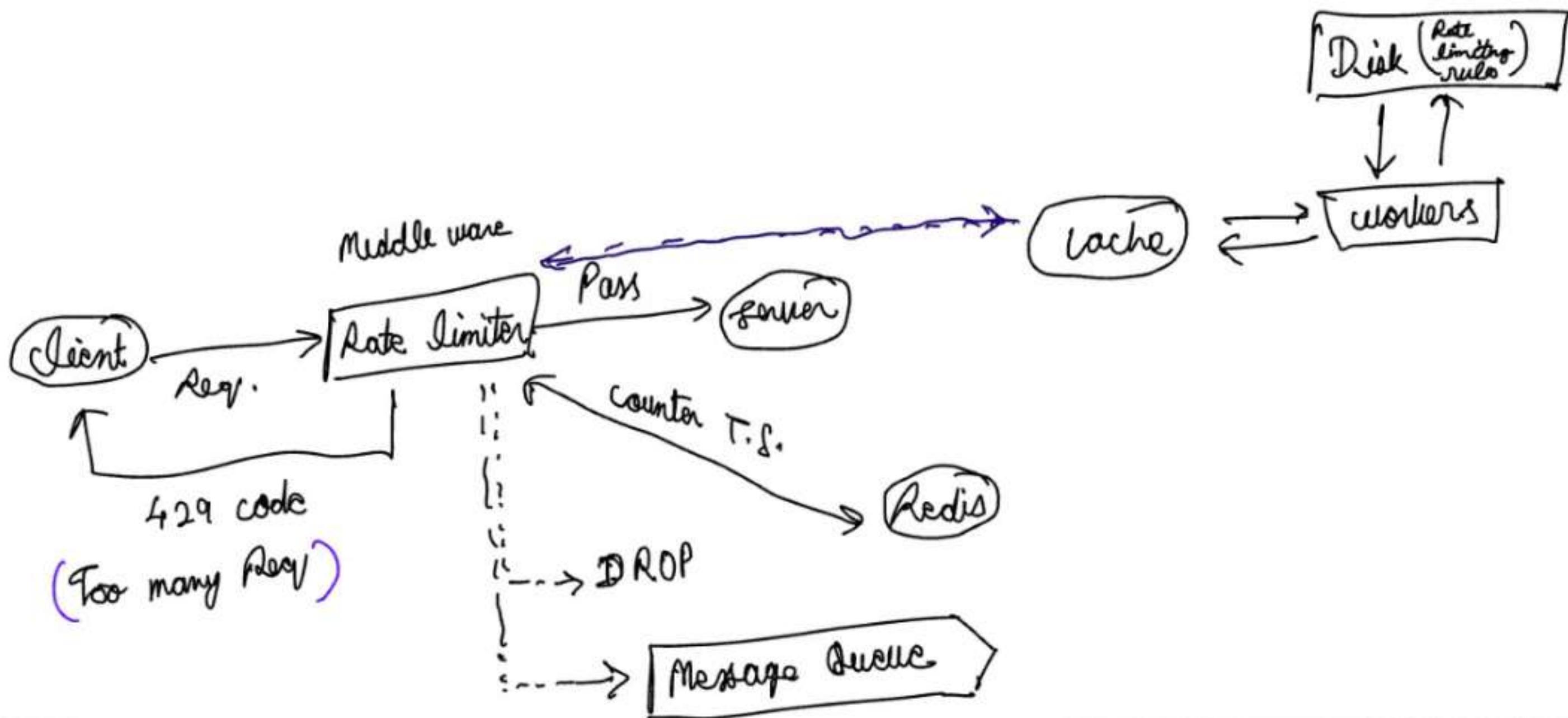
→ याद करो Rate limiter, counter credit जैसा होता है

→ Request rate limit हो जाते हैं

↳ 429 code send

↳ Drop it

↳ or enqueue it for later process



Q18 Design Rate limiter (3/4) part 3 → Rate limiter in Distributed Environment →

Q → Any challenges for Rate Limiter in Distributed Environment ?

A → There are mainly 2 challenges :-

↳ (i) Race Condition

↳ (ii) Synchronization Issue

Q → Can you elaborate ?

A → Race Condition → Rate limiter at high level

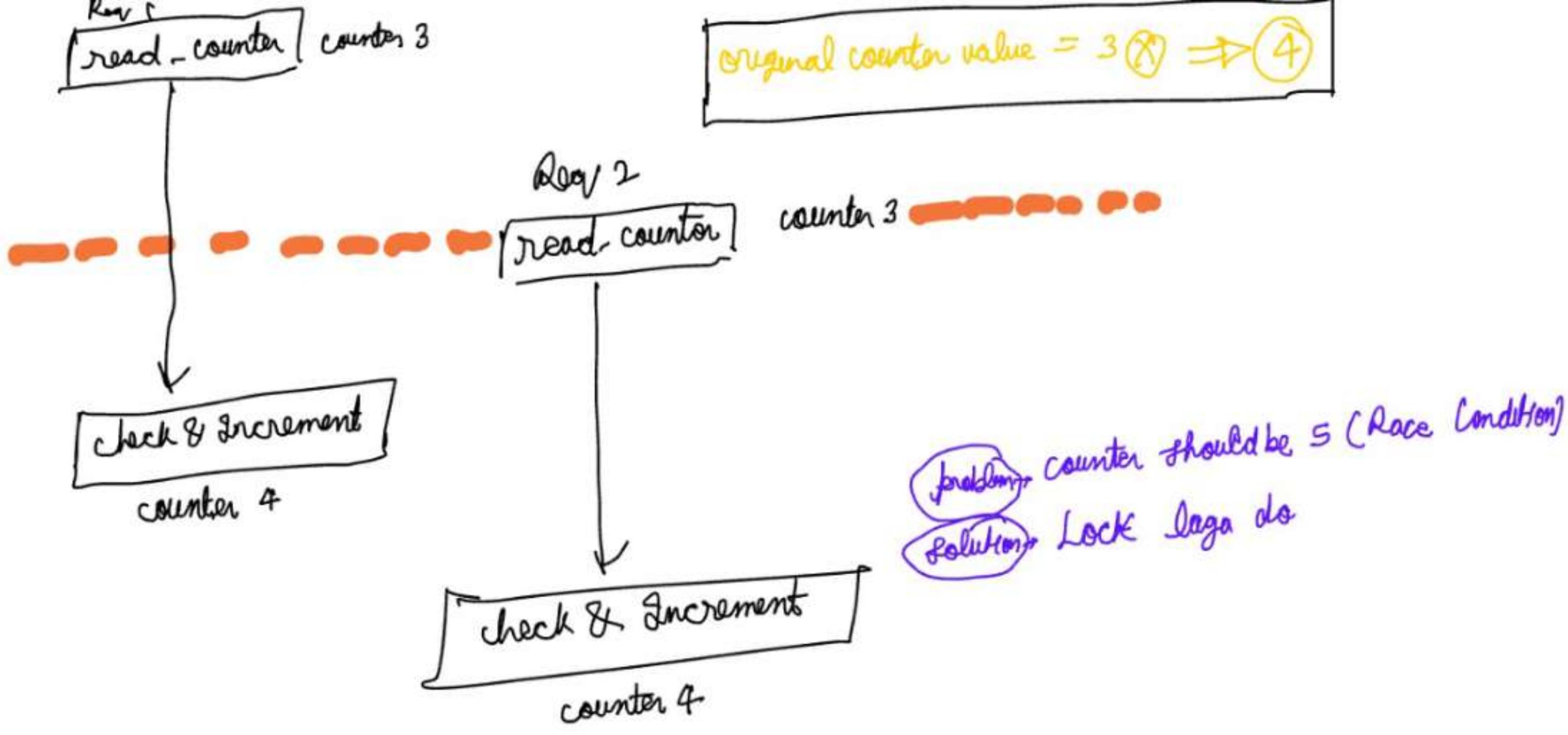
→ Read Counter from Redis

→ Check if $(\text{counter} + 1) > \text{Threshold}$

→ If Not, INC

If concurrently request visit then

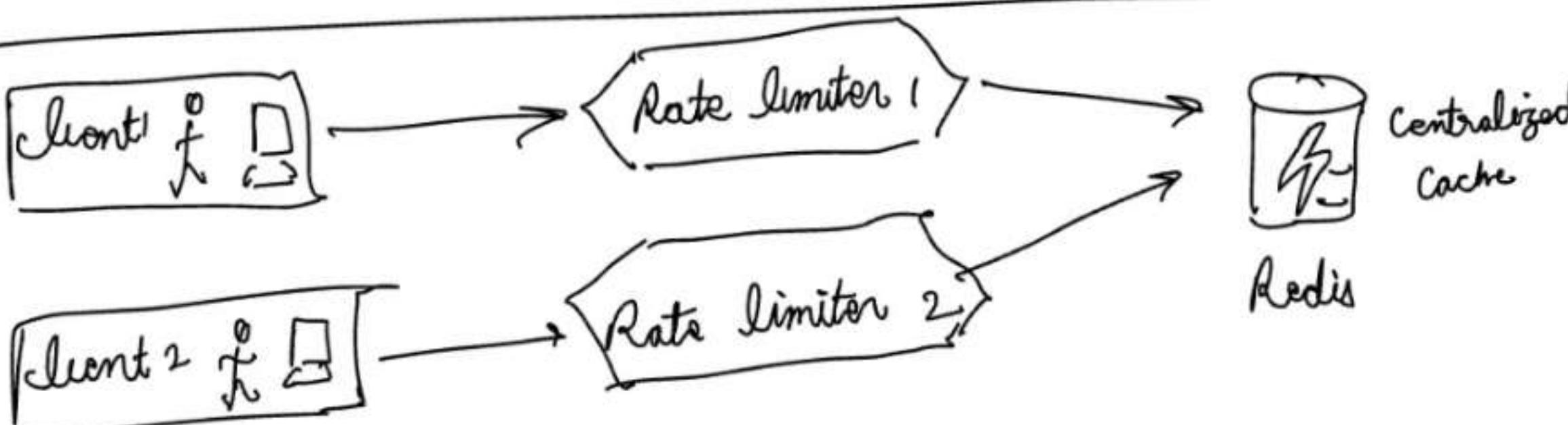
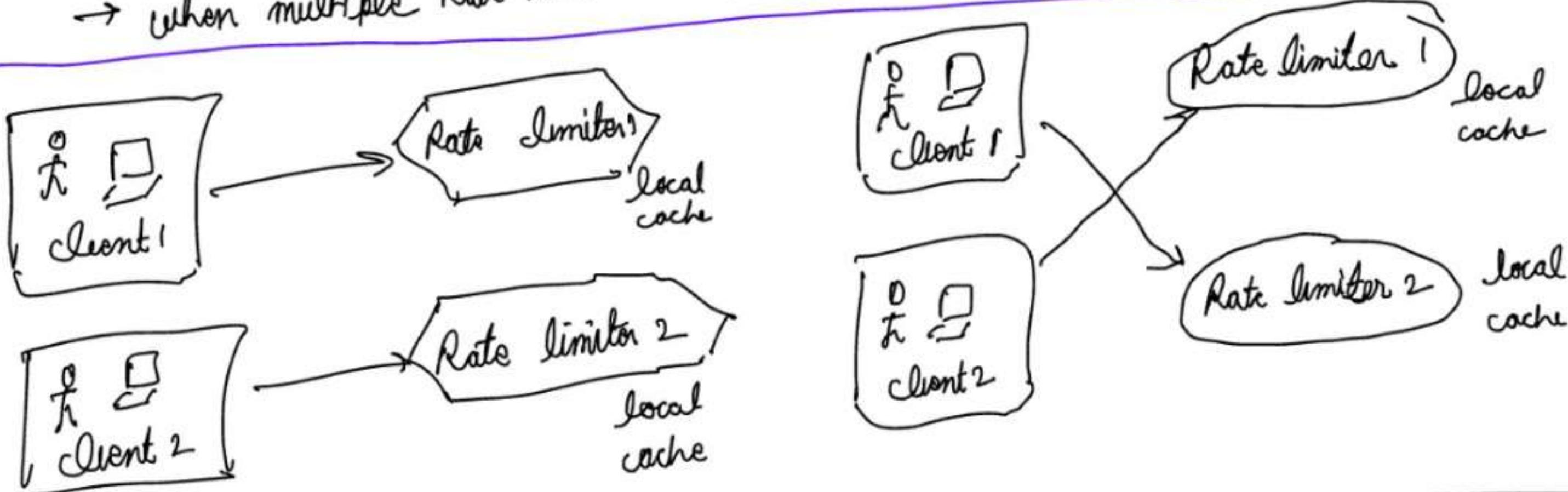
Original counter value = 3



Note → Locks are obvious sol", But it slows down the system.

* Synchronization Issue →

- Millions of users
- , Rate Limiter is not enough
- when multiple Rate Limiter used, synchronization must be there



* Synchronization Issue →

- millions of users
- , Rate Limiter is not enough
- when multiple Rate Limiter used, synchronization must be there

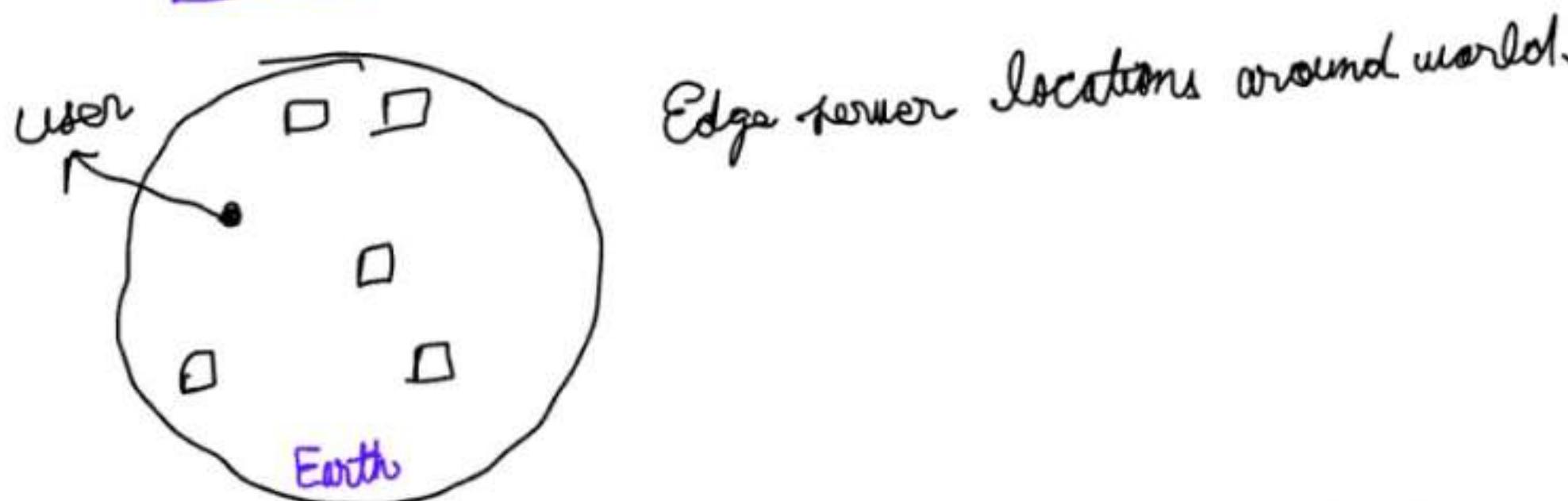
14 Design Rate Limiter (3/4) (part 4) Performance Optimization & Monitoring →

Q → Any areas where we can improve?

A → Multi Data Center Setup.

→ Better synchronize data.

→ Latency ↑; users far away from data center



Interviewer → How will you check if your rate limiter is effective?

Candidate → By monitoring metrics (gathering analytics data)

Mainly 2 things to ensure :-

* Rate Limiting Algorithm - Effective?

* Rate Limiting Rules - Effective?

Interviewer → Rate limiting rules are too strict?

Sudden increase in traffic in flash sales?

क्या करें?

Candidate → Relax the rules a bit bcz otherwise many valid request might be throttled.

Replace algorithm to support burst of traffic → **TOKEN BUCKET** टोकन बैकेट

↳ Burst of traffic supported समीक्षा

20 Design Rate Limiter (4/4) wrapup →

Additional talking points → If time allows!!

* Hard Rate Limiting ⇒ Request \leq Threshold (Always)

* Soft Rate Limiting ⇒ Can exceed threshold for a short period.

* Can avoid your requests being Rate limited?

Interviewer → As a client, what will you do to avoid your requests being Rate limited?

Candidate → * Use Client Cache → Avoid frequent calls to API.

* Do not send too many requests as per the threshold limit within a time frame.

* Code to catch exceptions.

* Add sufficient back off time to retry request call.

Recap → Rate Limiting algorithms

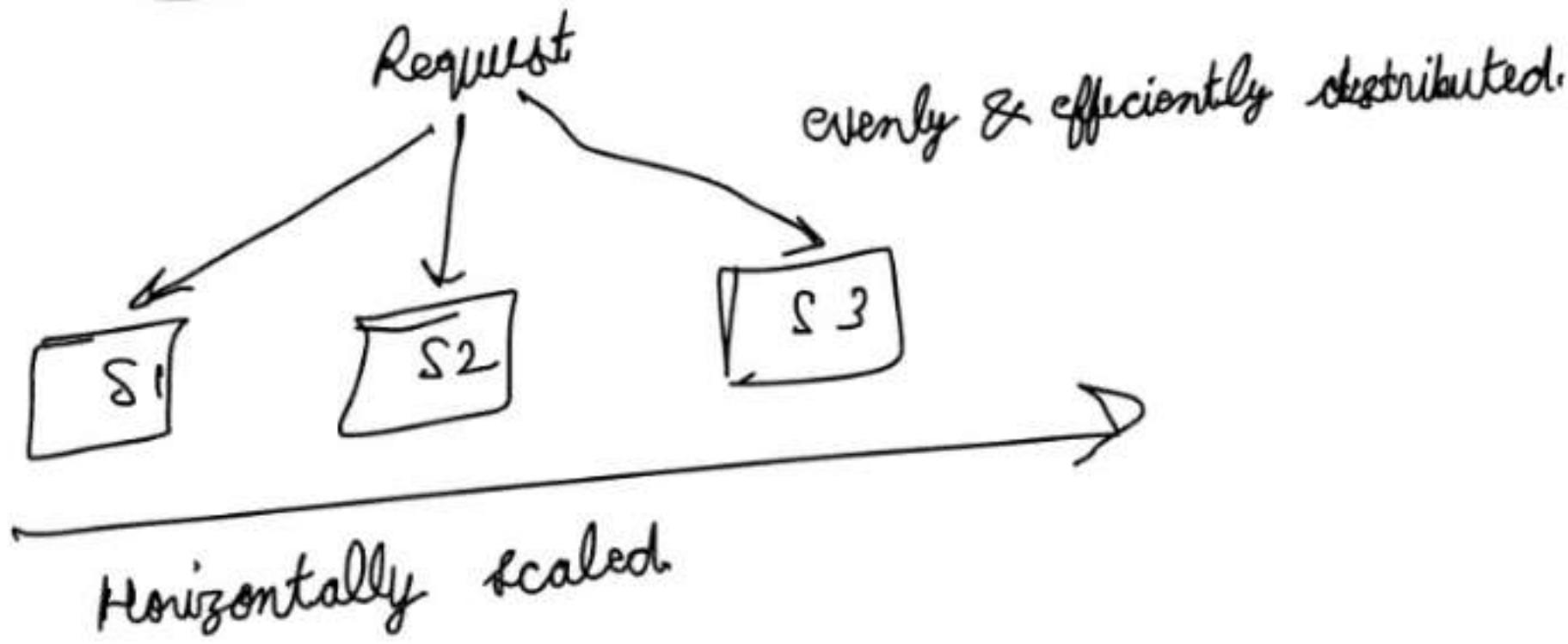
→ System architecture

→ Distributed Environments

→ Optimization & monitoring

Q1 Design Consistent Hashing (what's the problem? why need it?) →

* Remember Horizontal Scaling →



e.g. → Employee Record

Email-id	Details
xyz@gmail.com	~
abc@gmail.com	—
def@gmail.com	—
ghi@gmail.com	—

→ no of records ↑↑↑ Increasing ↑↑↑

→ can't store in one server

→ Distribute data in multiple servers

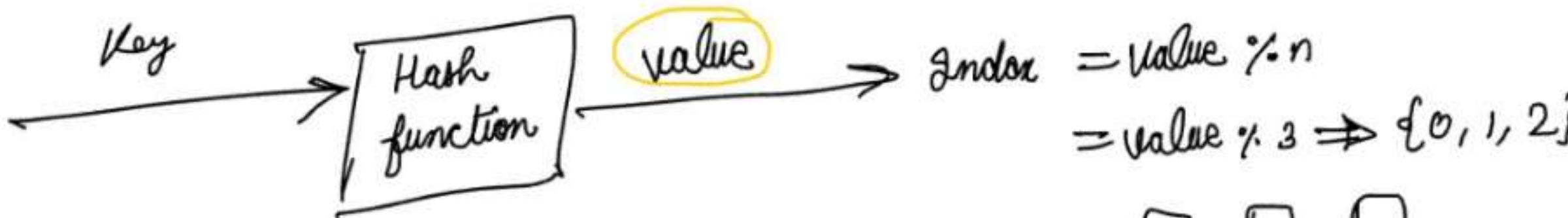
→ How to know which server to look up for particular email-id?

Problem :-

we will use email-id as the key

cache servers = n

server index = $\text{hash}(\text{key}) \% n$



Email (Key)	Hash	Server
john@example.com (S3)	$89 \% 3 = 2$	2 (S3)
mark@example.com (S1)	$30 \% 3 = 0$	0 (S1)
adam@example.com (S3)	$47 \% 3 = 2$	2 (S3)
smith@example.com (S2)	$52 \% 3 = 1$	1 (S2)
alex@example.com (S1)	$75 \% 3 = 0$	0 (S1)
bob@example.com (S2)	$22 \% 3 = 1$	1 (S2)

server [index = $\text{hash} \% 3$] ($N=3$)

S1 index = 0

S2 index = 1

S3 index = 2

Email (Key)	Hash	Server	Index = Hash % 2	$N=3/2$
john@example.com	$89 \mod 2 = 1$	1 (S2) <input checked="" type="checkbox"/>	s1 <input type="checkbox"/> index = 0	
mark@example.com	$30 \mod 2 = 0$	0 (S1) <input checked="" type="checkbox"/>	s2 <input type="checkbox"/> index = 1	
adam@example.com	$47 \mod 2 = 1$	1 (S2) <input checked="" type="checkbox"/>	s3 <input type="checkbox"/> cache miss	
smith@example.com	$52 \mod 2 = 0$	0 (S1) <input checked="" type="checkbox"/>		
alex@example.com	$75 \mod 2 = 1$	1 (S2) <input checked="" type="checkbox"/>		
bob@example.com	$22 \mod 2 = 0$	0 (S1) <input checked="" type="checkbox"/>		

→ Refreshing

22 Design Consistent Hashing How it works?

Previous Video → Rehashing ~~will fail~~ after 1 server failure

"Consistent Hashing" comes to the rescue → It's independent of no of servers

Email (Key)	Hash	Server	Index = Hash % 2
john@example.com	$89 \mod 3 = 2$	2 (S3) <input type="checkbox"/>	s1 <input type="checkbox"/> index = 0
mark@example.com	$30 \mod 3 = 0$	0 (S1) <input type="checkbox"/>	s2 John <input type="checkbox"/> index = 1
adam@example.com	$47 \mod 3 = 2$	2 (S3) <input type="checkbox"/>	s3 Crashed john <input type="checkbox"/> cache miss
smith@example.com	$52 \mod 3 = 1$	1 (S2) <input type="checkbox"/>	
alex@example.com	$75 \mod 3 = 0$	0 (S1) <input type="checkbox"/>	
bob@example.com	$22 \mod 3 = 1$	1 (S2) <input type="checkbox"/>	

→ Does not depends on no of servers

→ Use a hash function.

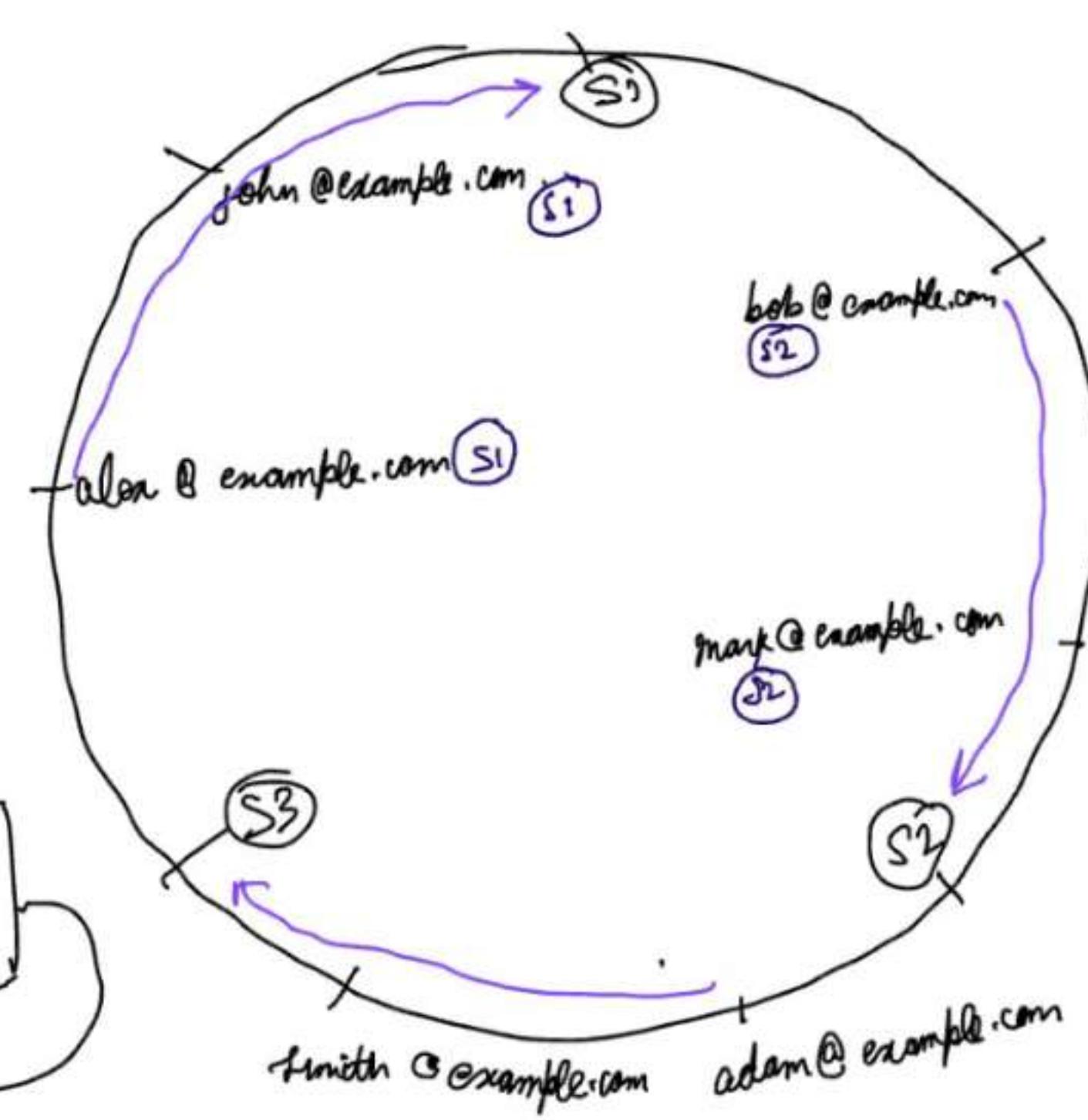
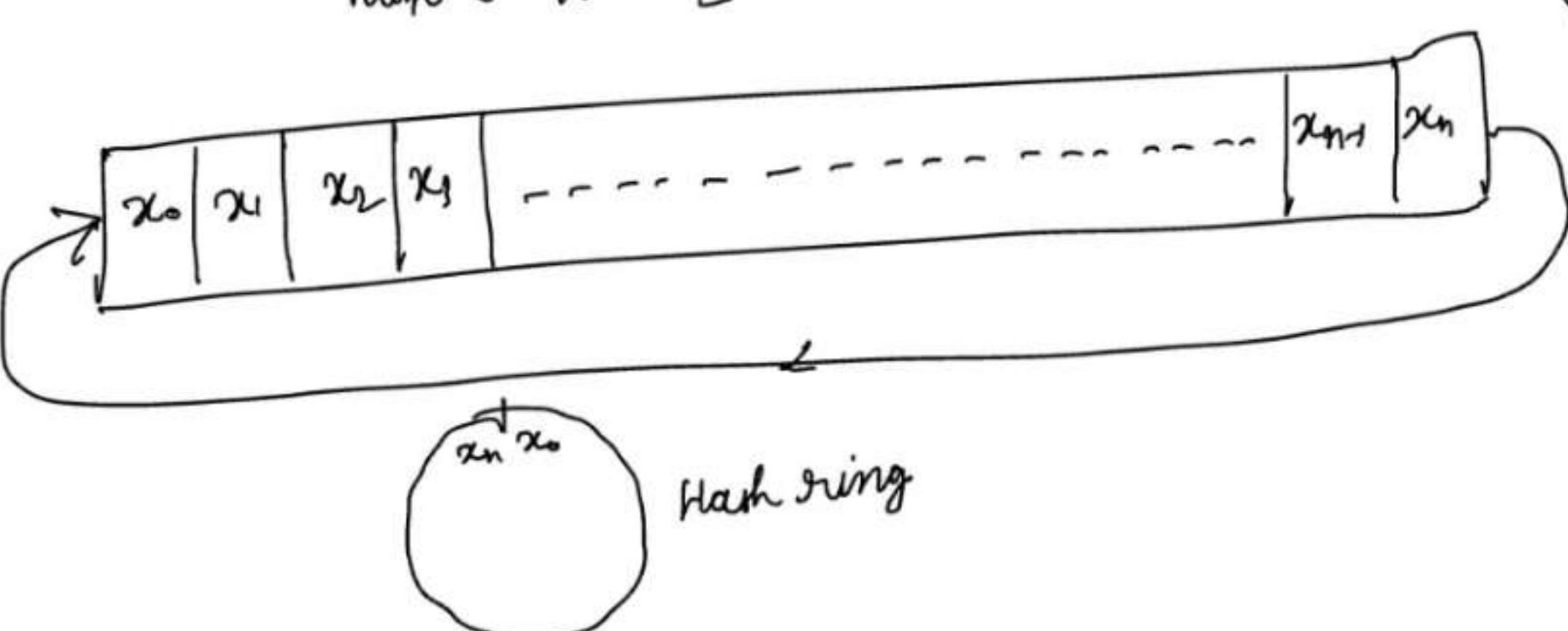
↳ maps servers in the Hash Ring

↳ maps keys in the Hash ring.

→ Locate key in ring using hash function

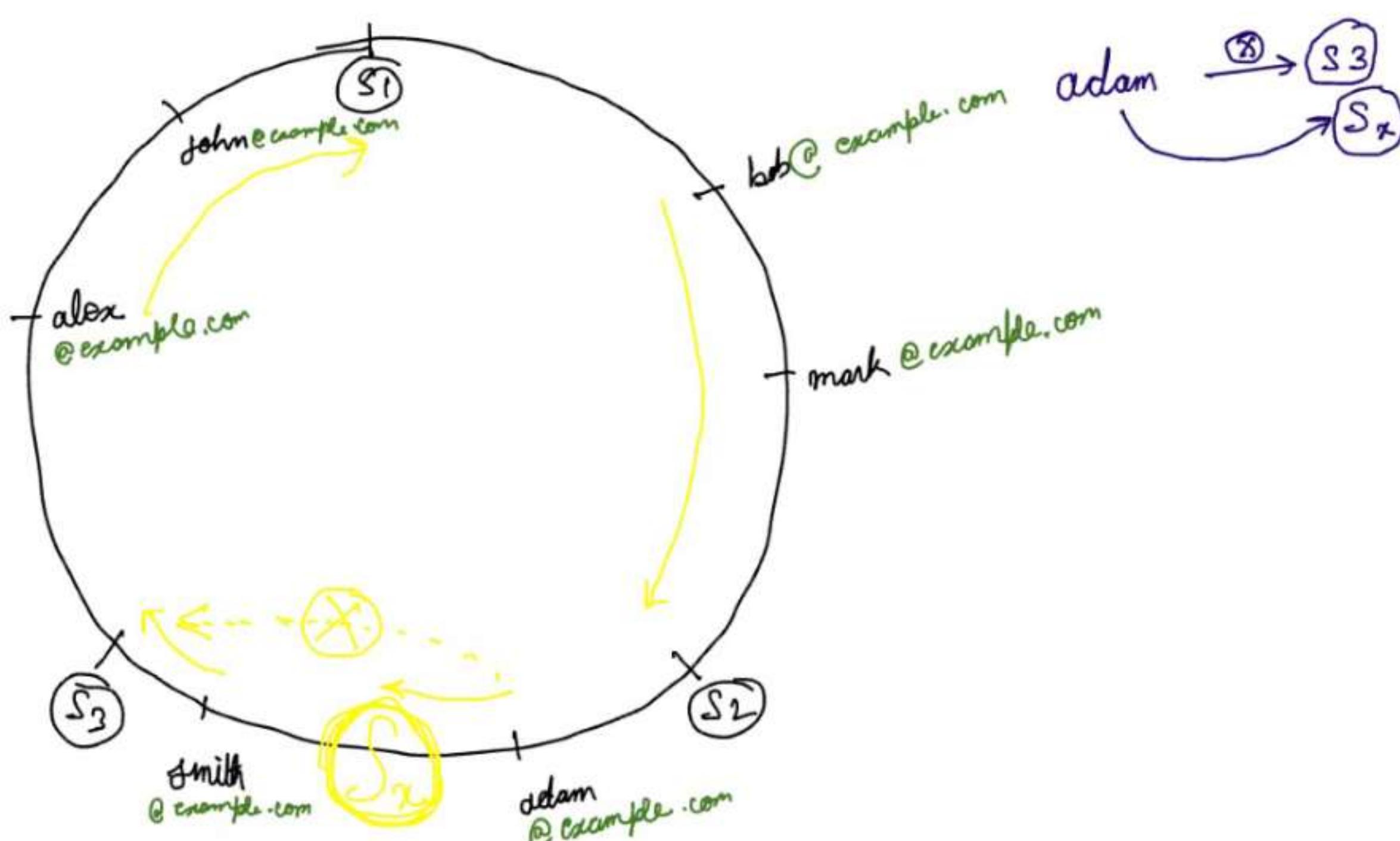
→ Go clockwise until you find a server.

$$\text{hash}(\text{key}) = [x_0, x_n]$$

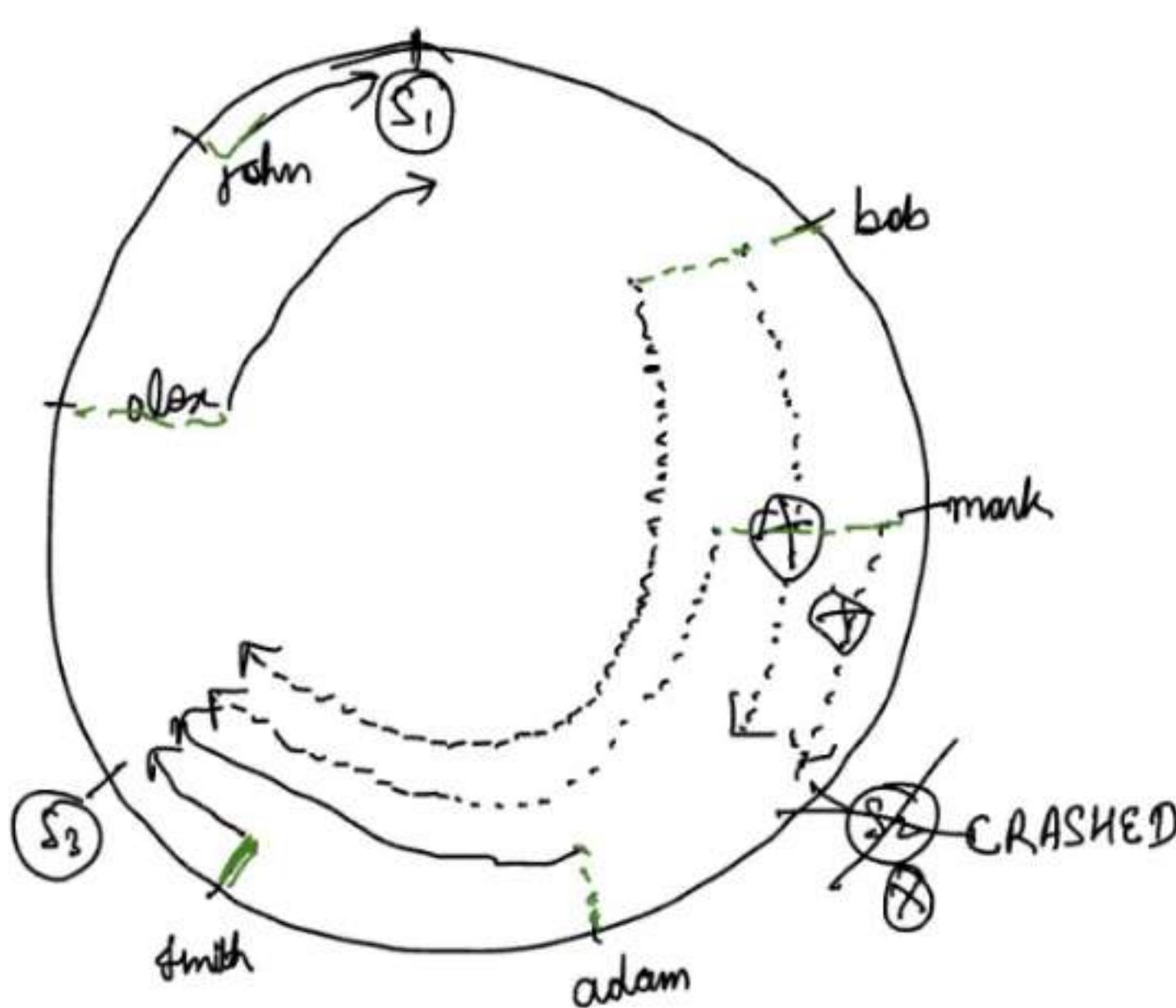


- Hash Server
- Add a Server
- Hash Key
- Remove a Server

* Adding a Server →



* Removing a Server →



S_2 Crash \rightarrow $S_1 \leftarrow S_2$



Q3 Design Consistent Hashing (Two Problems with Basic Approach) →

Two Issues → 😐

Recap → map servers and keys on to a ring. [Basically ring nahi hota ID space hota hai]

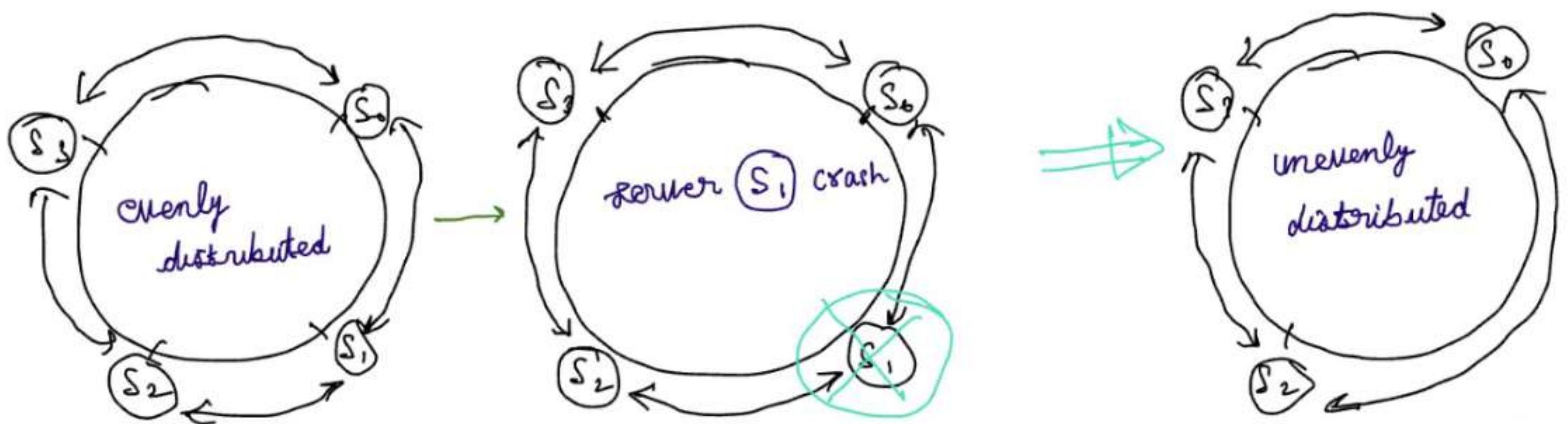
→ Using uniformly distributed hash function.

→ To find server, key to clockwise $\overbrace{S_1 S_2 \dots}$!

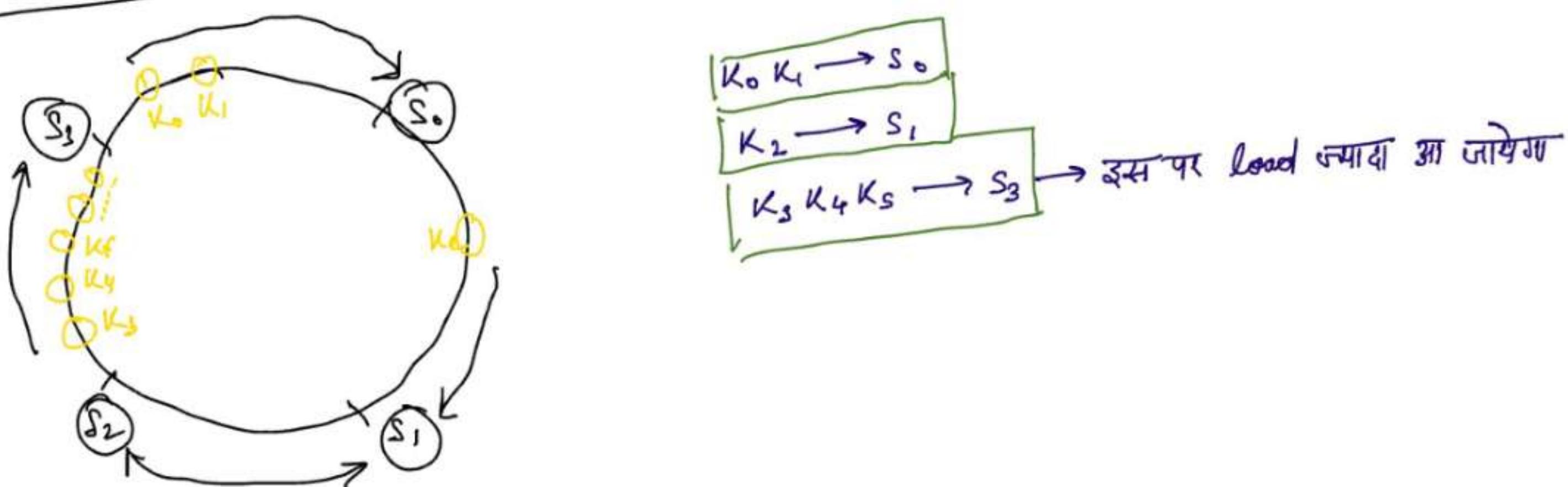
Interviewer → Can you point out problems above??

* 2 problems with above Basic Approach →

Problem ① Ring में अस्टी servers के बीच का partition हमेह same होना impossible है,
→ क्योंकि servers add & remove भी हो सकते हैं।



Problem ② Keys का Non-Uniform distribution भी हो सकता है। हमेही evenly distribute नहीं होगा



Q → What to do to solve those problems above?

A → Virtual Node

Q4 Design Consistent Hashing (Virtual Nodes) →

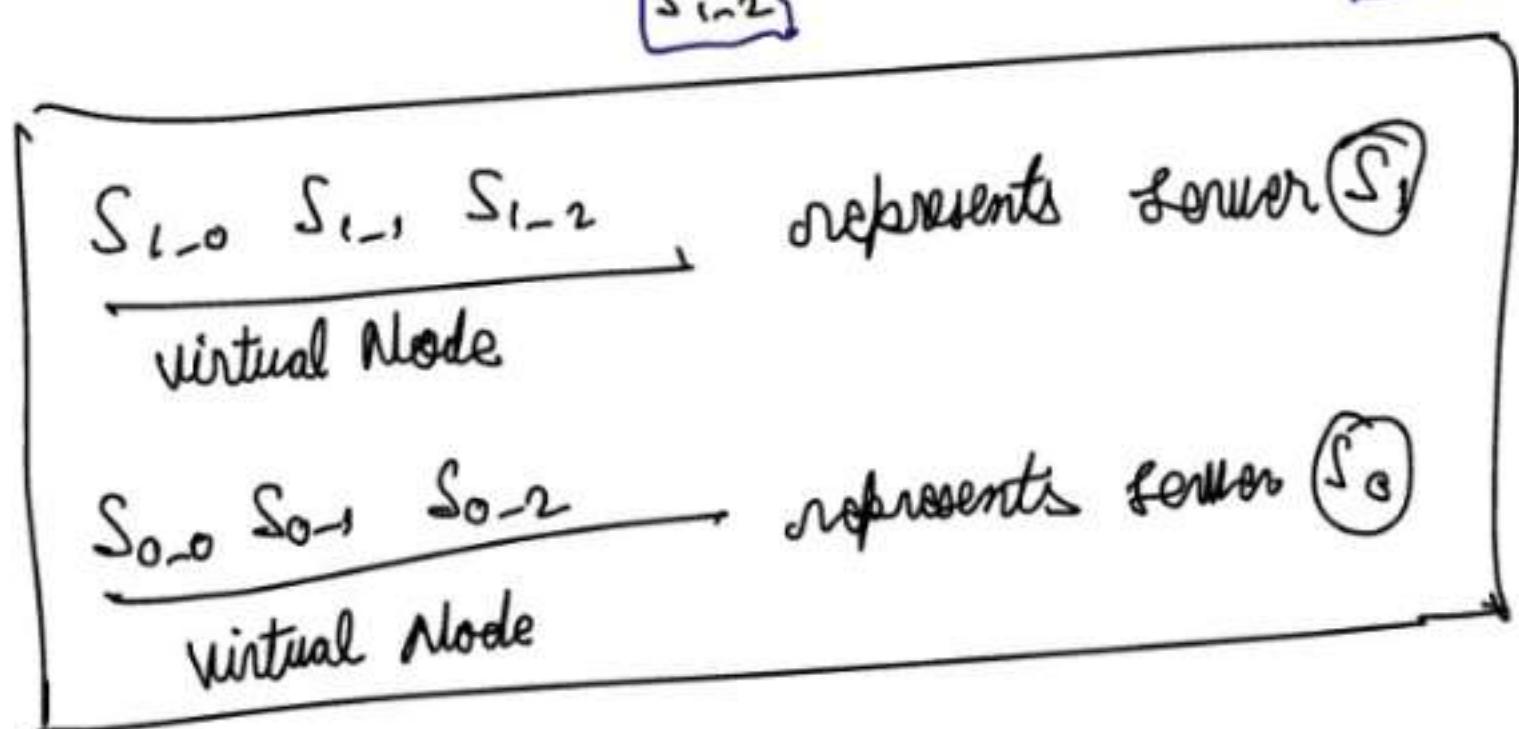
Virtual Nodes →

→ Refers to real node actually



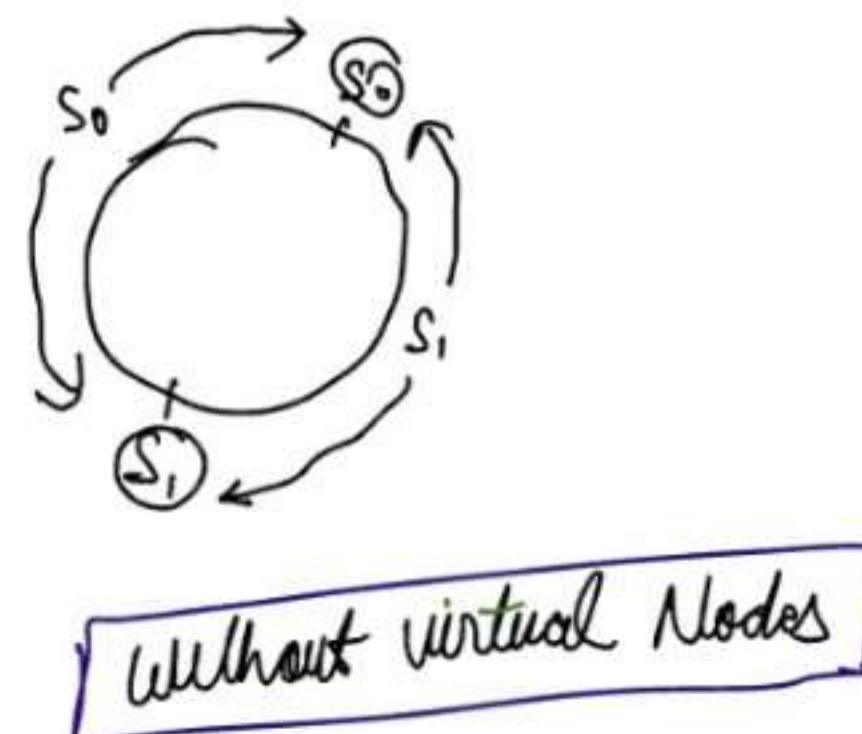
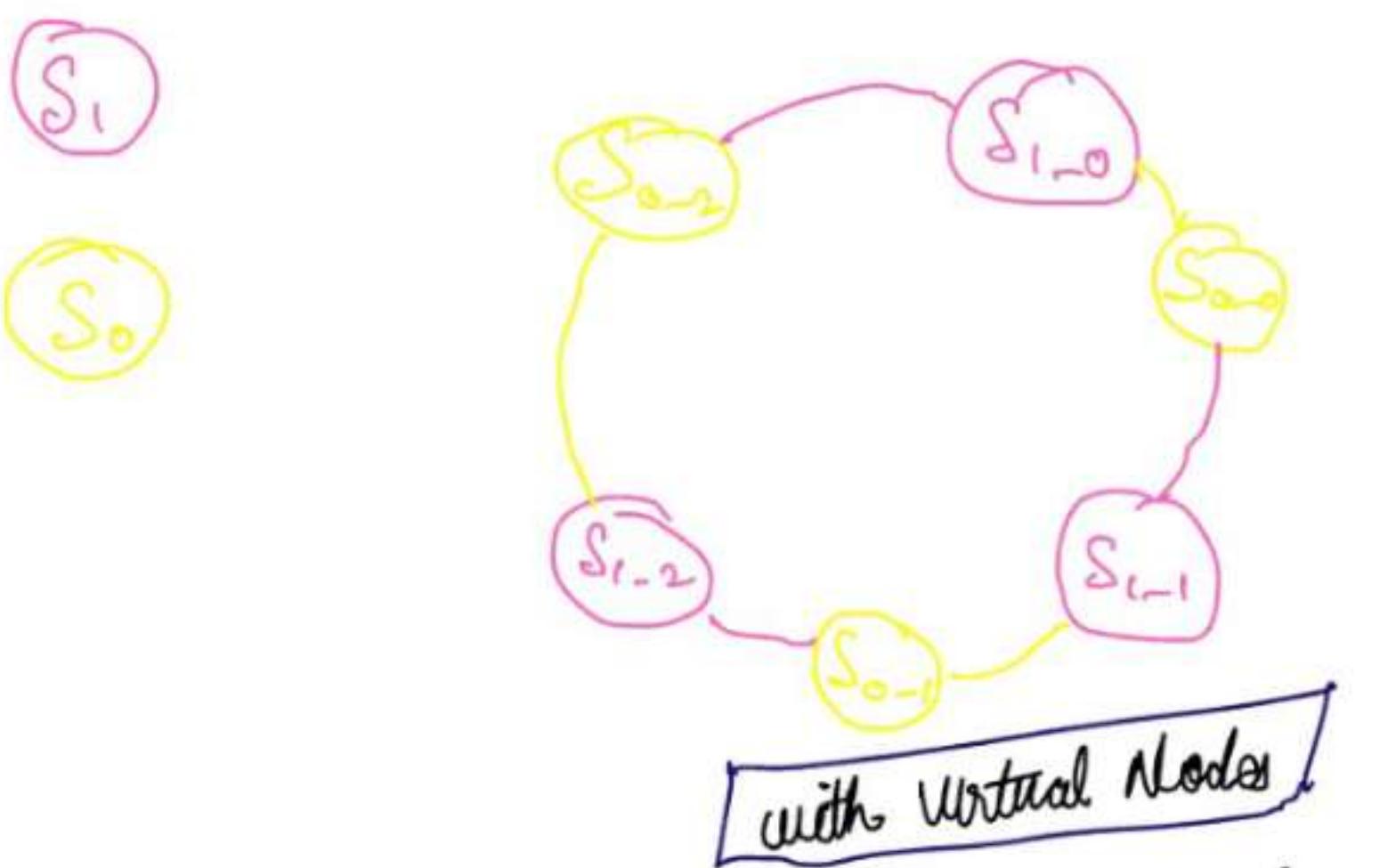
3 Virtual Nodes created

Real scenario mein virtual node bahut syada hota hai.

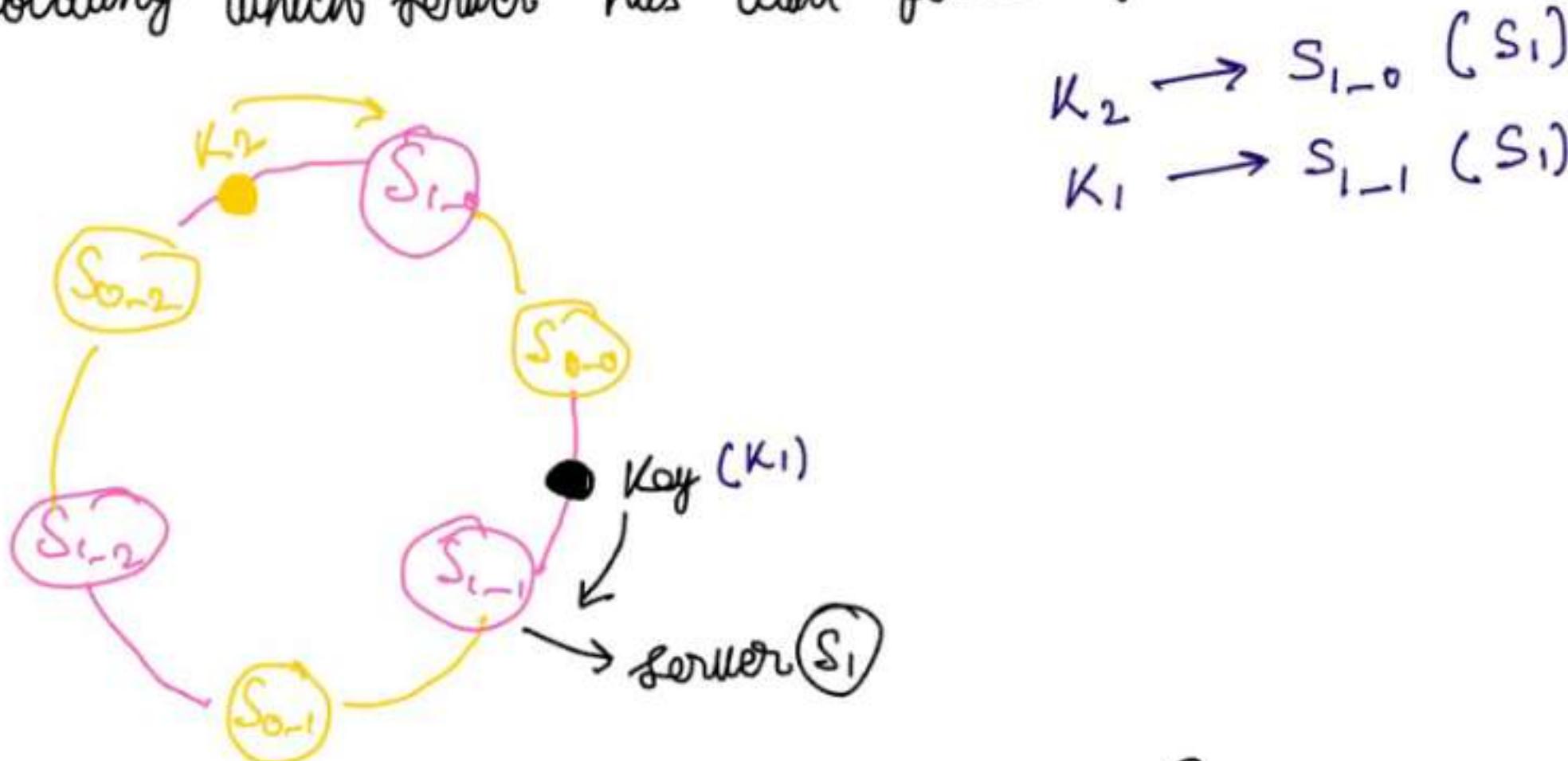


Note → Real World system \neq Virtual Nodes → Real world system mein virtual nodes bahut zyada hotे हैं।

→ Multiple small ranges \Rightarrow partitions \rightarrow



* Locating which server has data for a key.



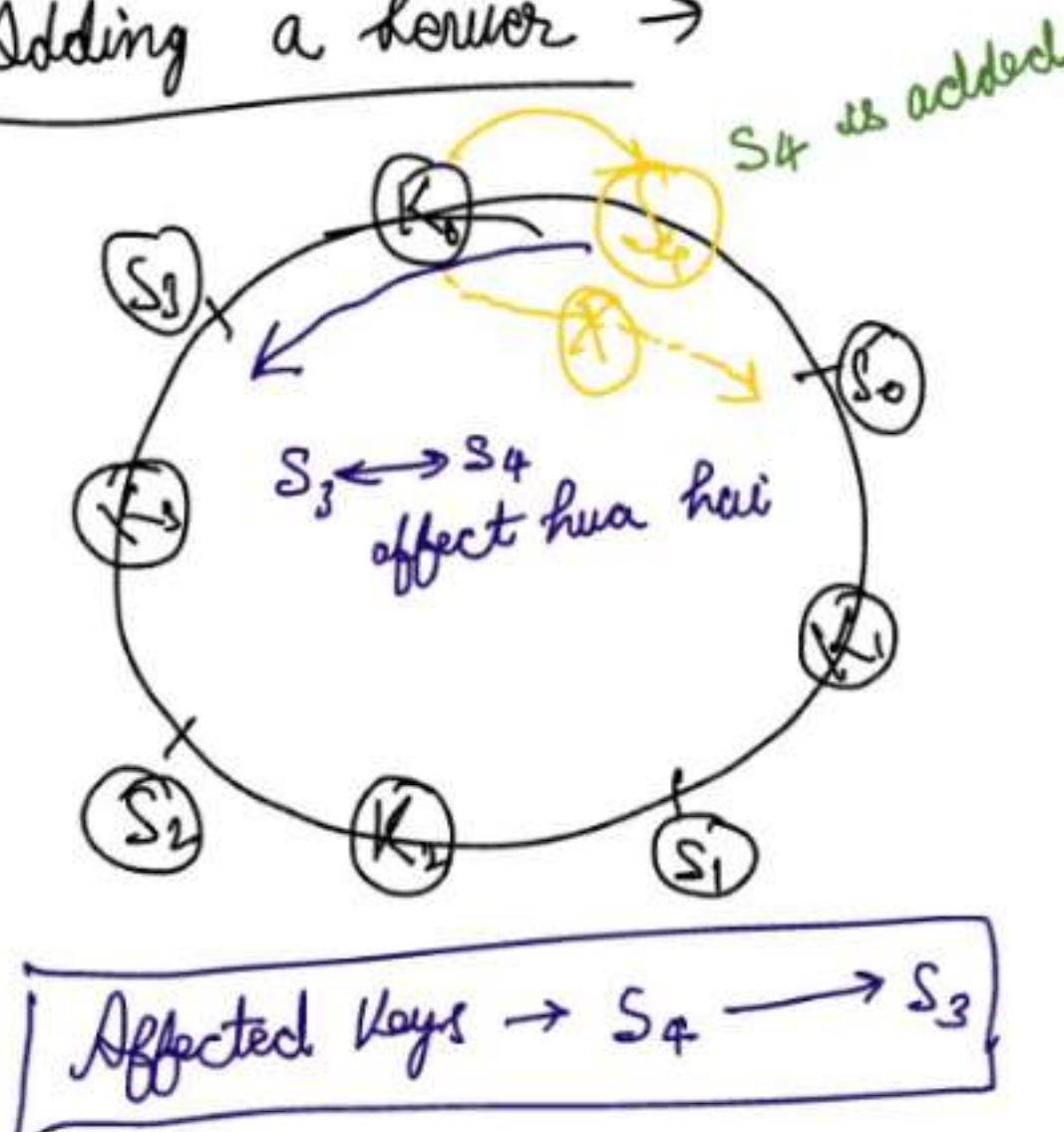
Interviewer \rightarrow How will you find affected Keys?

When :-

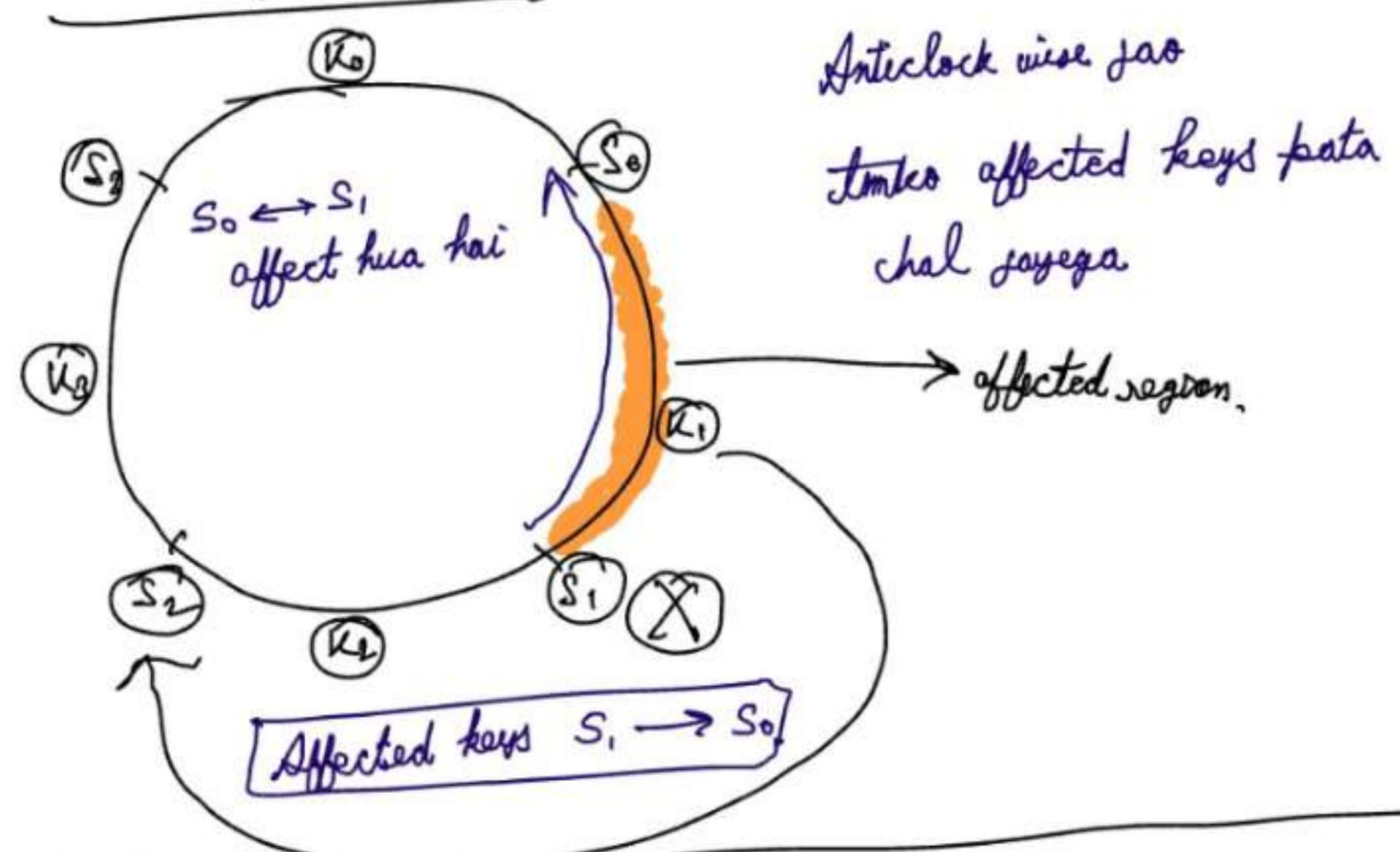
↳ Server is added

↳ Server is removed

* Adding a server \rightarrow



* Removing a server \rightarrow



(Q5) Design Consistent Hashing (Wrap up) \rightarrow

Note \rightarrow Always do a quick recap/summary for interviewer

Discord mein consistent hashing use hota hai.

- * We now know
 - why need?
 - how it works?
 - Benefits

server

- ↳ Minimum keys are redistributed (add, remove)
- ↳ Since data can now be evenly distributed
 - Easy to "Horizontally Scale".
- ↳ Avoid hotspot key problem.

- * Widely used in real-world systems →
 - Amazon Dynamo DB
 - Apache Cassandra
 - Discord Chat Application.

Q26 Design Key-Value Store →

- Also known as - Key value Database
- Non-relational Database (NoSQL Database)
- Key-Value pair → Requirement Gathering
- Key must be unique → Requirement Gathering
- Key
 - Plain Text → e.g. last_login_at
 - Hashed Key → e.g. → 259CFD25
- Value → strings, list, object etc.

Key	Value
145	Sachin
260	Amar
330	Skhar
210	Anthony

Interviewer → I want you to design a key-value store that supports!-

- put(Key, Value) // insert value associated with "key"
- get(key) // get "value" associated with "key"

Candidate → What should be the size of key value pair?

Interviewer → Less than 10KB
It should be able to store big data

Candidate → What about 'Availability', 'Scalability'?

Interviewer → **Availability** → HIGH, respond quickly

Scalability → Should be scalable to support large data set.

Note → Scaling should be 'Automatic' → Traffic be load se server add or delete hona chahiye
E.g., → addition / deletion of servers automatically based on traffic.

Candidate → What about consistency & Latency?

Interviewer → **Tunable Consistency**

Low Latency

S_1
Replica slave 1

S_2
Replica slave 2

S_3
Replica slave 3

consistency level

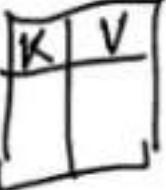
Level = 2

Master

M

Q7) Design Key-Value Store (part 2) single server key value store →

→ Easy to implement

→ Hash Table use करो 

→ will keep everything in memory

→ Memory access is fast but memory limited space constraint.

Interviewer → Any improvement you can do?

Candidate → Data को compress करो store करो

→ Frequently used data = Memory के लिए

→ less frequently used = Disk के store करो

Interviewer → Will it solve all our problem?

Candidate → No,

Even with above optimization

Single server will easily reach its capacity limit.

Interviewer → Then anything else we can do?

Candidate → We will use **Distributed Key-value store** (Serves many servers) → Spans data के multiple servers & distributes करता है

→ इसके लिए server पर ज्यादा load नहीं होता

Interviewer → Sounds good. Enlighten me more on this.

(28) Design Key Value store (Part 3) → Distributed Key Value store and CAP Theorem →

→ Distributed Hash Table

→ Distributed Key Value pairs across many servers. (कई DB servers अंतर्राष्ट्रीय स्तर पर हैं).

→ Distributed system design के लिए → we must know CAP Theorem

CAP Theorem ⇒

C → Consistency → All candidate see the same data at the same time no matter which node they connect to.

A → Availability → Any client requesting data gets response even if some of the nodes are down.

P → Partition Tolerance

 Partition → Communication break b/w 2 nodes.

 Tolerance → तभी तभी system continues to operate.

CAP Theorem → one of the 3 properties has to be sacrificed to support 2 properties.
"All 3 impossible"

→ Now a days, key-value stores are available

as →

(A) CP systems → sacrificed → A → Availability

(B) AP system → sacrificed → C → Consistency

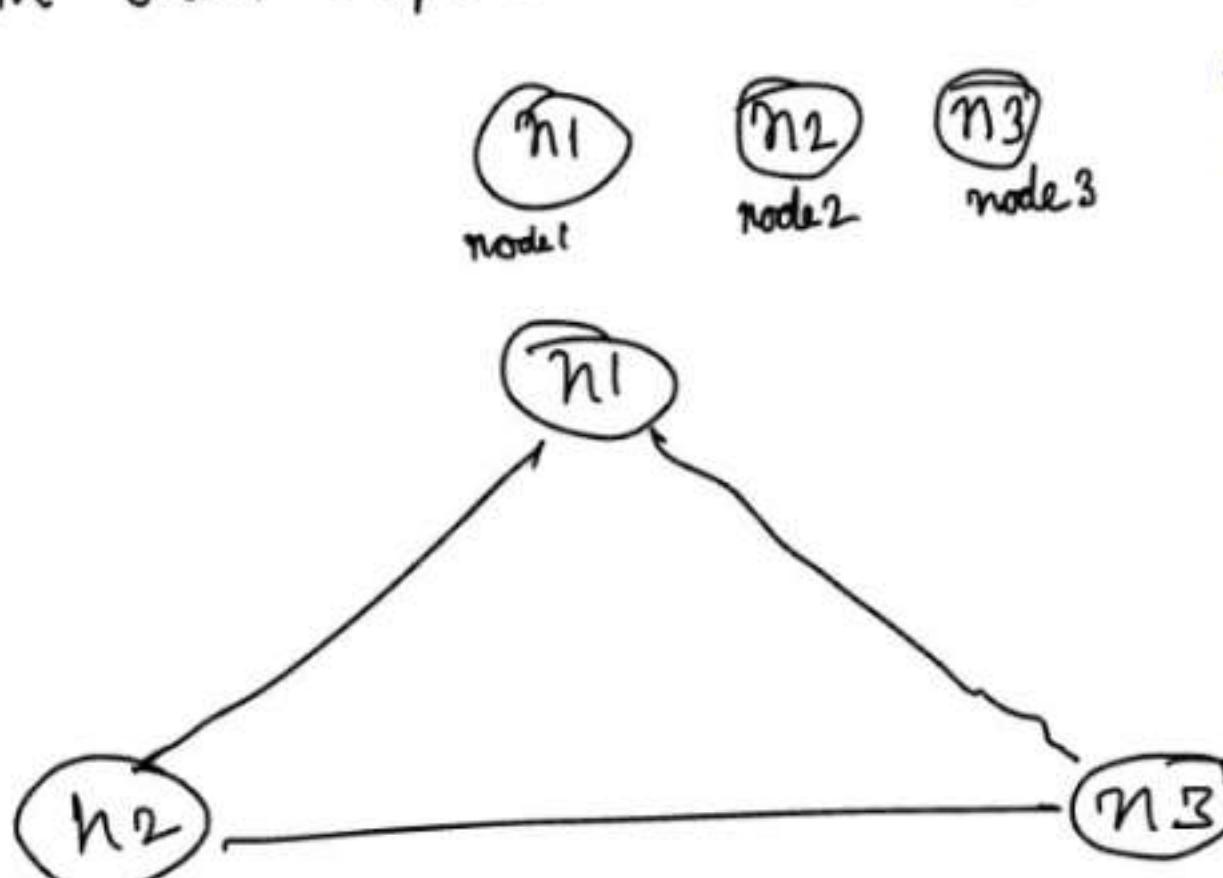
(C) CA system → sacrificed → P → Partition Tolerance

since n/w failure is unavoidable (C) is not possible → Real life में CA possible नहीं

* Example से समझें :-

→ Distributed System में data is usually replicated multiple times.

→ Assume data replicated on 3 replica nodes :-



Ideal situation → n/w partition वही समाज़

Data written on n1

replicated to n2 & n3

consistency & availability achieved

node1 = n1
node2 = n2
node3 = n3

जो data n1 में है वो data n2 & n3 में भी है।

C }
A }
P }

But real life में इसा नहीं होता

Apko P तो रखना ही पड़गा



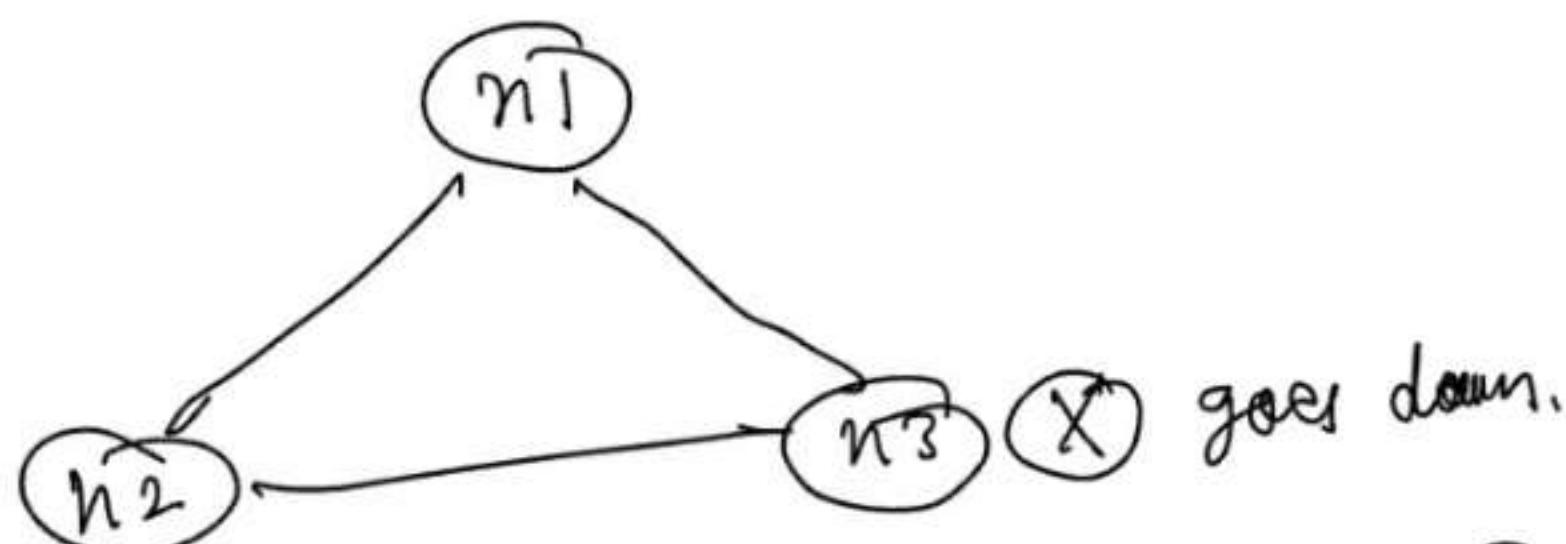
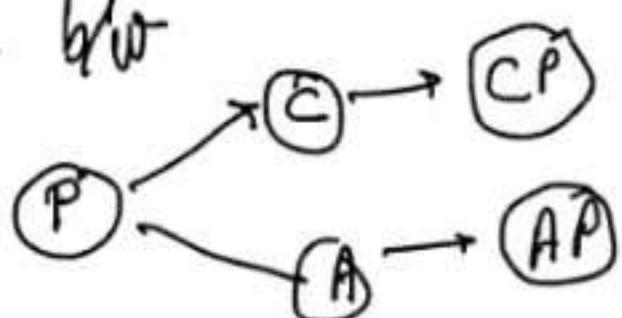
पहले 2 संभावना होंगी Real life में

Real World Distributed System →

- Partition can't be avoided
 - ~~पर्टिशन~~ Partition tolerance तो रखना पड़ेगा

Now you need to choose b/w

- ④ Consistency
 - ⑤ Scalability



- n_3 goes down, can't communicate to n_1 & n_2
 - clients writes on n_1 OR n_2 ; won't propagate to n_3
 - clients writes on n_3 , won't propagate to n_1 & n_2

- If we choose consistency (CP) →

- you want consistency
 - you will block all writes operation on n_1 & n_2 (Bank system) (CP system)
 - system unavailable हो जायगा

C.g → Bank System

- ↳ up to date Balance info,

↳ Fix inconsistency & return error until resolved.

- If we choose Availability (CA) →

- system keeps accepting reads

→ stale data

→ Wires will be done on n1 & n2

- Unites will be synced to n3 when resolved.

→ will be synced to $\textcircled{N_3}$ when needed

→ will be synced to n₃ when resource
Note → GAP ENT correct combo as per your need is a crucial step for distributed NoValueOne.

Note → GAP का correct answer वाला है।

→ so discuss / ask your interviewer की उनकी कौन सा प्राइवेट AP
CA → Not possible in Real life.

29 Design Key-Value Store (Part 4) Data Partition and Data Replication →

Core Components & Techniques to build Key-Value Store →

- Data Partition
- Data Replication
- Consistency
- Inconsistency Resolution
- Handling failures
- System Architecture Diagram
- Write Path
- Read Path.

Data Partition & Data Replication →

* Data Partition →

→ Large Application

→ Infeasible to fit data in single server.

→ Split data (smaller partition) → **multiple servers**

→ Two challenges :-

↳ Distribute data across multiple server evenly (data to evenly distribute करना)

↳ Minimize data movement when nodes are added or removed. (data movement कम करना चाहिए)

→ What, who helped us to fix above challenges ??.

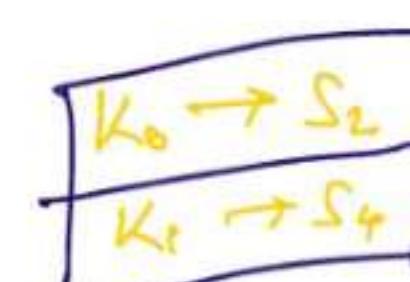
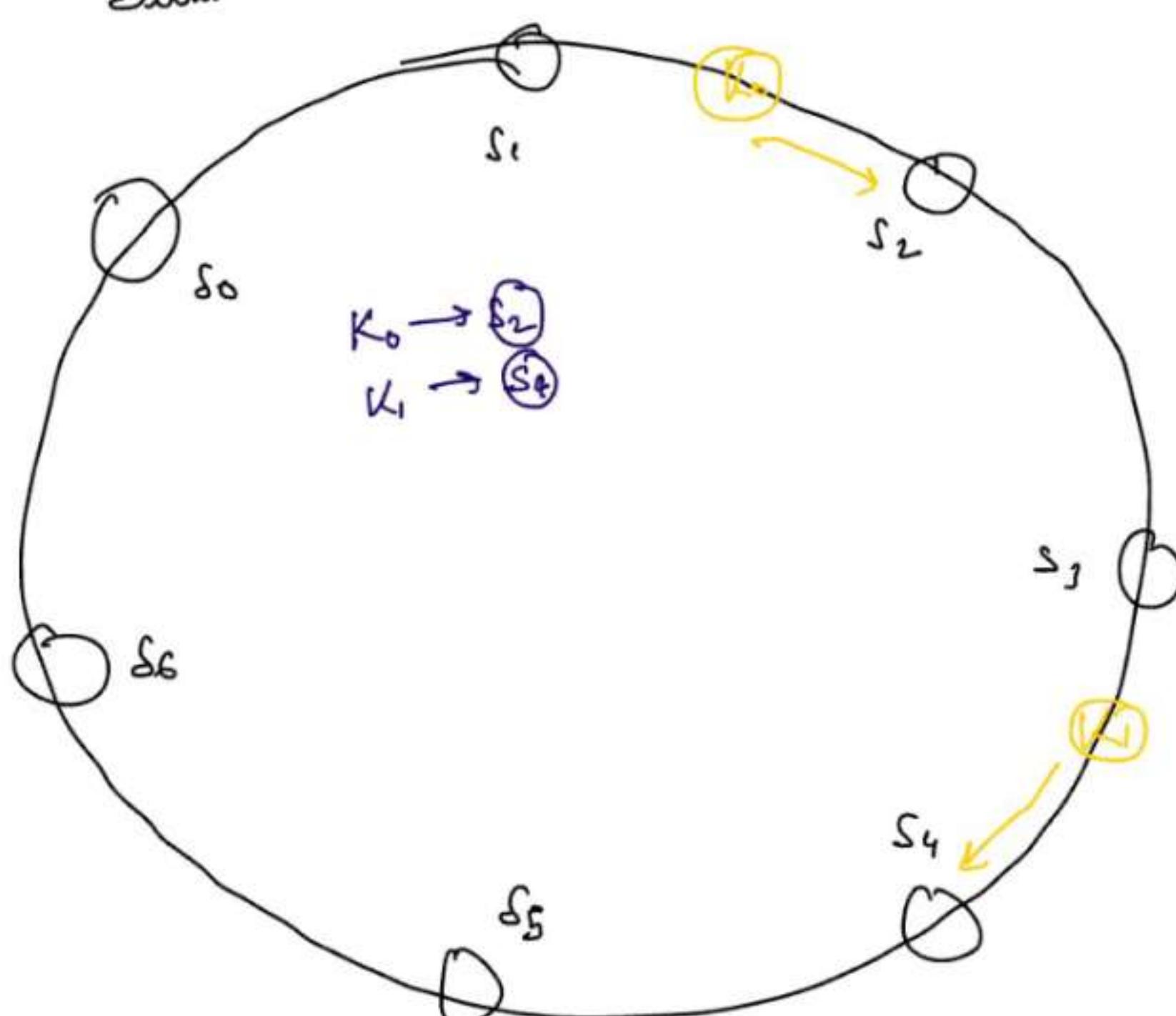
Yes 😊 → **Consistent Hashing**

Derive क्या है अंडा. →

↳ server placed in Hash ring.

↳ key is hashed in same ring.

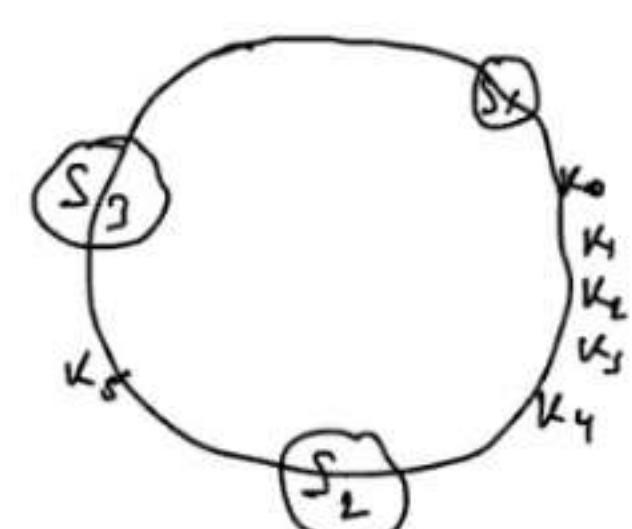
Data → clockwise ओरी



Benefits →

→ वर्ता server add / remove होने पर तब से automatic scaling कर सकते हैं

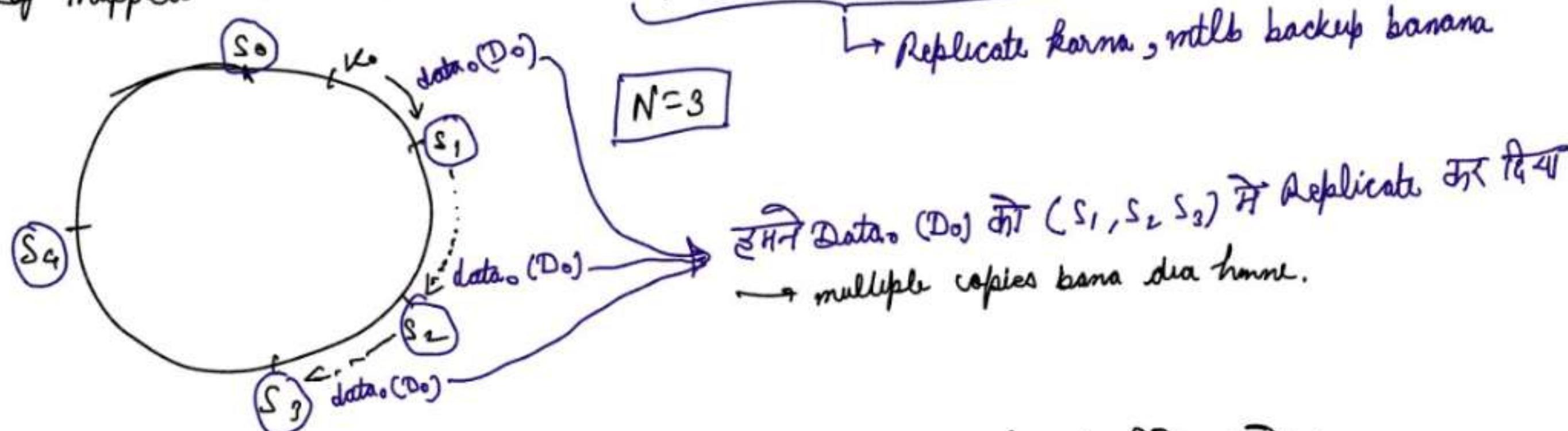
→ virtual node solved uneven distribution in Hash ring



* Data Replication →

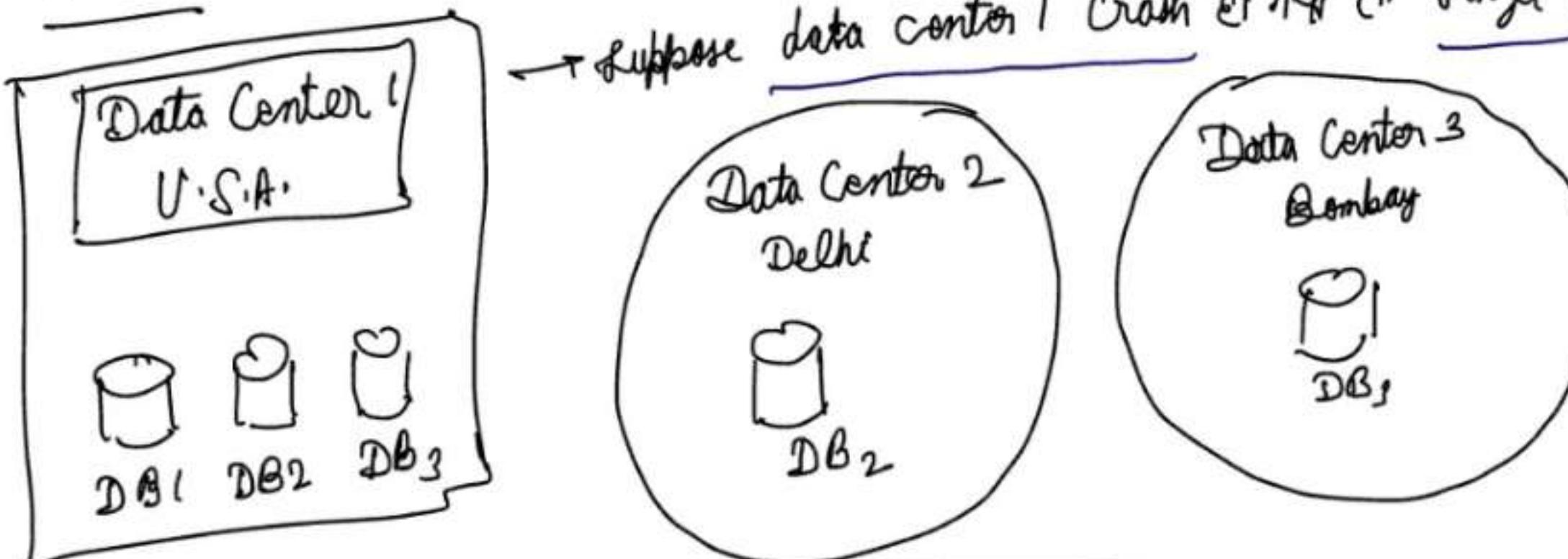
→ For high availability & reliability, data must be replicated to N servers.
configurable failover as per our need

→ Key mapped → walks clockwise → replicate to N servers.



→ Replicate करो तो मला मला data center (different location) में करो।

कैसी?



⑩ Design Key-Value Store (Part 5) (Consistency, Quorum, Tunable Consistency) ⇒

→ Tunable Consistency

→ Consistency →

→ Data replicated in multiple nodes में

→ consistency का उमान देना पड़ता था

→ synchronized होना चाहिए।

N: no of replicas
W: Write quorum of size W
R: Read quorum of size R

मेरे Quorum क्या है?
→ एक तरह का Agreement की W no of ACKआवाजाएं।

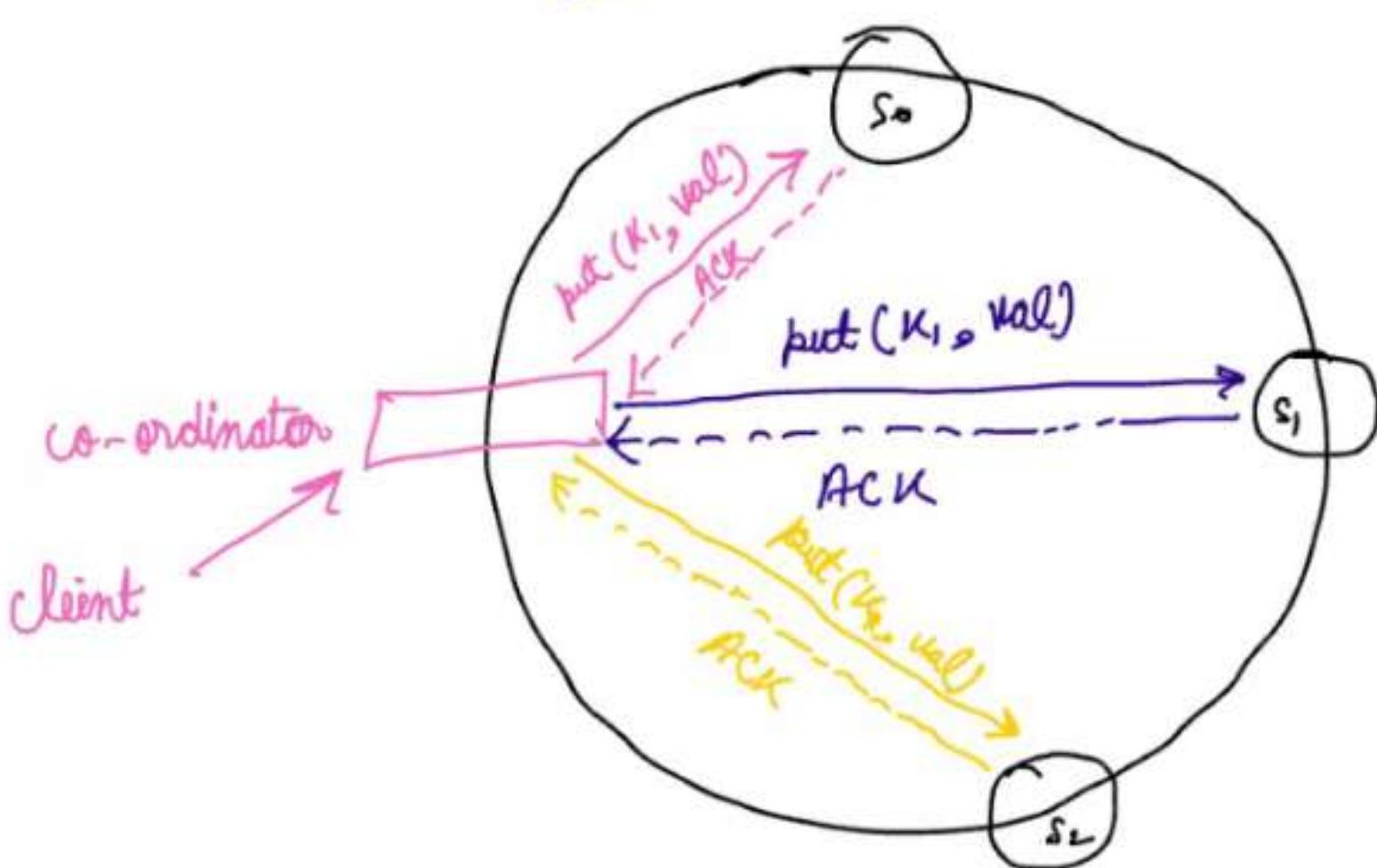
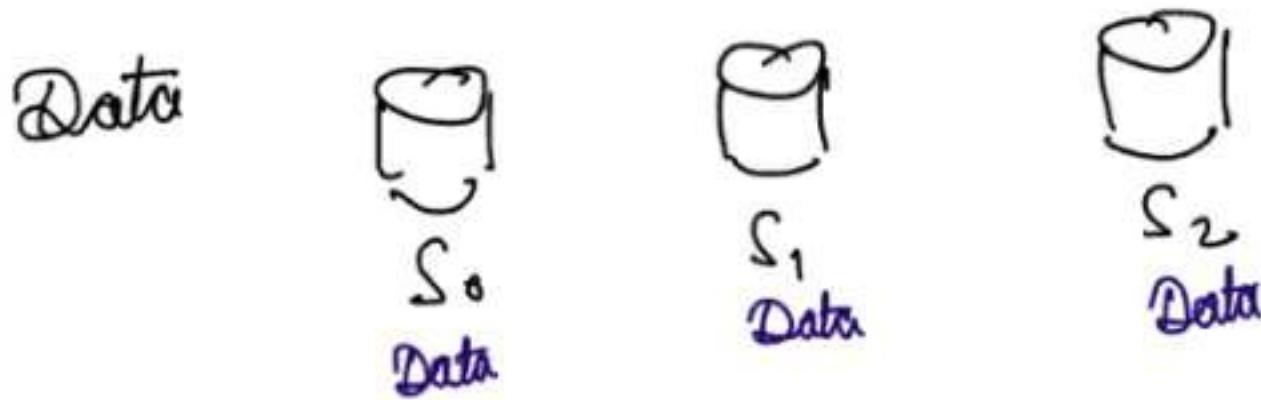
"W": For write operation to be successful कम से कम "W" replicas से ACKआवाजाएं for write operation

"R": For read operation to be successful कम से कम "R" replicas से Response का इच्छा करो।

Example से समझो :-

Let's take $N=3$

We will replicate our "Data" in 3 servers (S_0, S_1, S_2)



$W=1 \rightarrow$ यदि S_0 से ACK मिला तो इस write operation को successful मानो।
 Replica Nodes. we will not wait for S_1 & S_2 का ACK
 co-ordinator → ये एक तरह का MiddleWare है b/w clients & servers

The Configuration of (W, R) & N is a typical tradeoff b/w Consistency & Latency.

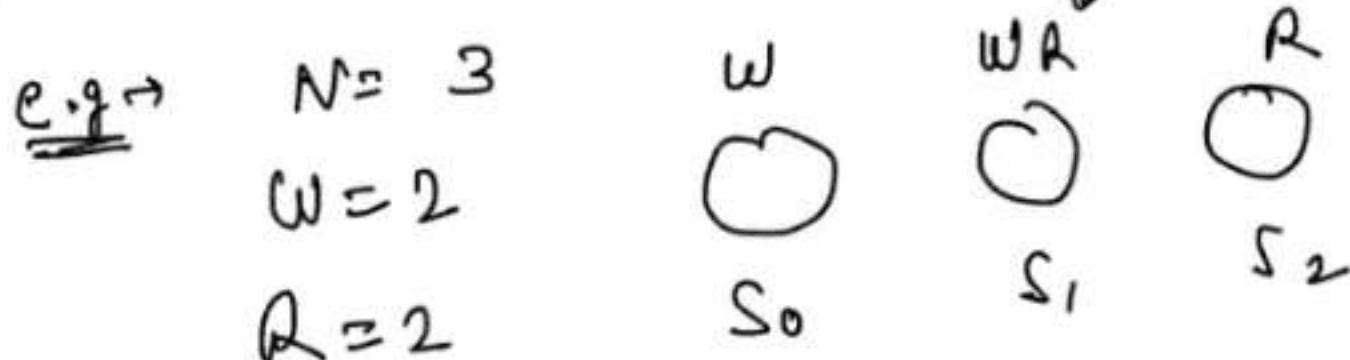
$W=1$: write operation is returned quickly bcz बाकी replicas का wait नहीं करेगा

$R=1$: Read operation is returned quickly bcz बाकी replicas का wait नहीं करेगा

Low Latency (quick response) But Not good Consistency

$W \text{ or } R \geq 1 \Rightarrow$ consistency तो मिल हो जायेगा But latency बढ़ेगा क्योंकि more replicas की response का wait अब भी पड़ेगा।

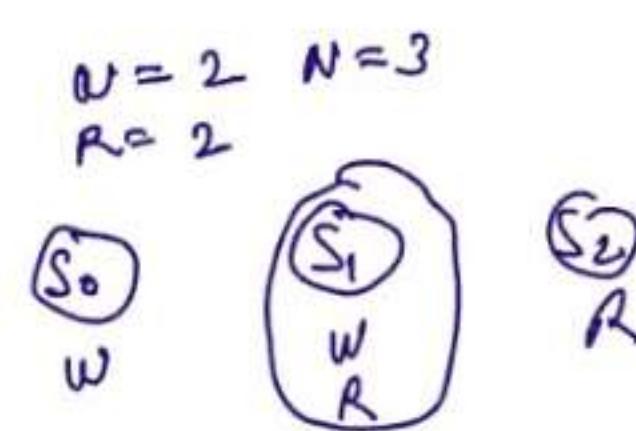
$(W+R) > N \Rightarrow$ Strong consistency Guaranteed क्यों?



$$W+R > N$$

$$2+2 > 3 \checkmark$$

→ कम से कम एक overlapping node हो भी हो जायेगा



$W+R > N \rightarrow$ इस case में कम से कम एक overlapping node जरूर मिलेगा जिसमें Read(R) भी हुआ हो और WRITE(W) भी हुआ हो।

How to configure $[N | W|R]$ to fit our use case.

* $R=1 \quad \} \text{ Fast Read}$
 $W=N$

* $W=1 \quad \} \text{ Fast writes}$
 $R=N$

* $W+R > N \quad \} \text{ strong consistency guaranteed}$
 $\text{e.g. } W=R=2 \quad [N=3]$

* $W + R \leq N$

strong consistency not guaranteed

→ Depending on requirement you can tune W, R & N to achieve your desired consistency.

↳ Tunable consistency

Interviewer → Do you know about "Consistency Models"?

Candidate → It says "~~दर्जा~~" level of consistency ↗!
i.e. → Degree of data consistency.

- * Strong Consistency → Client never sees out of date data. (updated data always)
जिसका एक operation को एकल updated data को करा
- * Weak Consistency → Subsequent read operation may not see latest data always.
There are high chances ki updated data nahi milega
- * Eventual Consistency → Given enough time, all updates propagated to replicas & they are updated.
जोकि किसी item data, it will replicate updated data to all server (node).

Interviewer → Do you think "strong consistency" is good?

Candidate → "strong consistency" के लिए we will have to block read/write until current write is updated to all replicas. → latency ↑

Not good for highly available system.

Interviewer → What will you use for "Key-Value" store?

Candidate → We should use "Eventual Consistency".

C.g. → Dynamo & Cassandra also adopted Eventual Consistency.

Q1) Design Key-Value store (Part 6) Inconsistency Resolution & Vector Clocks →

→ Inconsistency Resolution →

→ Replication gives high availability.

→ But causes inconsistencies among replicas

→ Replication give high availability

→ But causes inconsistencies among replicas.

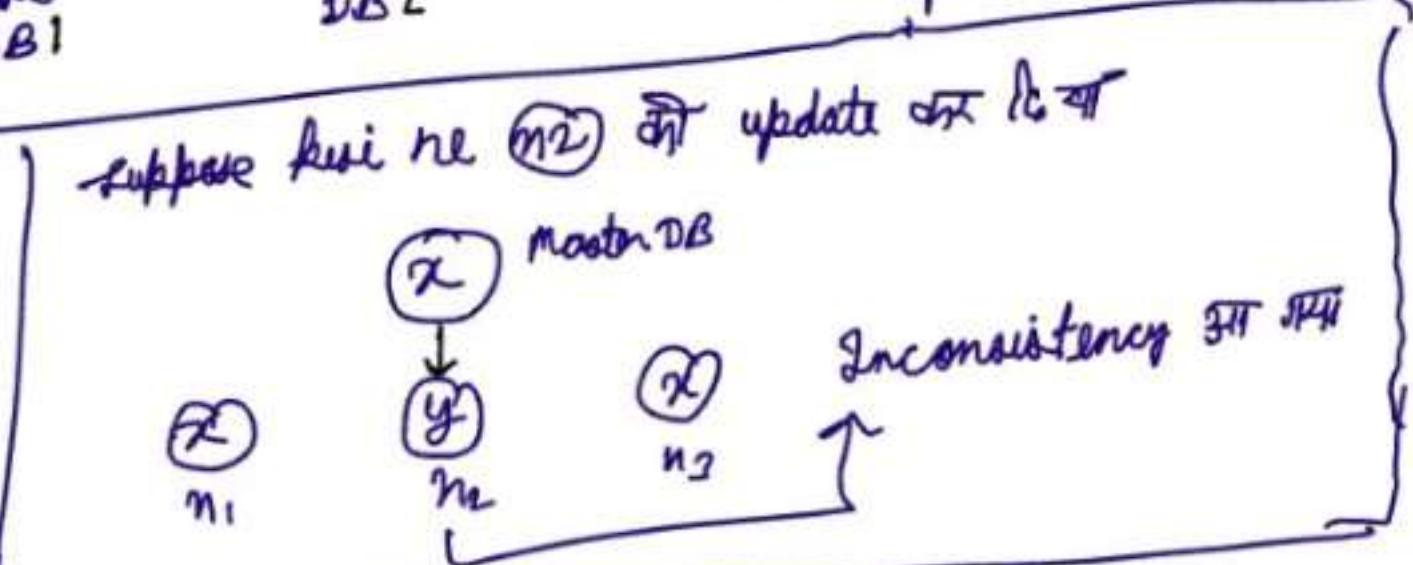
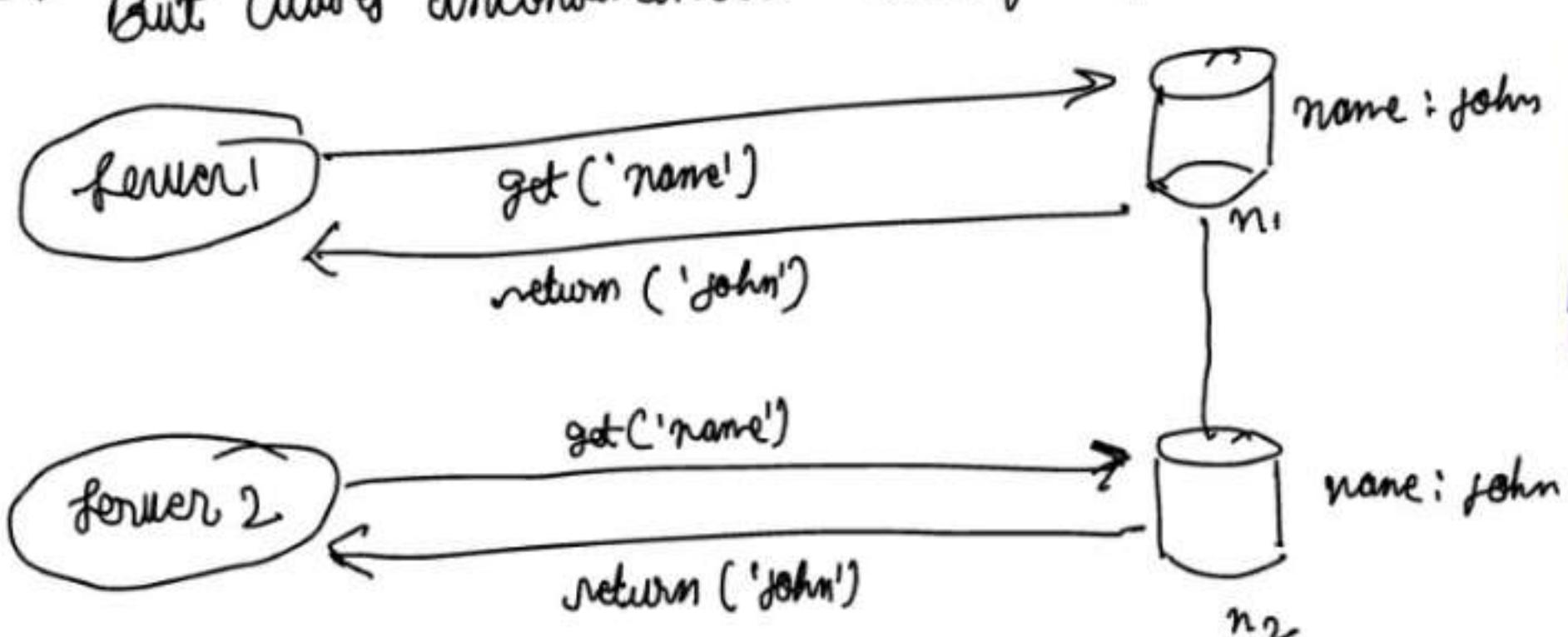
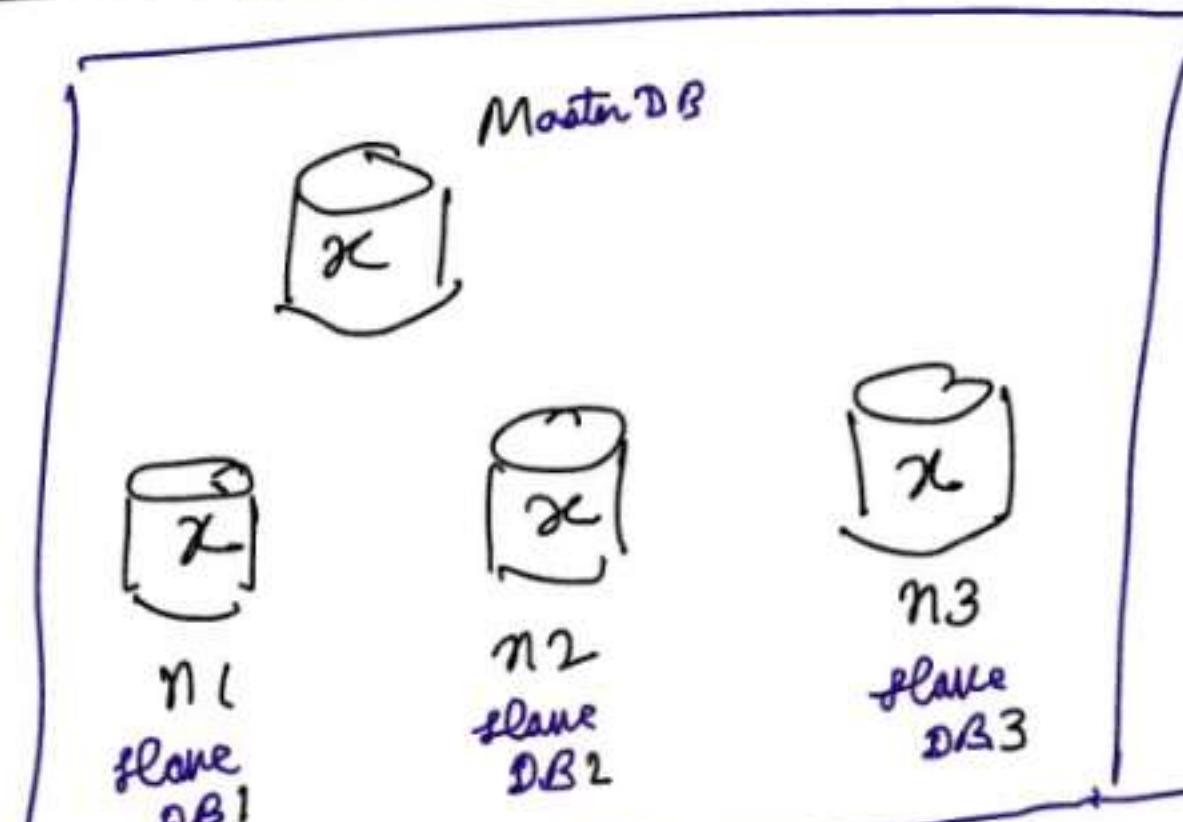


fig: n₁ & n₂ are replica nodes

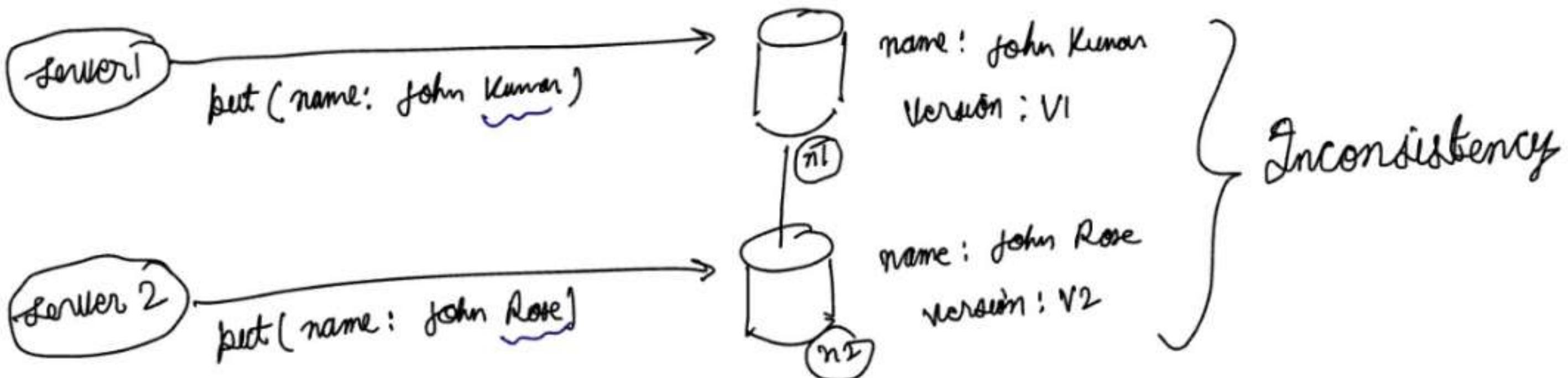


fig: n₁ & n₂ are replica nodes

→ conflict आएगा

→ somehow आरे conflict resolve कर पाय by knowing past versions.

Vector clock → To resolve conflict.

A pair [server, version]

associated with a data item.



S_x

$D_1([S_x, 1])$

S_x

$D_2([S_x, 2])$

S_x

$D_2([S_x, 2] [S_y, 1])$

S_y

→ conflict आएगा

$D_4[(S_x, 2) (S_z, 1)]$

S_z

$D_5[(S_x, 3) (S_y, 1) (S_z, 1)]$ → conflict resolved.

* Ancestor → $D[(S_0, 1) (S_1, 1)]$

Ancestor
(Parent)

→ $D[(S_0, 1) (S_1, 2)]$ → S_1 no update करेगा $(S_1, 1) \rightarrow (S_1, 2)$

जो जद में आया होगा $(S_1, 2)$

* Sibling Conflict →

- Disadvantage →
- Complexity → Client needs to implement conflict Resolution logic.
 - [Server, Version] pairs grow rapidly
 - Threshold set कर लेंगे and limit exceed कर तो oldest pair को remove कर लेंगे.
But यही है कि resolution की inefficiency भी सही है।
- But according to Dynamo Paper, Amazon has not faced such problem yet in production.

③ Design Key - Value store (Part 7) : (Handling Temporary failures) (Gossip protocol) (Hinted)

Handling Failure →

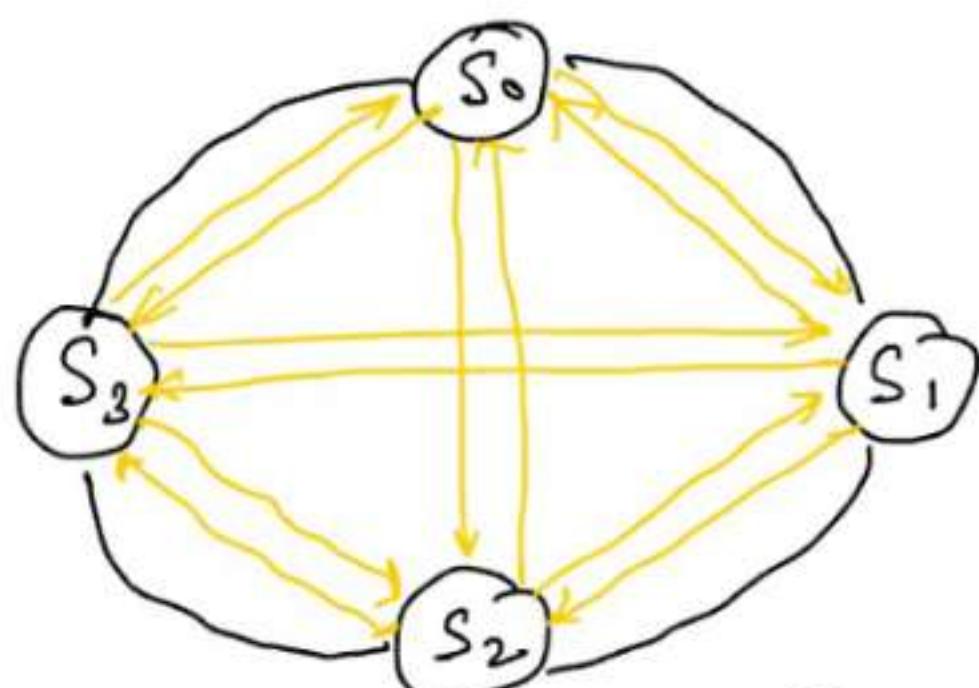
- Temporary
- Permanent.

→ Handling failure →
large system → failure is common

* Failure Detection → In a Distributed System



जबकि, शब्द बारे में पता हो.
km se km 2 log batayenge ki server down hai
main tabhi yakin karungi.



All to All Multicasting (सबको, सबके बीच में पता हो)

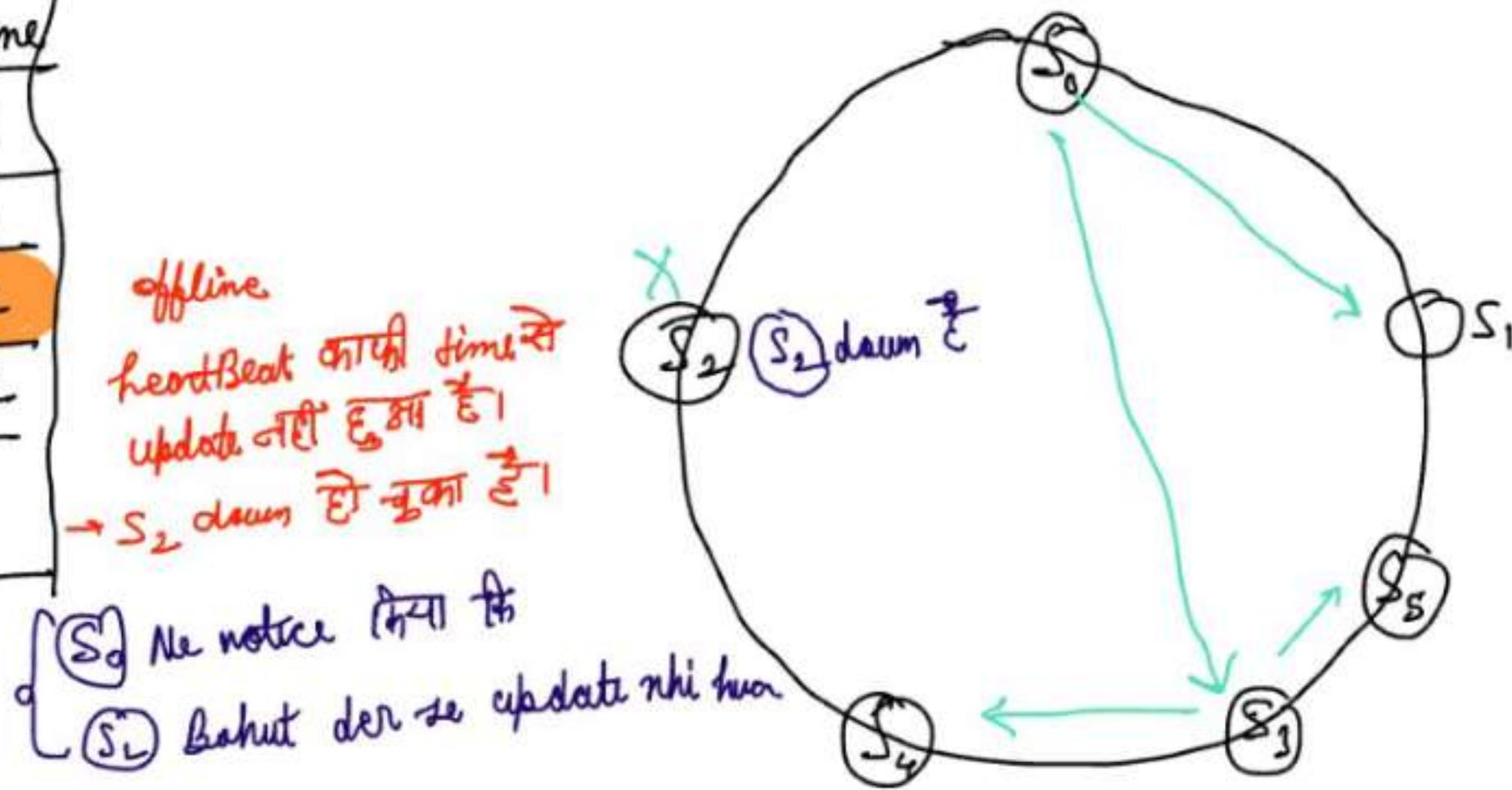
* Gossip Protocol →
Every Node has a membership list [MemberID, HeartBeat Counter, Time]

Member ID	HeartBeat Counter	Time
0	10 23 2	12:00:01
1	10 22 4	12:00:10
2	9 000	11:58:02
3	10 23 7	12:00:02
4	10 23 4	12:00:34

S0's Member list.

offline
heartbeat का time से update नहीं हुआ है।
→ S2 down हो चुका है।

(S0) Ne notice किया है।
(S1) Bahut der se update nahi hua.



Handling Failures \Rightarrow

① Temporary Failures \rightarrow

\hookrightarrow strict Quorum \rightarrow Read/writes blocked.

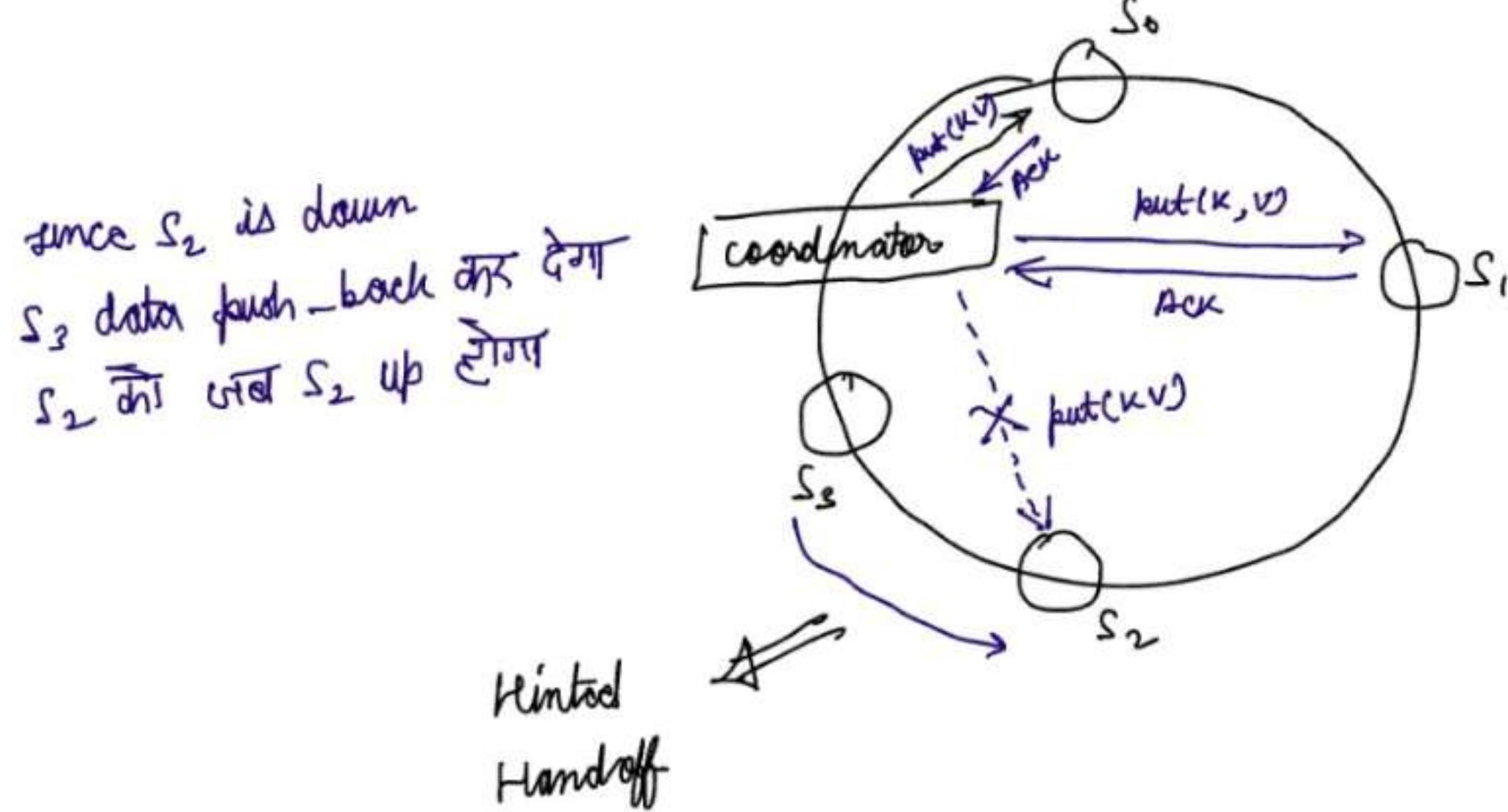
\hookrightarrow Sloppy Quorum \rightarrow Used to improve availability (जोड़ी से छूट है)

Hash Ring \hookrightarrow First W servers to write \rightarrow पहले W servers में write कर दे
First R servers to Read. \rightarrow पहले R servers में Read allow कर दे

When server is down, some other server will process request temporarily.

Hinted Handoff } Once it is up, changes will be pushed back to achieve consistency.

since S_2 is down
 S_3 data push-back कर देगा
 S_2 की ओर S_2 up होगा



③ Design Key-Value Store (Part 8) \rightarrow

Handling Permanent Failures | Anti-Entropy | Merkle Tree \Rightarrow

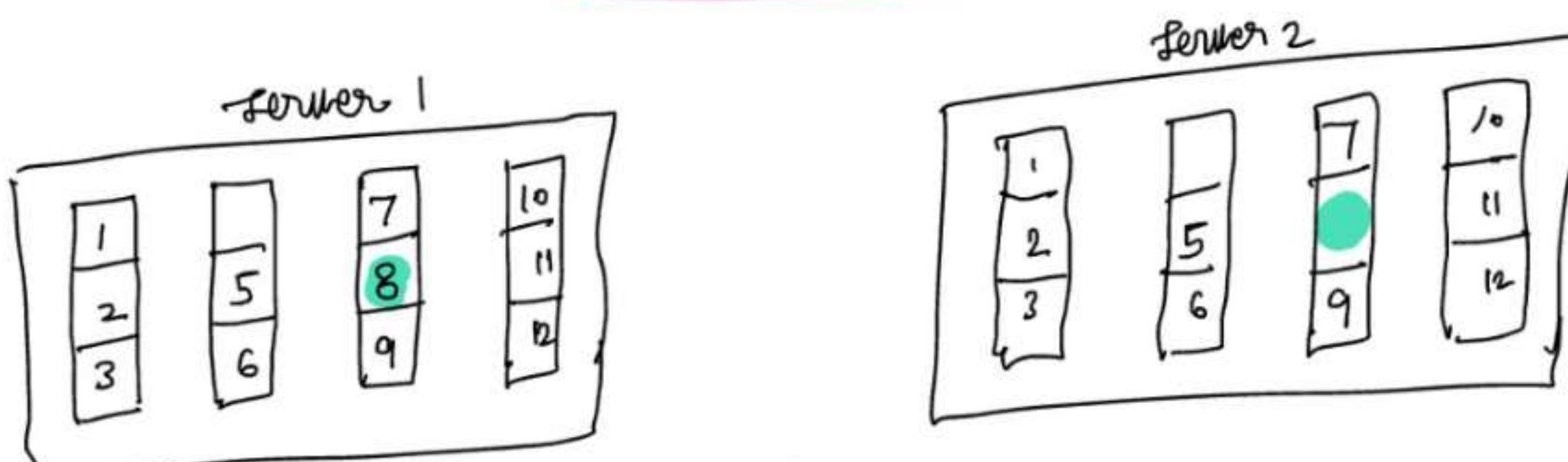
* Permanent Failure \Rightarrow

\rightarrow We implement "Anti Entropy" protocol to keep replica in sync.

Anti-Entropy \rightarrow

\hookrightarrow compare each piece of data on replicas and updating them.

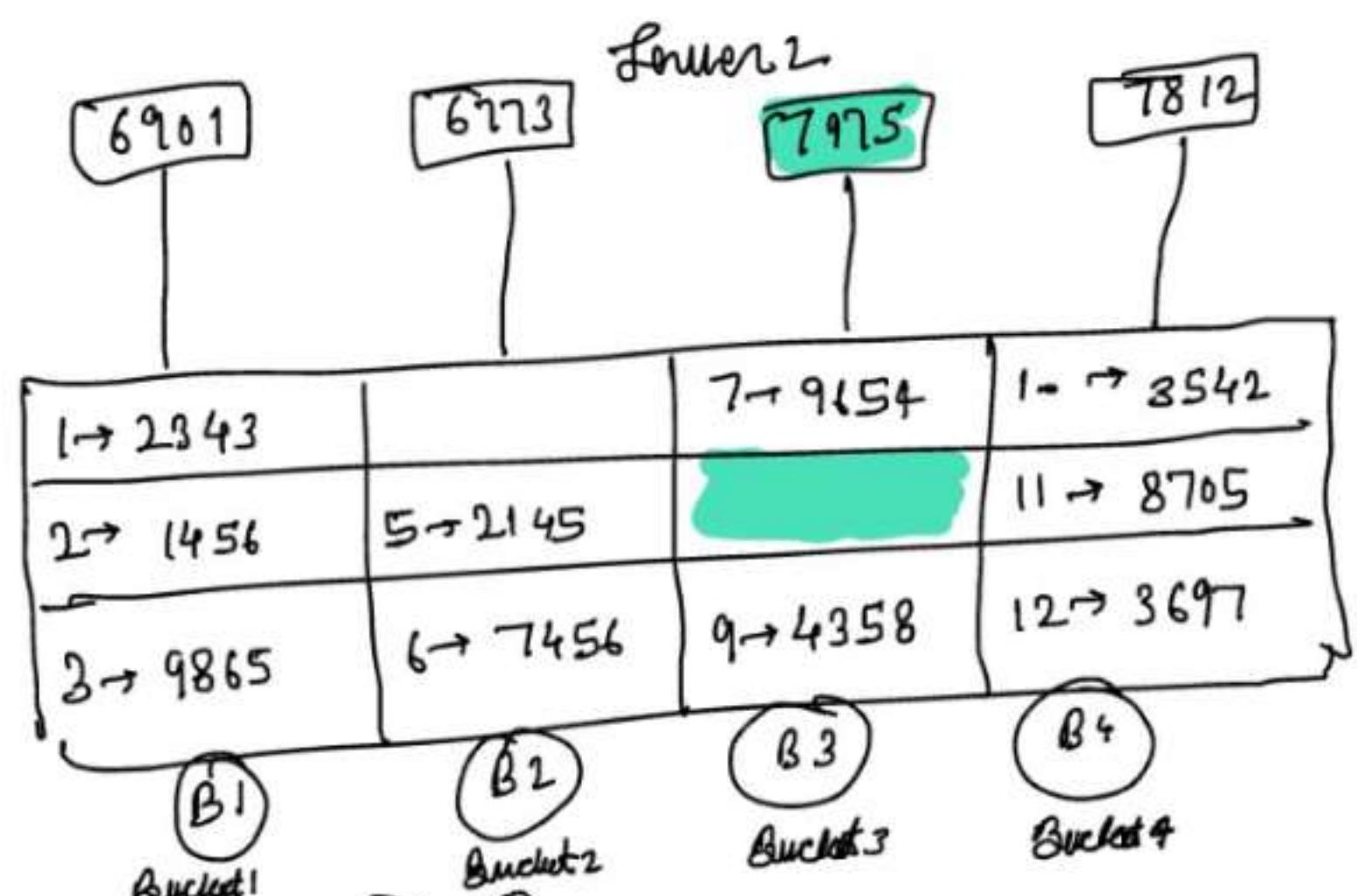
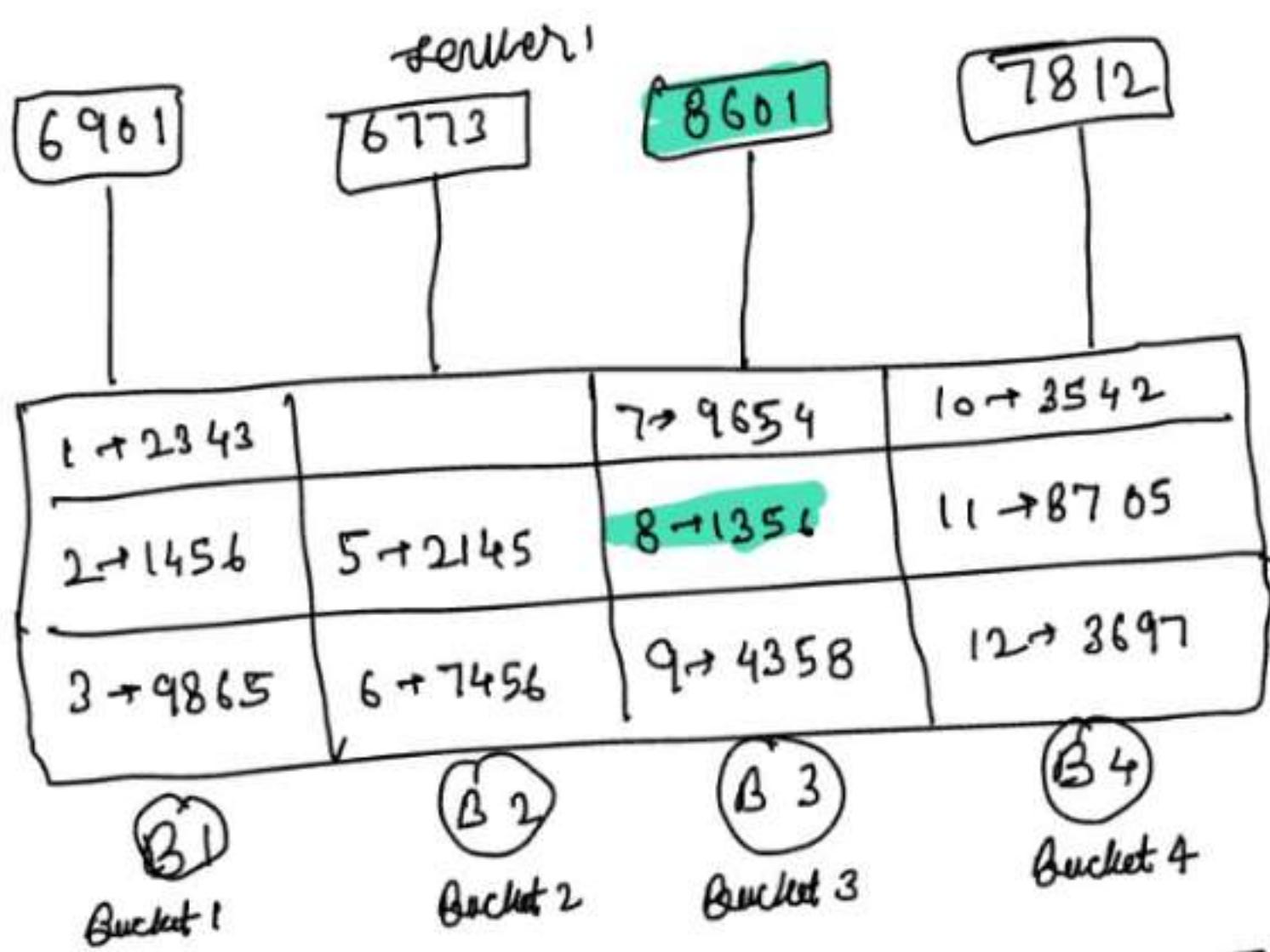
\hookrightarrow We use "**MERKLE TREE**"



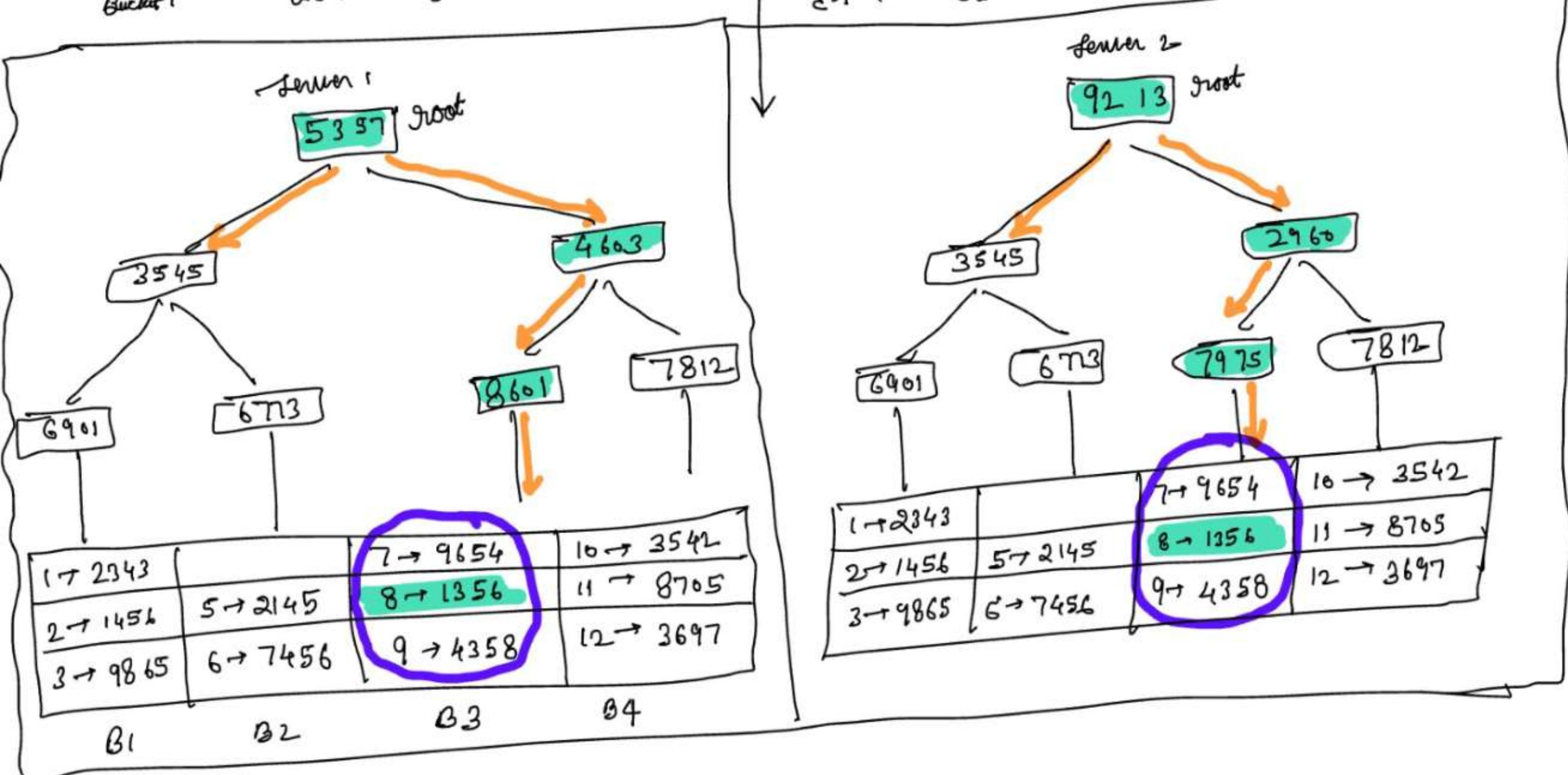
\downarrow सारे Keys को Hash कर देते हैं।

Server 1				Server 2			
1 → 2343		7 → 9654	10 → 3542	1 → 2343		7 → 9654	10 → 3542
2 → 1456	5 → 2145	8 → 1356	11 → 8705	2 → 1456	5 → 2145		11 → 8705
3 → 9865	6 → 7456	9 → 4358	12 → 3697	3 → 9865	6 → 7456	9 → 4358	12 → 3697

हर Bucket के लिये single Hash Node होगा



हम इस Tree बनाते हैं।



Benefit →

Amount of data to be synchronized \propto difference b/w two replicas

Handling Data Center Outage →

- Power Outage
- Network Outage
- Natural Disaster

We should replicate data across multiple data centers.

U.S.

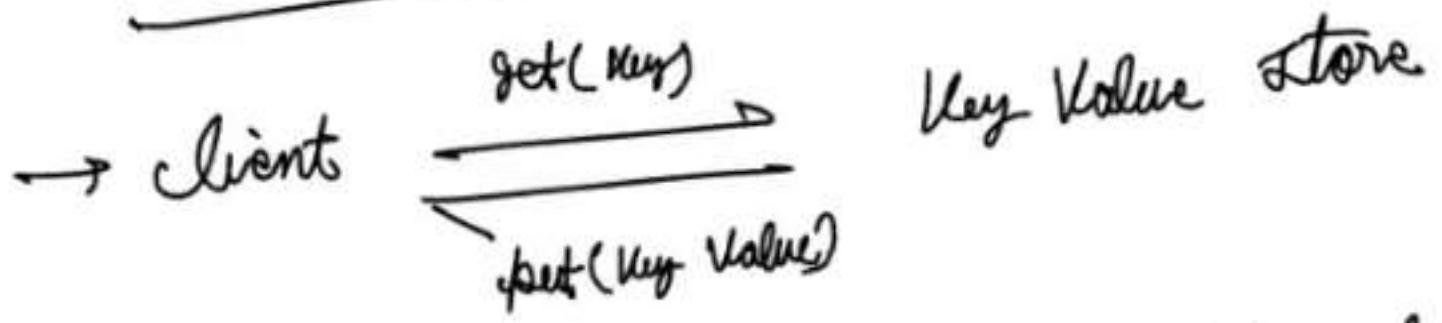
India

England

34

Design Key Value store (Part 9) →

System Architecture Diagram →



→ co-ordinator node : proxy b/w client & key-value store.

→ Nodes distributed on a ring → Consistent Hashing.

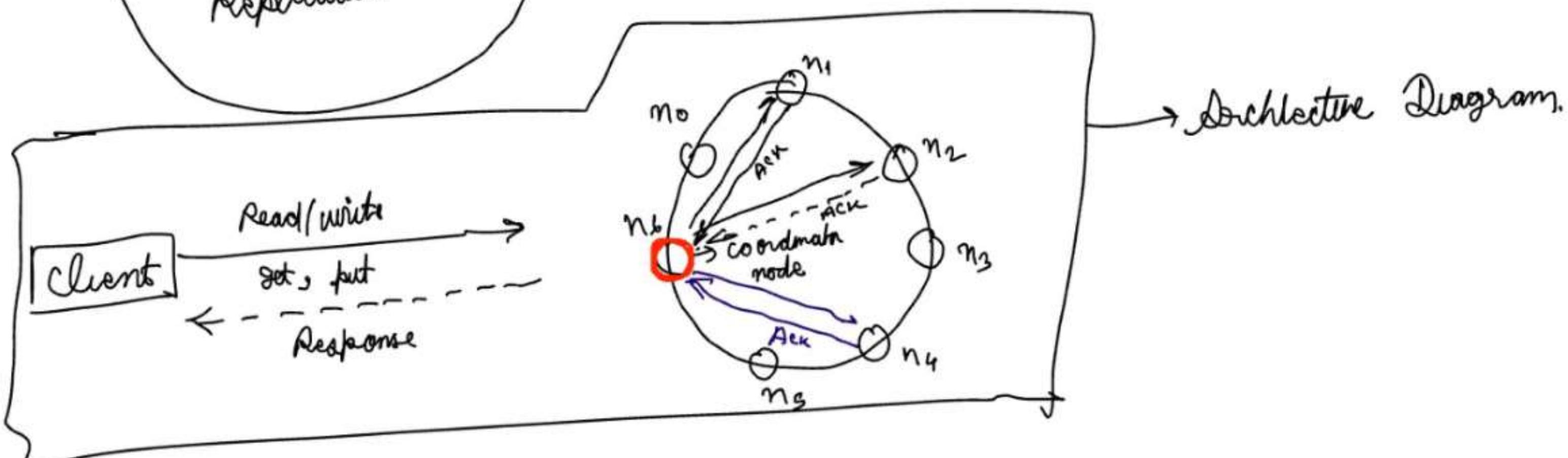
→ Add / Remove Node → Automatic.

→ Data Replicated in Replicas

→ No single point of failure - Multiple Nodes

Node responsibilities

- Client API
- Failure detection
- Failure Handling
- Conflict Resolution
- Replication



Arc. Diagram →

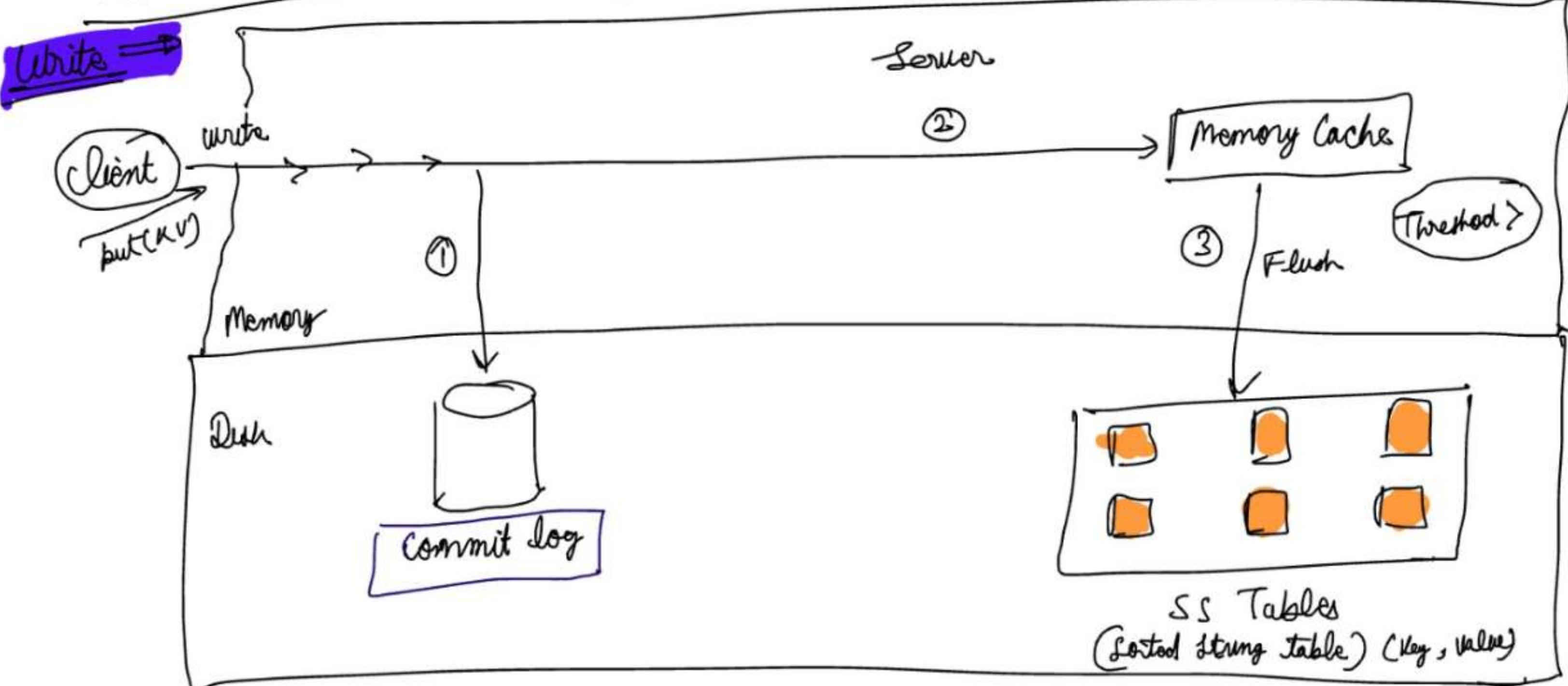
Read → - node (Read Path)

write → - node (Write Path)

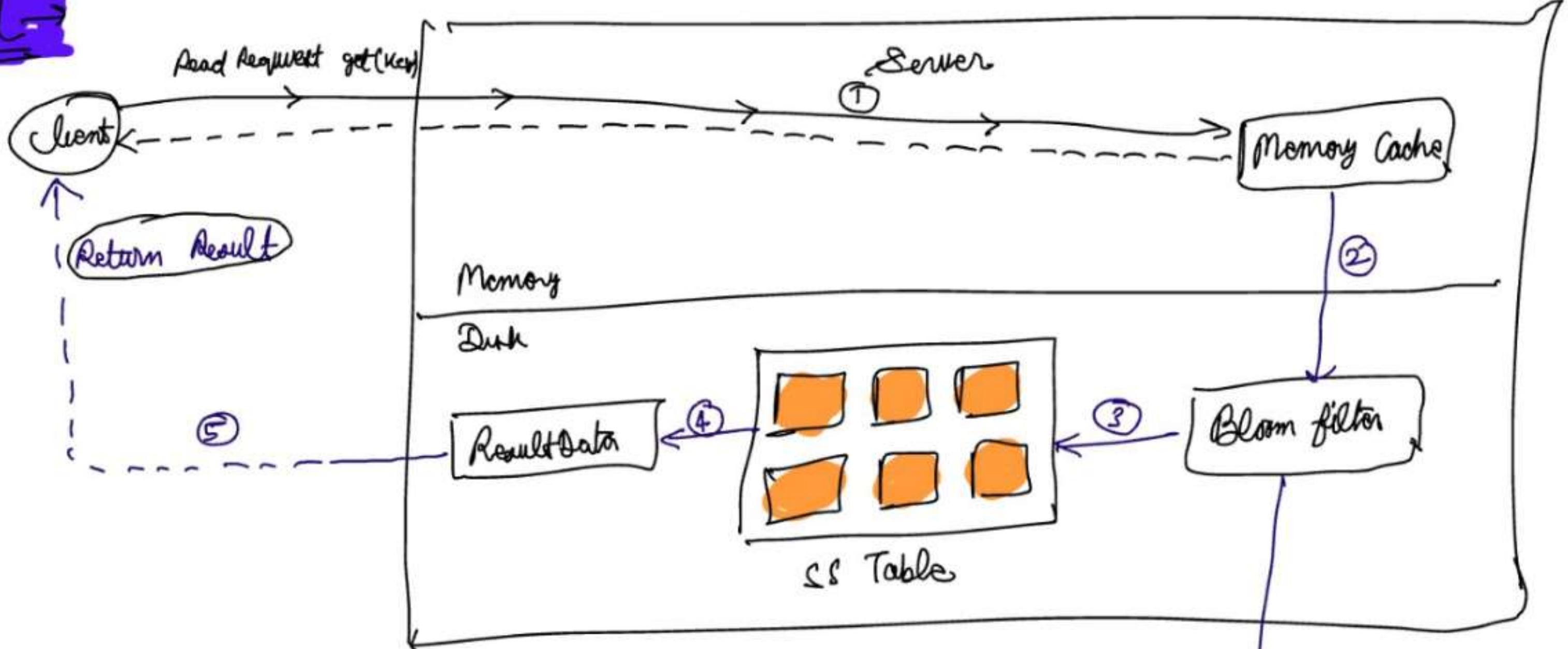
Real life example → Cassandra.

35 Design Key Value Store (Part 10) →

Read Path Write Path → Similar to Cassandra.



Read



इसकी वजह से पता - बल्टी हो जाती है कि कौन से SS table में डाटा है।

Goal/problems

Ability to store big data

Technique

Use consistent Hashing to spread the load across servers

High availability reads

Data replication Multi-data centre setup

High available writes

Versioning and conflict resolution with vector clocks

Data set partition

consistent Hashing

Incremental scalability

consistent Hashing

Heterogeneity

consistent Hashing

Tunable consistency

quorum consensus

Handling temporary failures

slappy quorum & hinted handoff

Handling permanent failures

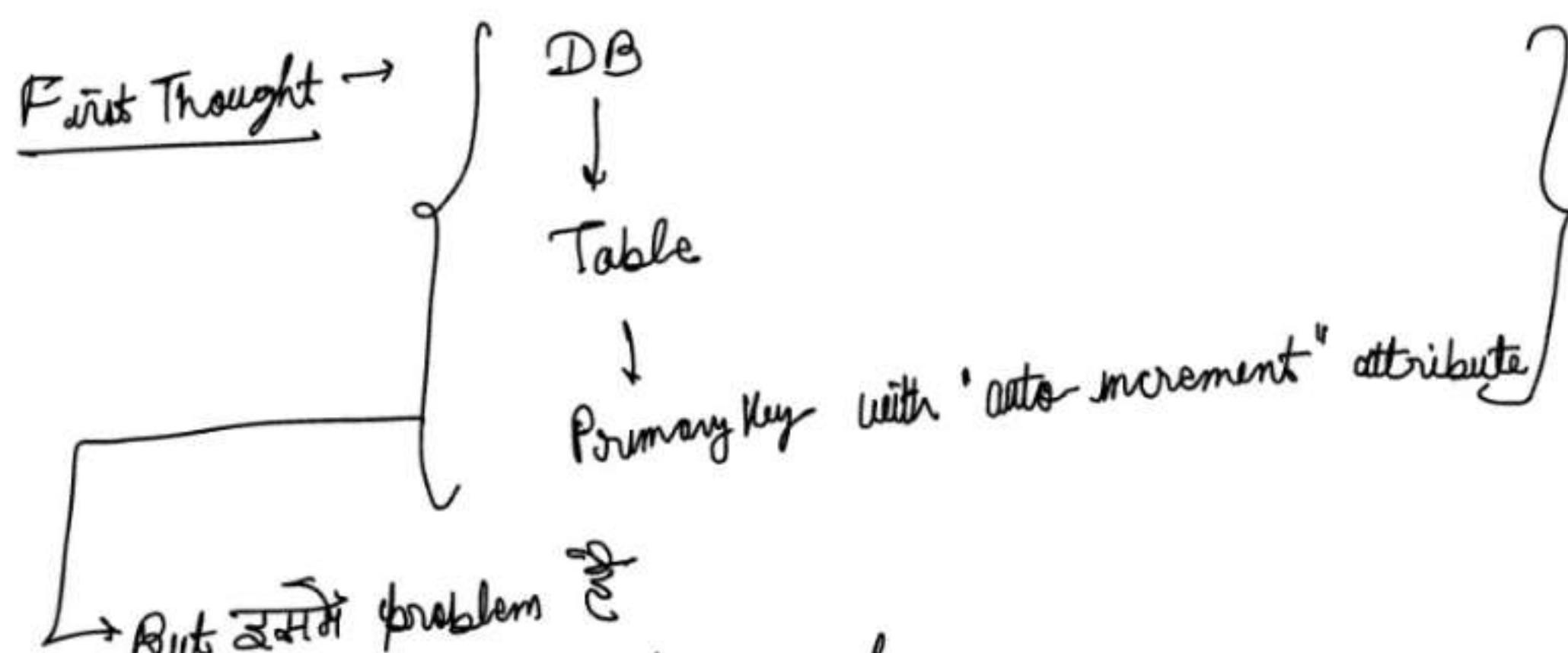
Merkle tree

Handling data center outage

Cross-data centre replication

36 Design a unique ID generator in distributed system (Part 1) →

Understand the problem →



example →

user_id
1 3 2 3 4 8 6 2 9 1
1 4 3 9 2 1 8 6 3 2
1 5 0 0 1 9 2 0 3 4

step ① Understand the problem and Establish design scope

Candidate → What are the characteristics of Unique IDs?

Interviewer → IDs must be unique and sortable.

Candidate → For each record, does ID increment by 1?

Interviewer → Increment ~~is~~ करेगा Time के अनुसार but not necessarily by 1.

$$\boxed{(\text{Evening ID}) > (\text{Morning ID})}$$

Candidate → Do IDs contain only numerical values?

Interviewer → Yes

Candidate → ID की length का requirement \rightarrow ?

Interviewer → ID should be fit in 64-bit.

Candidate → What is the scale of the system?

Interviewer → System should be able to generate 10 000 IDs per second

Requirement →

- ID must be unique
- Numerical values only
- 64 bit
- ordered by date.
- 10 000 unique IDs / sec

37) Design a unique ID generator in Distributed Systems (Part 2) Multi Master Replication →

Step ② Propose high level design and Get $\text{Guy} \rightarrow$

Multiple options can be used to generate unique IDs in distributed systems

→ (i) Multi Master Replication.

(ii) UUID (Universally unique Identifier)

(iii) Ticket Server

(iv) Twitter snowflake Approach (popular)

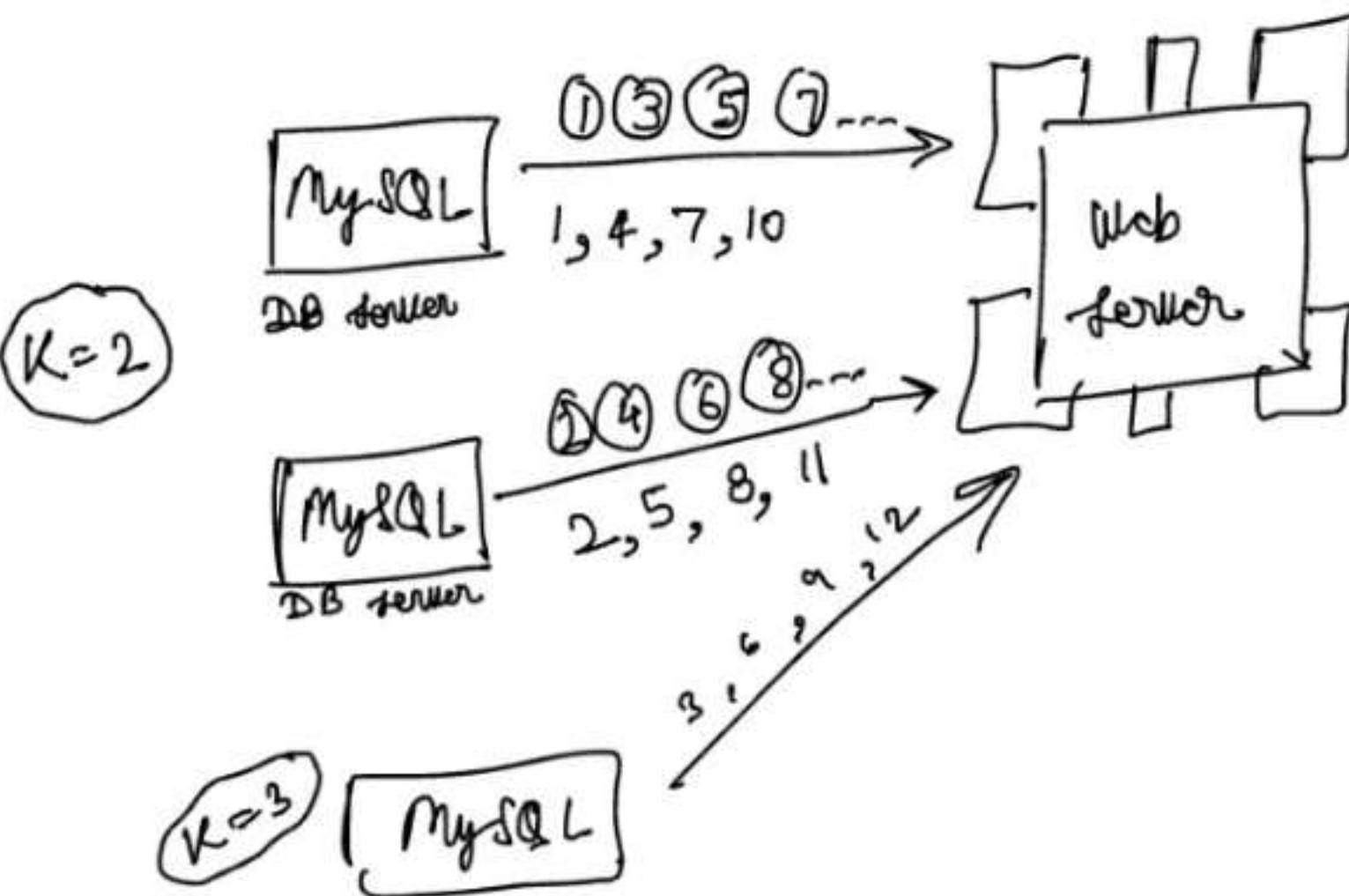
(i) Multi Master Replication →

→ It uses Database (auto-increment) feature.

→ ID को 1 से 0T increment करते, it increments by 'K'.

$$K = \text{no of DB servers in use}$$

→ Let's see a diagram with $K=2$.



This solves some scalability issue क्यों?

Major Drawbacks

→ How to scale with multiple Data centers

→ When a server is added OR Removed
क्या करें कैसे करें?

(will not scale well)



multiple data centre main scaling difficult to jata hai

③8) Design a Unique ID Generator in Distributed System (Part 3) : UUID →

UUID →

- Easy way to obtain Unique IDs.
- 128-bit number.
- Very low probability of getting collisions.

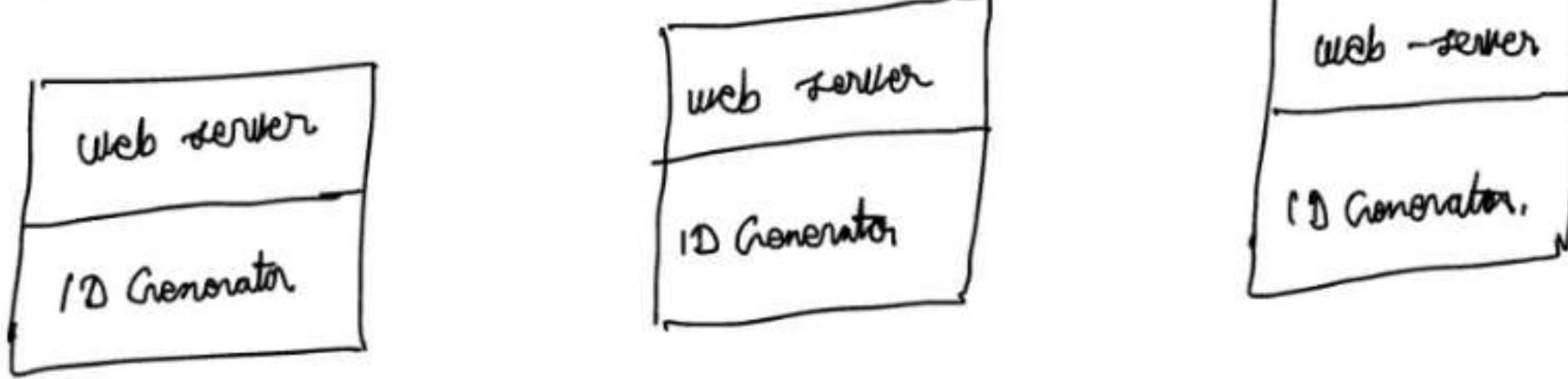
Fact : 1 billion UUIDs/second for 100 years probability of single collision will be 50%.

- Can be generated independently without co-ordinating with servers.

Eg →

123e4567-e89b-12d3-a456-426614174000 → numeric part →
non-numeric part →

Diagram →



Pros :-

- Generation is simple.
- No co-ordination within servers required.
- So, no synchronization issues.
- Easy to scale.

Cons :-

- 128 bits long. Our requirement was 64-bits.
- IDs could be non-numeric.

③9) Design a Unique ID Generator in Distributed System (part 6) →

Ticket Server →

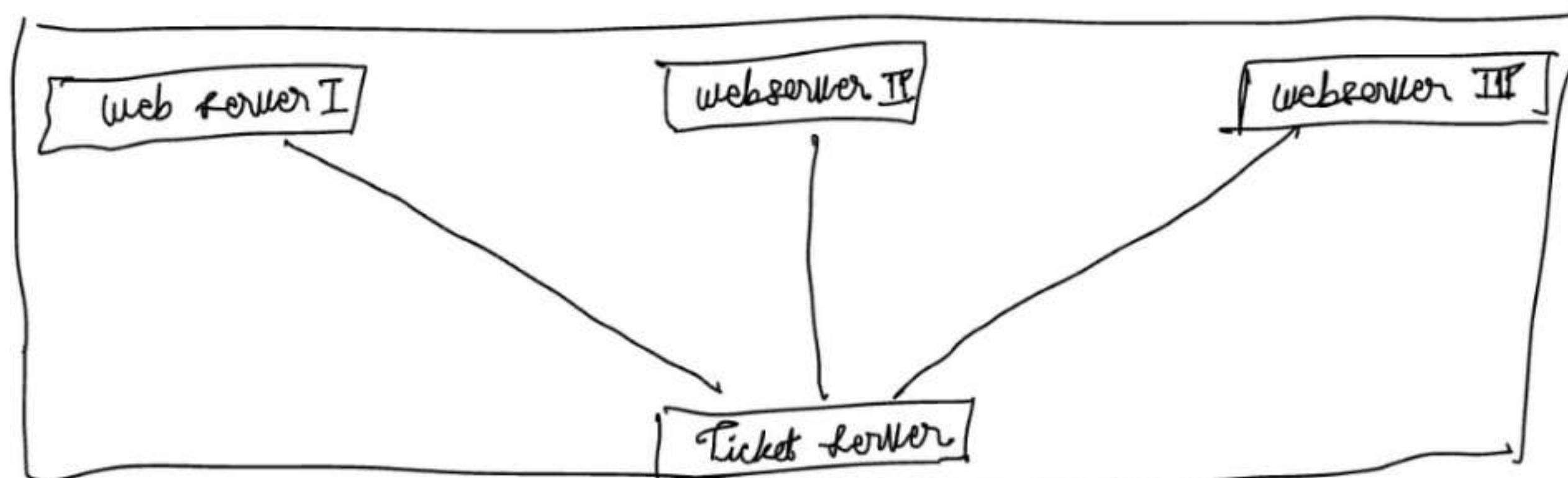
★ Ticket Server →

- Developed by "FLICKER"
- To generate Distributed Primary keys.

→ क्या है इसका Idea?

- Use a centralized auto increment feature in a single database server (Ticket Server)

Lets see this with a Diagram →



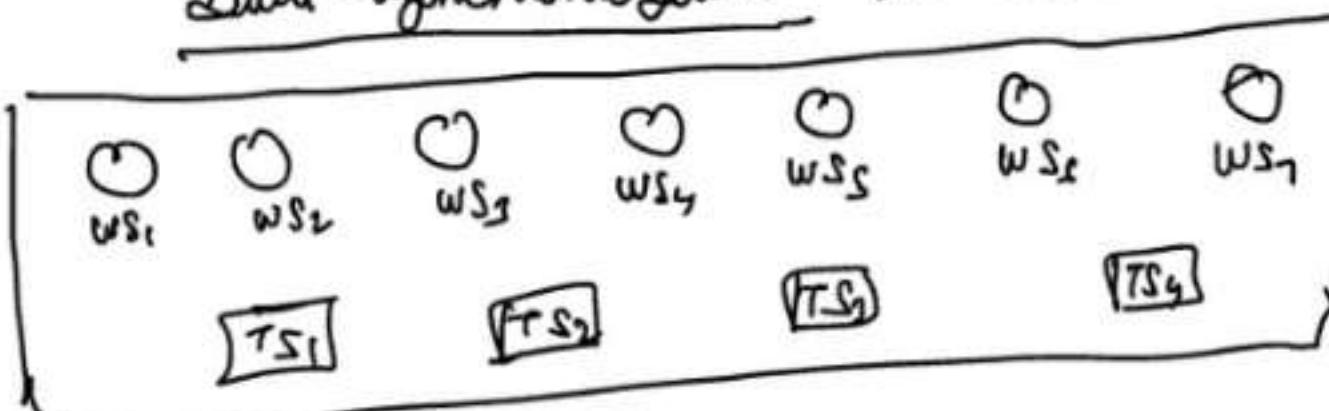
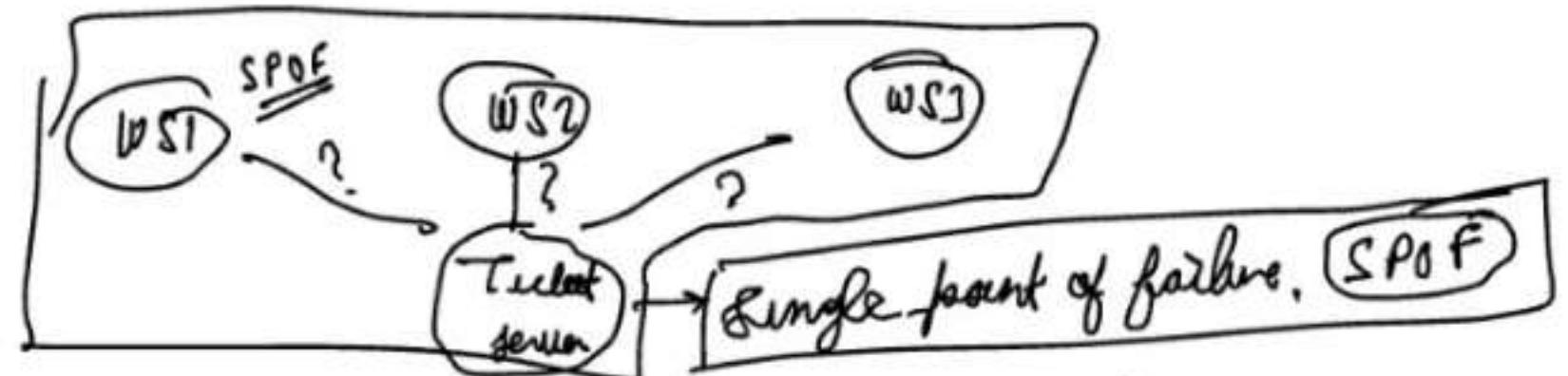
Pros →

- Numeric IDs
- Easy to Implement
- works for small to medium scale applications.

Cons :-

- single point of failure - Ticket Server
- To avoid SPOF (single point of failure)
 - we can setup multiple ticket server. But this will introduce challenges like

'Data synchronization' b/w Ticket Servers.



yaahan par SPOF se Bachne ke liye ek se zyada ticket server laga lia hai lekin iss wale approach me Data synchronization ka issue aayega.

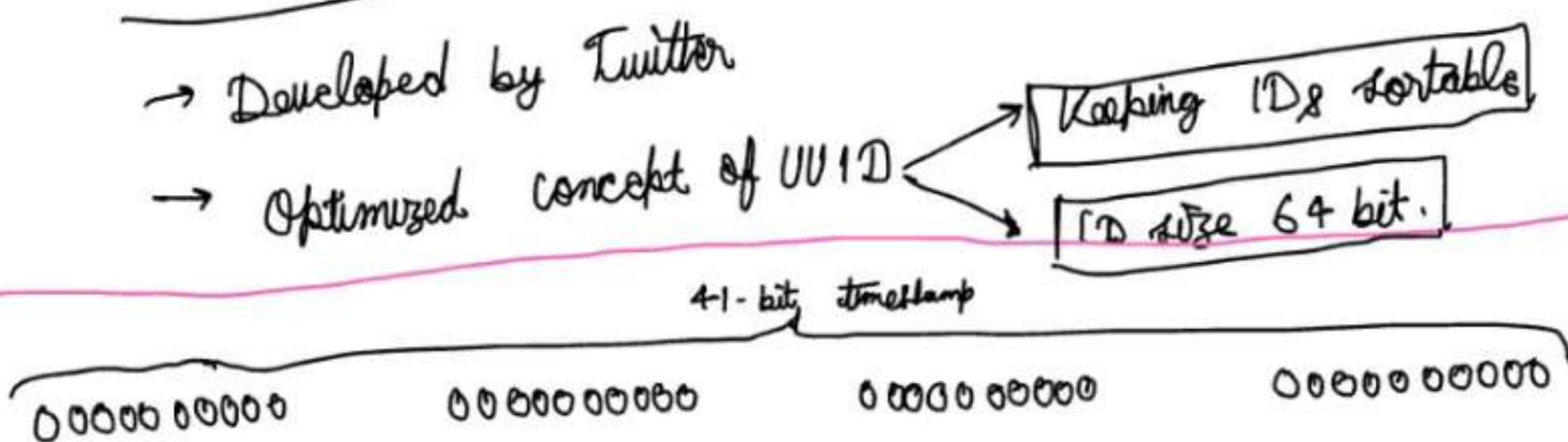
Q40 Design a Unique ID Generator In Distributed System →

Part ⑤: Twitter snowflake Approach →

* Twitter Snowflake Approach →

→ Developed by Twitter

→ Optimized concept of UUID



0000000000 00
12 bit sequence no
extra bit

→ 41 Bit → Epoch timestamp in milliseconds

→ 10 Bit → Machine ID - $(2^{10}) = 1024$ machines accommodated)

→ 12 Bit → Counter bits - Generated by local counter for every machine.

0 --- 10 --- 2¹² 0 --- 10 --- 2¹²

→ 1 Bit → Extra bit for future use. Default set to 0.

* Advantages of Twitterflake \Rightarrow

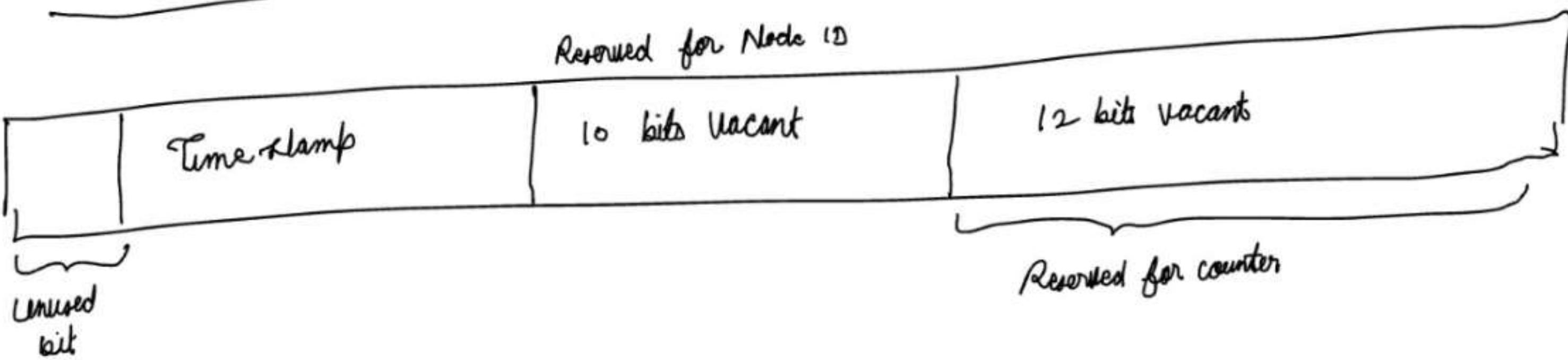
- Sortable by time \rightarrow They have 2 timestamp in the beginning.
- 64-bit \rightarrow less storage size requirement.
- Scalable \rightarrow 1024 machine
- Highly available \rightarrow 1024 machine & each machine generates 2^2 (4096) unique IDs each millisecond.

* Disadvantages of Twitter Snowflake \Rightarrow

- we have to maintain the machine IDs
- Generated IDs are not random like UUIDs
 - \hookrightarrow Predictable \rightarrow security issues.

(4) Design a Unique ID Generation in Distributed System \rightarrow

Part 6: Design Deep Dive \rightarrow

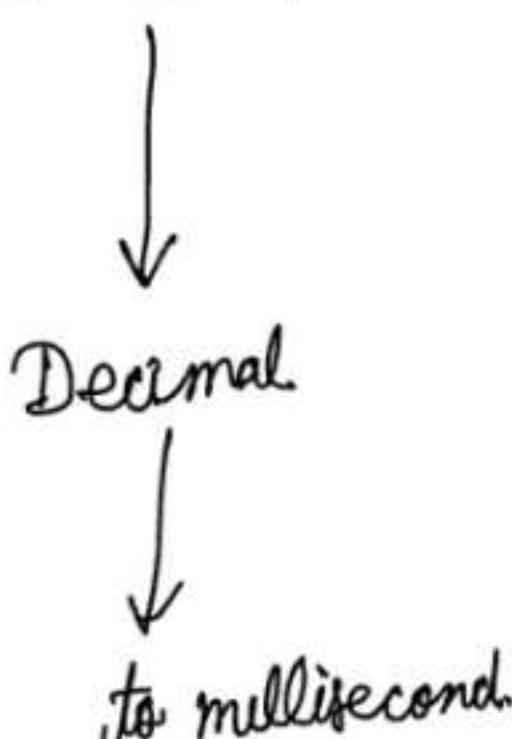


- * Node ID \rightarrow chosen at the startup time.
(10 bits) Generally fixed when the system is running. Use MAC Address — unique for every machine.

* Timestamp (41-bits)

\hookrightarrow As timestamp grows with time, IDs are sortable by time.

\hookrightarrow time stamp (binary) w.r.t epoch.



e.g. Aug 09 2022 17:34:31 UTC

* Maximum timestamp with 41 bits →

$$2^{41} - 1 = 2199023255551 \text{ (ms)}$$

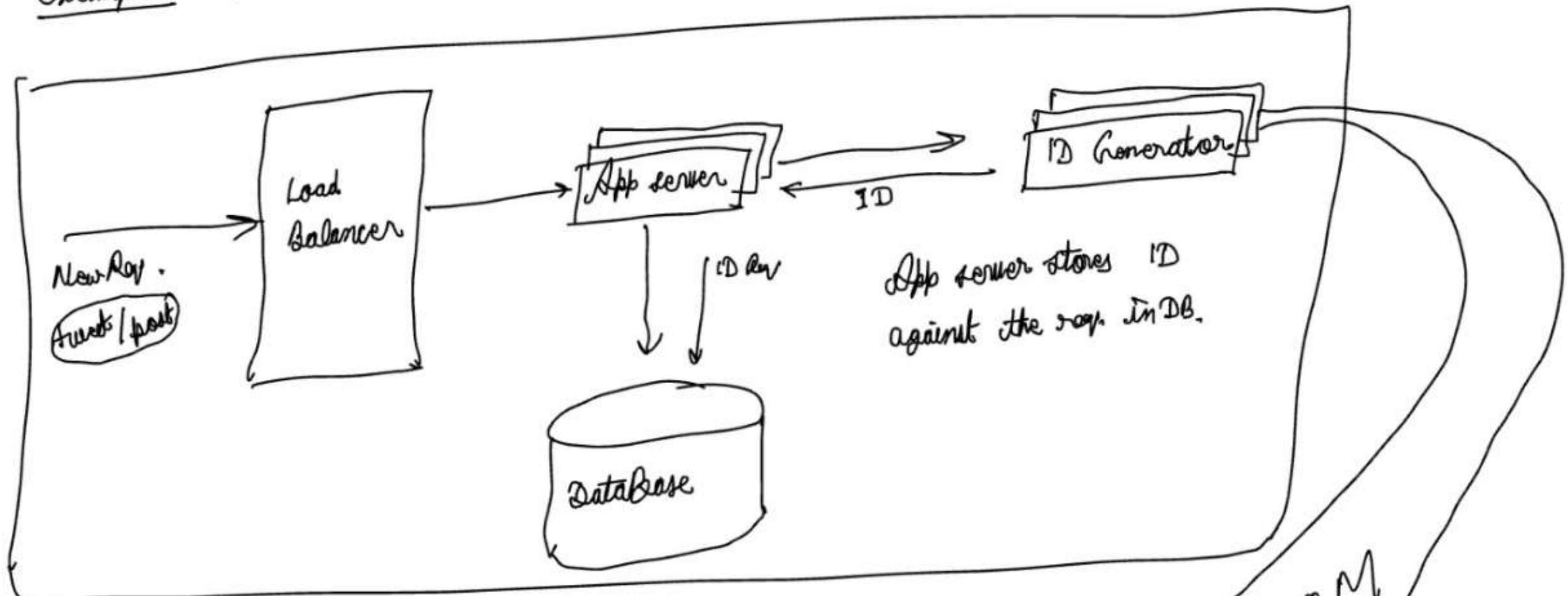
≈ 69 years

So after 69 years we will need new epoch time OR some other technique to migrate IDs.

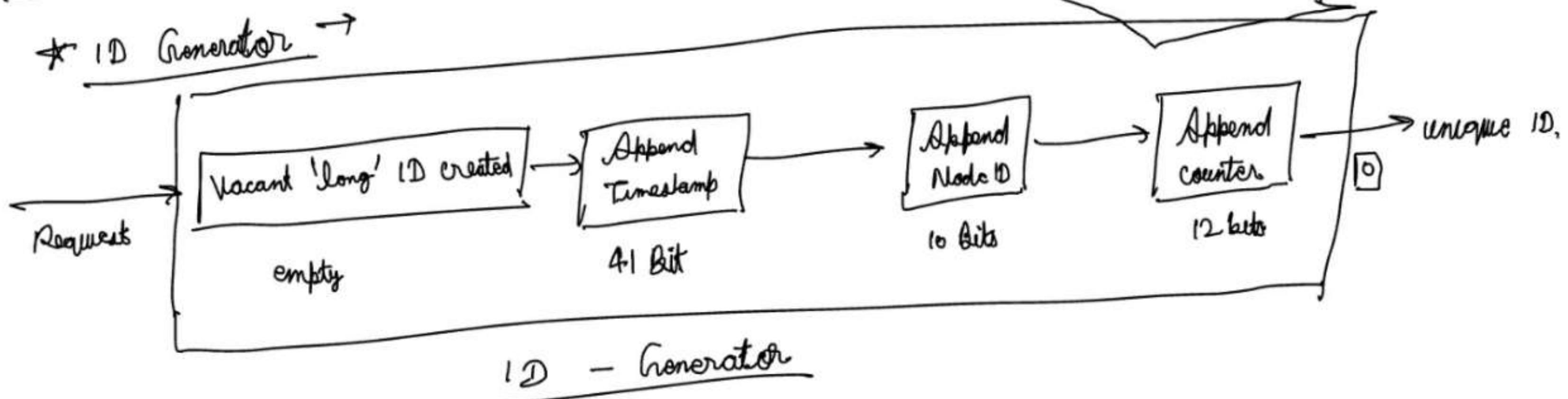
* Counter (12 Bits) → 12 Bits = $2^{12} = 4096$ combination
= 4096 IDs per millisecond

* System Diagram →

Example → A social networking application calls ID generator to assign unique IDs to each post.



ID generation की Zoom करके देखते हैं →



Q2 Design a Unique ID Generator in Distributed System (Part 2) (Wrap up) →

What we discussed so far →

- Multi Master Replication
- UUID
- Ticket Server
- Twitter Snowflake → settled for this
 - ↳ supports all our usecases
 - ↳ Scalable in a distributed Environment
- If you have extra time you can discuss few additional points!:-

① Clock synchronization → We assume ID generation servers has the same clock. Might not be true when a server is running on multiple cores.

"Network time protocol"

② Fewer Sequence numbers but more timestamp bits →

→ sounds like a good idea for long-term applications.

③ High availability → ID generator is a mission-critical system, so it must be highly available

Q3 Design a URL Shortner ⇒ (Hulu, Microsoft)

Part 0 : (Understand the problem & Establish Design Scope)

Step 0 → Understand the problem & Establish Design Scope →

Candidate → Can you give an example how URL shortner works?

Interviewer → Assume URL :

$\boxed{\text{https://www.interview_ds_algo.com/?v=course\&c=loggedin\&V=v3\&l=long \rightarrow \text{original URL}}$

Your service should create an alias with shorter length.

Example → $\boxed{\text{https://tinyurl.com/y7keocwxyj}}$ → shortened URL

↓
if you click this it should redirect to original URL.

Candidate → What is the traffic volume?

Interviewer → $\boxed{100 \text{ millions}}$ URLs are generated / day.

Candidate → How long is the shortened URL?

Interviewer → As short as possible

Candidate → what characters are allowed in shortened URL?

Interviewer → Can be a combination of
number: (0 - 9)
character: (a-z) (A-Z)

Candidate → Can shortened URLs be updated or deleted?

Interviewer → For simplicity, assume that they cannot be updated or deleted.

Summary and Basic use cases →

- URL shortening → Long URL → Short URL
- URL redirecting → Short URL → Redirects to Long URL address
- High availability → Scalable → Fault Tolerance

Next: Back of the envelope estimation →

44 Design a URL shortener →

Part ②: Back of the envelope estimation →

→ Write operation → 100 million URLs/day, generate 100.

→ Write operation per second: $100 \text{ million} / 24 / 3600 = 1160 \text{ URLs}$

→ Read operation → Assume that users read 10 times the write operations.

$$\text{So, Read operations/second} = 1160 * 10 = 11600 \text{ Read.}$$



→ कह तक चलेगी मेरी service ⇒

→ Assume → 10 years

→ So we must support

$$= 100 \text{ million} * 365 * 10 = 365 \text{ billions record.}$$

$\underbrace{\quad}_{\text{per day}}$ $\underbrace{\quad}_{\text{days in a year}}$ $\underbrace{\quad}_{\text{years}}$

Assume average URL length = 100

Let, 1 byte for each character

$$\text{So, total} = \underbrace{100}_{\text{URL length}} * 1 = 100 \text{ bytes per URL}$$

→ Storage requirement for 10 years →

We have to support 365 billion URLs over 10 years

→ per URL = 100 bytes.

→ Total storage → (365 billion) * (100 bytes) = 365 TB.

10 years ≈ 365 TB

④ Design a URL shortener (Part 3) →

Step ① → Proposed a High level Design and Get Buy-in !!! →

① API end points →

② URL Redirection →

③ URL shortening flows →

① API end points → Our service needs 2 API end points! :-

(i) → URL को short करने के लिये

POST : api/v1/data/shorten

↳ Request Param :

{ longURL : <long URL String> }

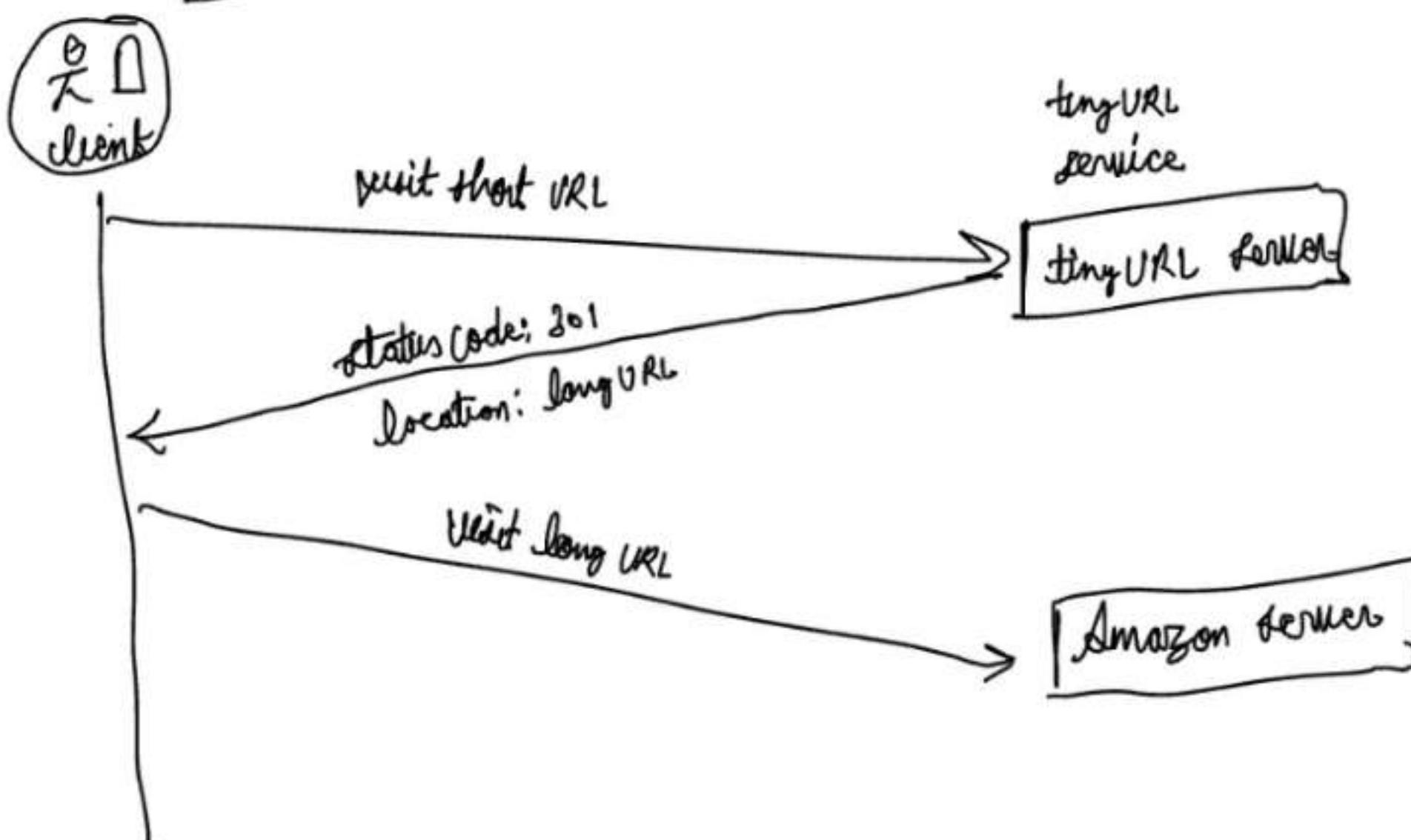
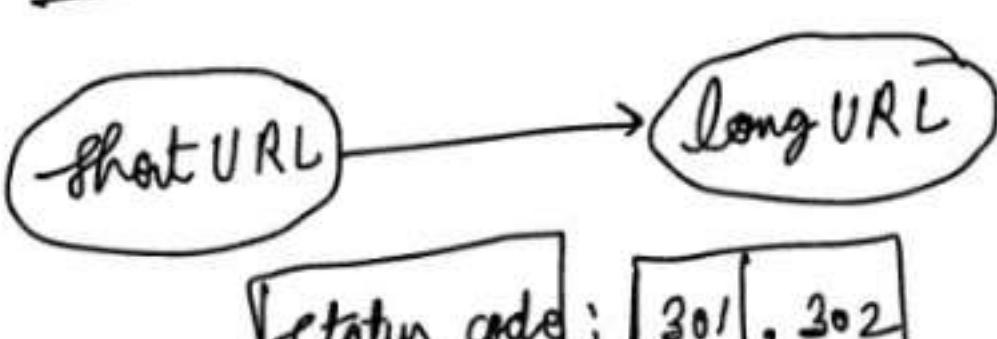
↳ Returns <short URL>

(ii) short URL को long URL में redirect करने के लिये →

GET : api/v1/shortURL

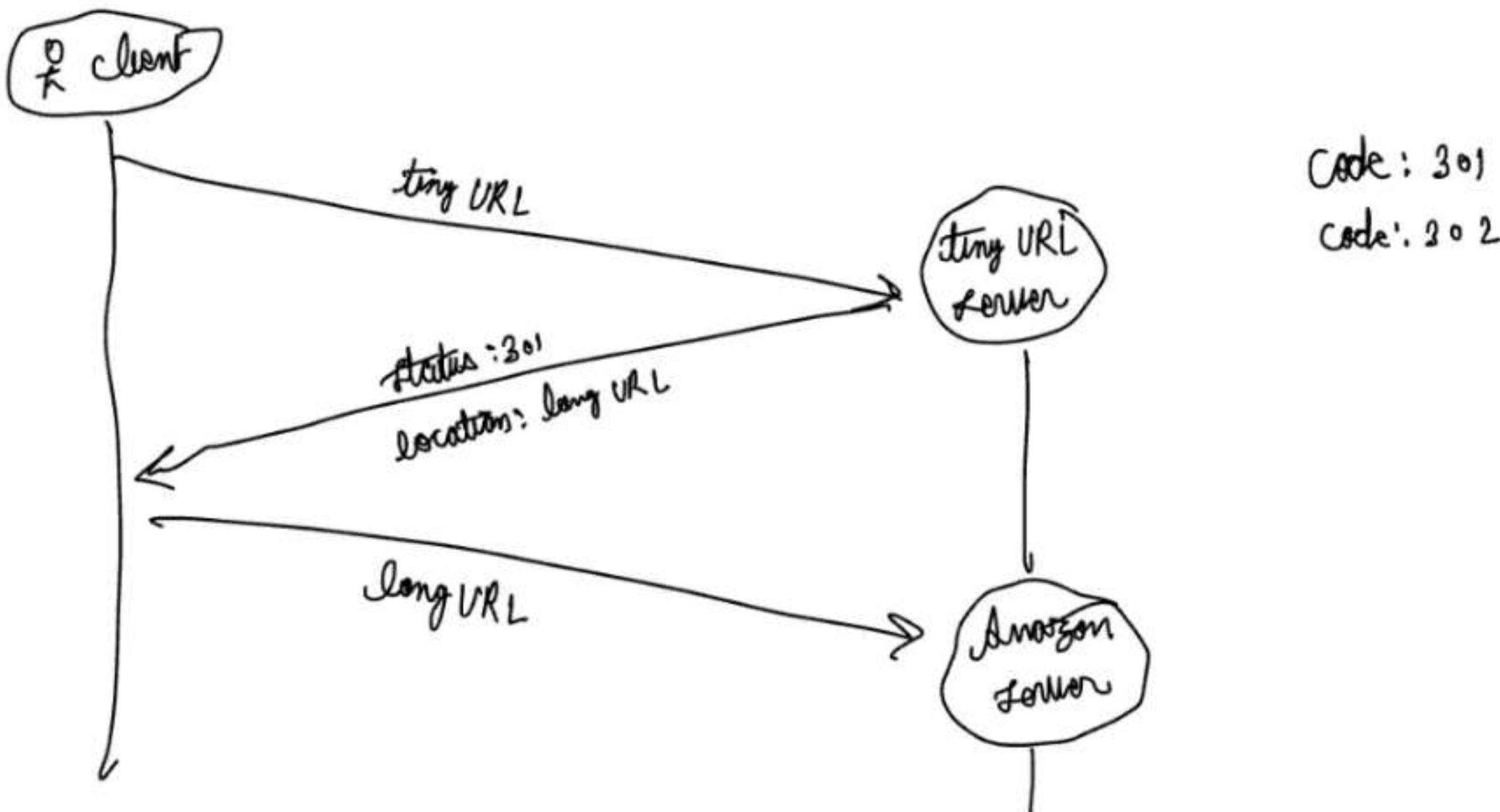
↳ Returns long URL for redirection →

② URL Redirection ⇒ कैसे काम करता है Browser में?



46 Design a URL shortener →

Part ④ : step ① 301 Redirect Vs 302 Redirect Status code ⇒



* 301 Redirect →

↳ Requested URL (tiny URL) is permanently moved to long URL.

↳ Browser caches the response

↳ फिर से click करेगा तो tinyURL में, तो tinyURL service के पास request नहीं जायेगा अब।
अब direct long URL से ही redirect होगा cache से

* 302 Redirect →

↳ Requested URL (tiny URL) is "Temporarily" moved to long URL.

↳ फिर से click करेगा तो tiny URL में तो tinyURL service के पास request जायेगा।

Pros and Cons of (301/302) redirect ⇒ कब कौन क्या use करें? ⇒

tinyURL service

301 → cache

302 → no cache

↓
Load कर दें

↑
तुम याहों की उसका request कर सको

↑
so 301 sounds good

* temporarily तुम्हें किसी और page से redirect कराना है

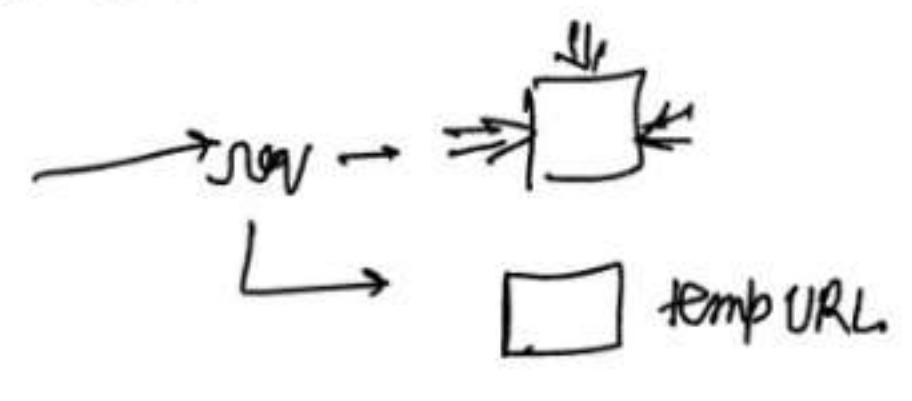
↳ suppose your website is under maintenance
↳ then u want to redirect other URL

→ also analytics के लिए ज़रूरी है।

↳ track click rate

↳ source of click .

↓
302 founds good

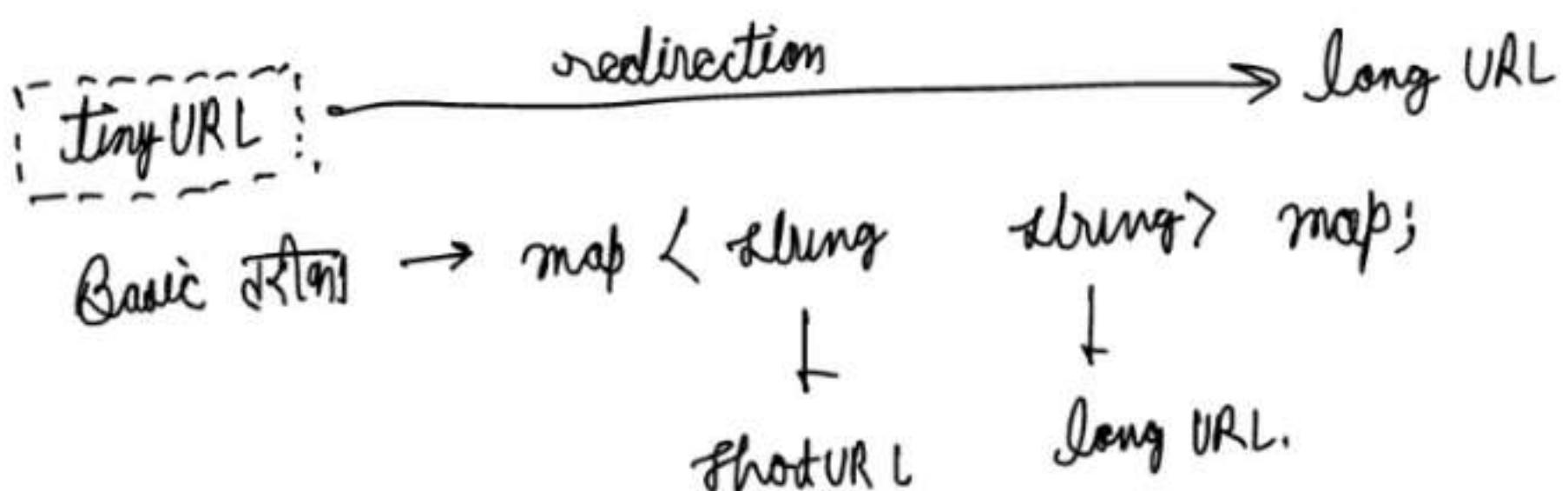


④ Design a URL shortener →

Part 5: Step ②: URL Shortening (Hash function) ⇒

③ URL Shortening Flows ⇒

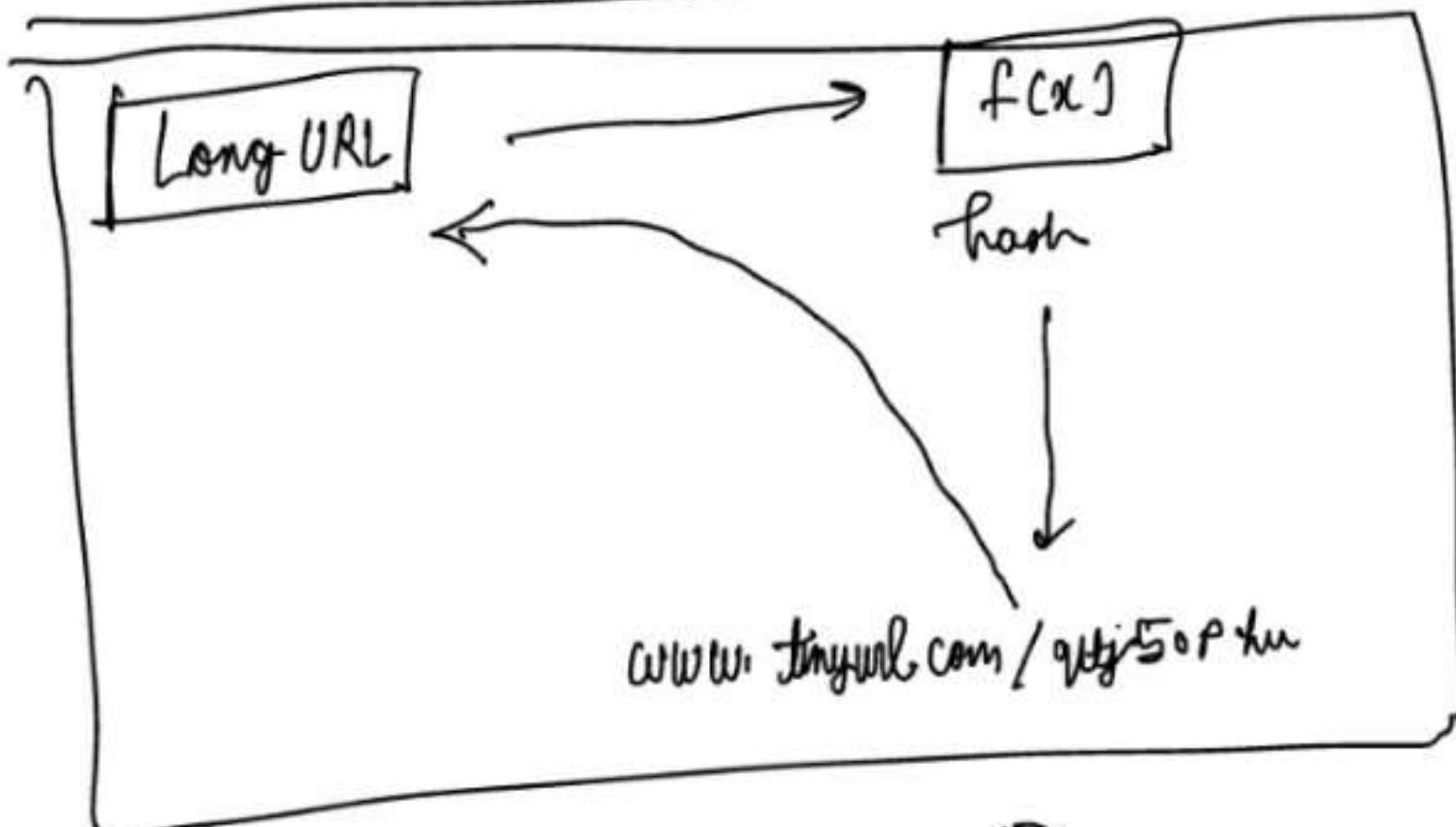
Before we start point ③ एक बात बताओ → point ① 'URL redirection' के बारे में !-



③ URL shortening flow →

example! - www.tinyurl.com/{hash,value}

We need a hash function $f(x)$ such that :-



→ Hash function कैसा होना चाहिए?

→ Each long URL must be hashed to one hash value.

→ Each hash value can be mapped back to long URL.

④ Design a URL shortener (Part 6) →

Step ③ Design Deep Dive (Data Model) →

→ URL shortening

→ URL redirection

Now, we will deep dive inside →

→ Data Model → Database कैसा लगायें?

→ Hash function →

→ URL shortening →

→ URL Redirection →

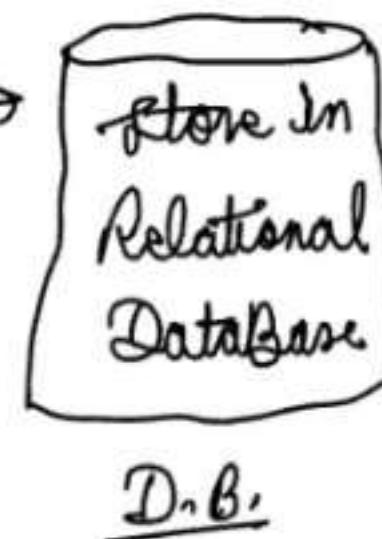
Data Model → Hash Table हो जाती है जो $\langle \text{Key}, \text{Value} \rangle$ प्रिय database रूपों चाहिए अब?

↳ Key, Value

↳ Hash Table → not feasible for real world systems . → क्यों ??? → मी^० Memory resources limited होगा Hash table one

So better option is →

< short URL, Long URL >



D. B.

Interviewer → Table Design कैसा रखेंगे?

Candidate → 3 columns रूपाने

`id` = primary key and auto-increment

PK

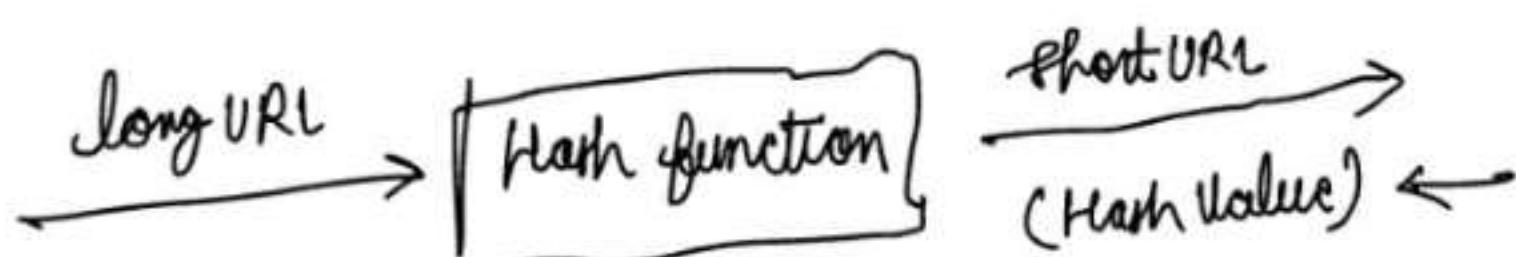
<u>Id</u>	<u>shortURL</u>	<u>longURL</u>

URL table

④9 Design a URL shortner (Part 1) →

Design Deep Dive (Hash function) \Rightarrow

* Hash function →



Interviewer → Hash Value की length क्या रखोगे?

$$\underline{\text{Candidate}} \rightarrow \text{Hash Value} = [0^{-9} \quad a^{-2} \quad A^{-2}]^T$$

$$\rightarrow \text{Total possibilities} = 10 + 26 + 26 = 62 \text{ possibilities}$$

Hash Value length = n

$$\text{combination} = \binom{62}{n}$$

— — — — — — — —
62 62 62 62 ————— 62

$$\text{Actual URL support} \approx 365 \text{ billion.} \Rightarrow 62^n \geq 365 \text{ billion, } n = ??$$

N	Maximal no of URLs
1	$62^1 = 62$
2	$62^2 = 3844$
3	$62^3 = 238,328$
4	$62^4 = 14,776,336$
5	$62^5 = 916,132,832$
6	$62^6 = 56,800,235,584$
7	$62^7 = 3,521,614,606,208 = \approx 3.5 \text{ Trillion}$

so $n=7$ is more than enough to store 365 billion URLs

Interviewer → What hash function you will be using?

Candidate → ① "Hash + Collision Resolution"

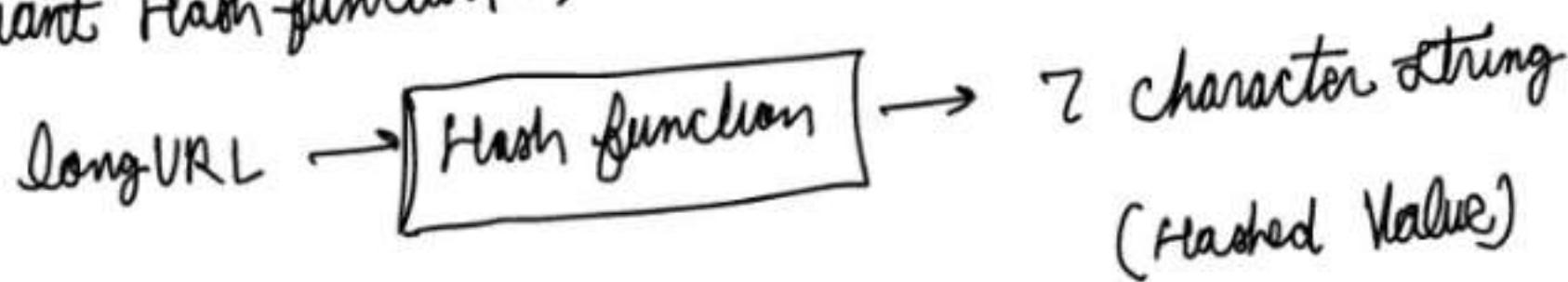
② "base 62 conversion"

⑤ Design a URL shortener (Part 8) →

step ③ Design Deep Dive (Hash + Collision Resolution) ⇒

* Hash + Collision Resolution →

→ We want Hash function →



→ straight forward solution → Use well known hash functions like →

- CRC 32
- MD5
- SHA-1

} These hash function available in market.

Long URL :- (https://en.wikipedia.org/wiki/System_design)

Hash function	Hash Value
CRC 32	5cb54054
MD 5	5a62509a84df9cc03fe1230b9df8b84e
SHA 1	0 ceae7916c06853901d9ccbefbfcaf4de57ed85b

But सबसे दोता (CRC 32) का शॉट Length > 7 है।

What to do ??

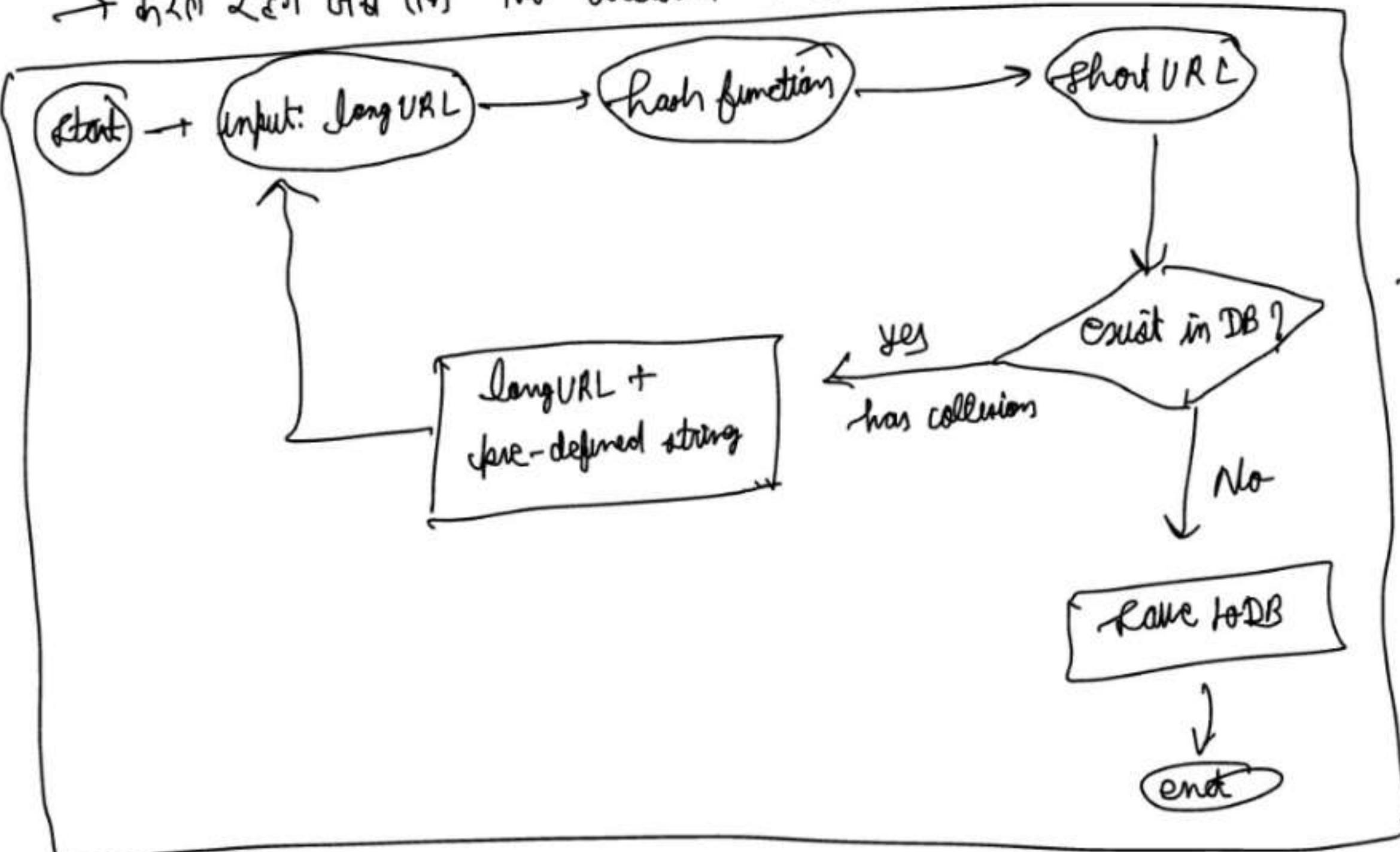
→ Starting की 7 character उत्तर नहीं (X)

e.g. → $\frac{5c654054}{5c654059}$ different है But 2 URL के 7 character same हैं
→ Hash collision.

मिल क्या करेंगे → ??

→ Hash collision हुआ हो long URL में एक 'Pre-defined string' append करें फिर से try करें।
→ करेंगे ऐसी तरफ "No collision" नहीं हो जाये।

(longURL + predefined string)



→ problem → डाटा-डाटा DB में memory होना।

solution → Bloom filter

Bloom filter ki help se quickly
bata kar skta hai ki koi value
present hai ki nahi किसी set में

51 Design a URL shortener

Part④ : Step③ Design Deep Dive (Base 62 conversion) →

* Base 62 Conversion →

Hashed Value

0-9	→ 10
a-z	→ 26
A-Z	→ 26
62	

Take Unique Number → example → auto-incremented id of a MySQL table.

e.g. → 11157 → convert this to Base 62. (0 1 2 ... 9 a b c ... z A B C ... Z)
(0 1 2 9 10 11 12 35 36 37 38 61)

62	11157	59	X ↑
62	179	55	T ↑
62	2	2	2
62	0		

62 characters.

https://tinyurl.com/2TX → shortURL

10	134	4
10	13	3
1.	1	1
	0	

$$(10^2 \times 1) + (10^1 \times 3) + (10^0 \times 4) \rightarrow \text{Decimal}$$

$$(62^2 \times 2) + (62^1 \times 15) + (62^0 \times 59) = 11157$$

why \rightarrow www.tinyurl.com/2Tx

why not? \rightarrow www.tinyurl.com/11157 \rightarrow (11157 > 2Tx \rightarrow it's too small; ultra archive)

Hash + collision Resolution

Base 62 Conversion.

Fixed short URL length

The short URL length is not fixed - It goes up with ID.

It does not need a unique ID generator

This option depends on a unique ID generator.

Collision is possible & must be resolved

Collision is impossible because ID is unique.

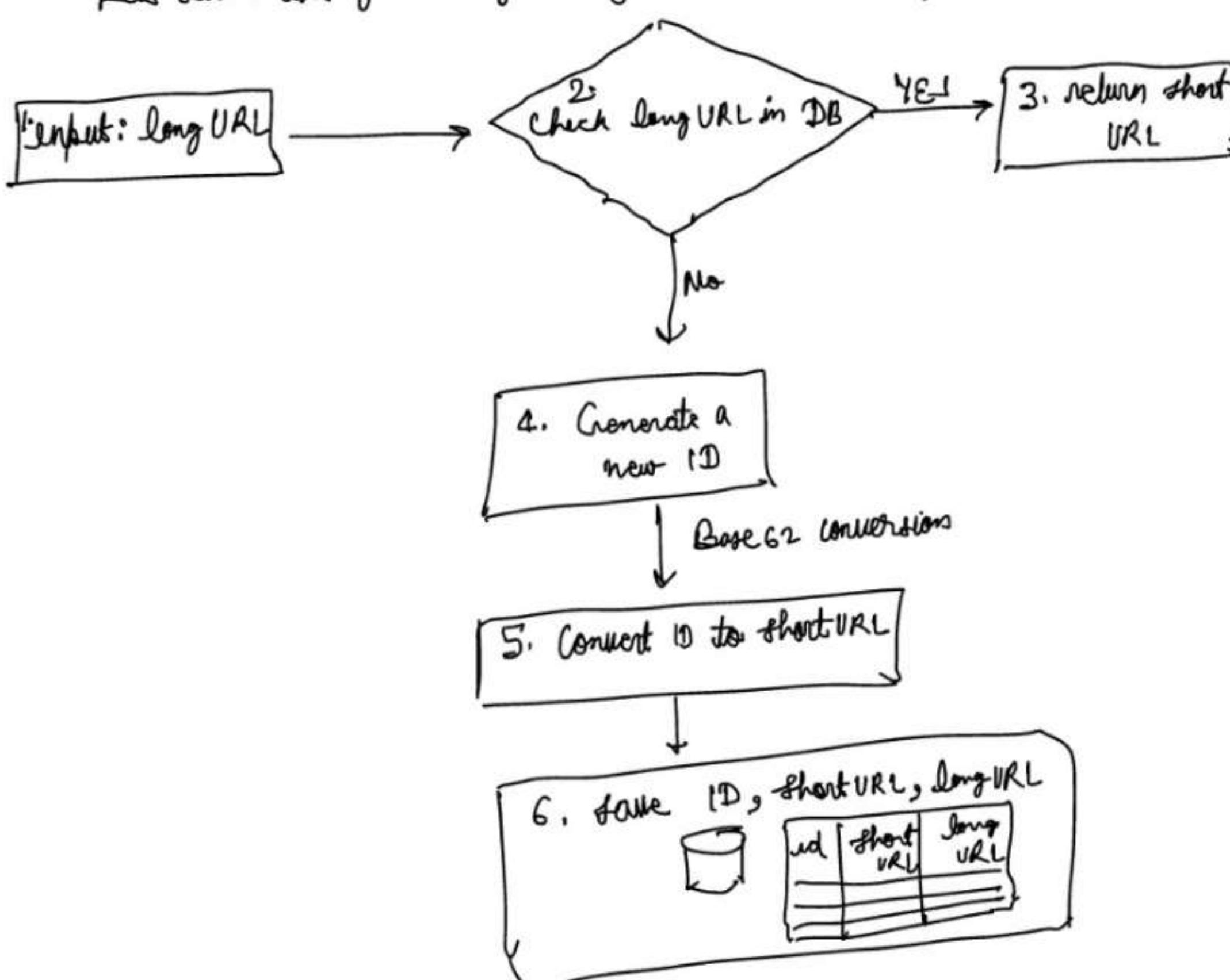
It is impossible to figure out the next available short URL because it does not depend on ID,

It is easy to figure out the next available short URL if ID increments by 1 for a new entry. This can be security concern.

52 Design a URL shortener (Part 10) \Rightarrow

Step ③ Design Deep Dive (URL shortening flow) \rightarrow

Let's draw the flow diagram for URL shortening flow \rightarrow



example \rightarrow

long url: - https://en.wikipedia.org/wiki/System_design

unique ID generator returns ID: 20138653

Convert ID to short URL using: "mten2f3"

Base62 conversion

ShortURL: www.denyal.com/mten253

id	short URL	long URL
20138653	mten253	https://en.wikipedia.org/wiki/System_design

⑤ Design a URL shortner \Rightarrow

Part ⑪ Step 3 & 4: Design Deep Dive and WRAP UP

(URL Redirecting Deep dive) \rightarrow

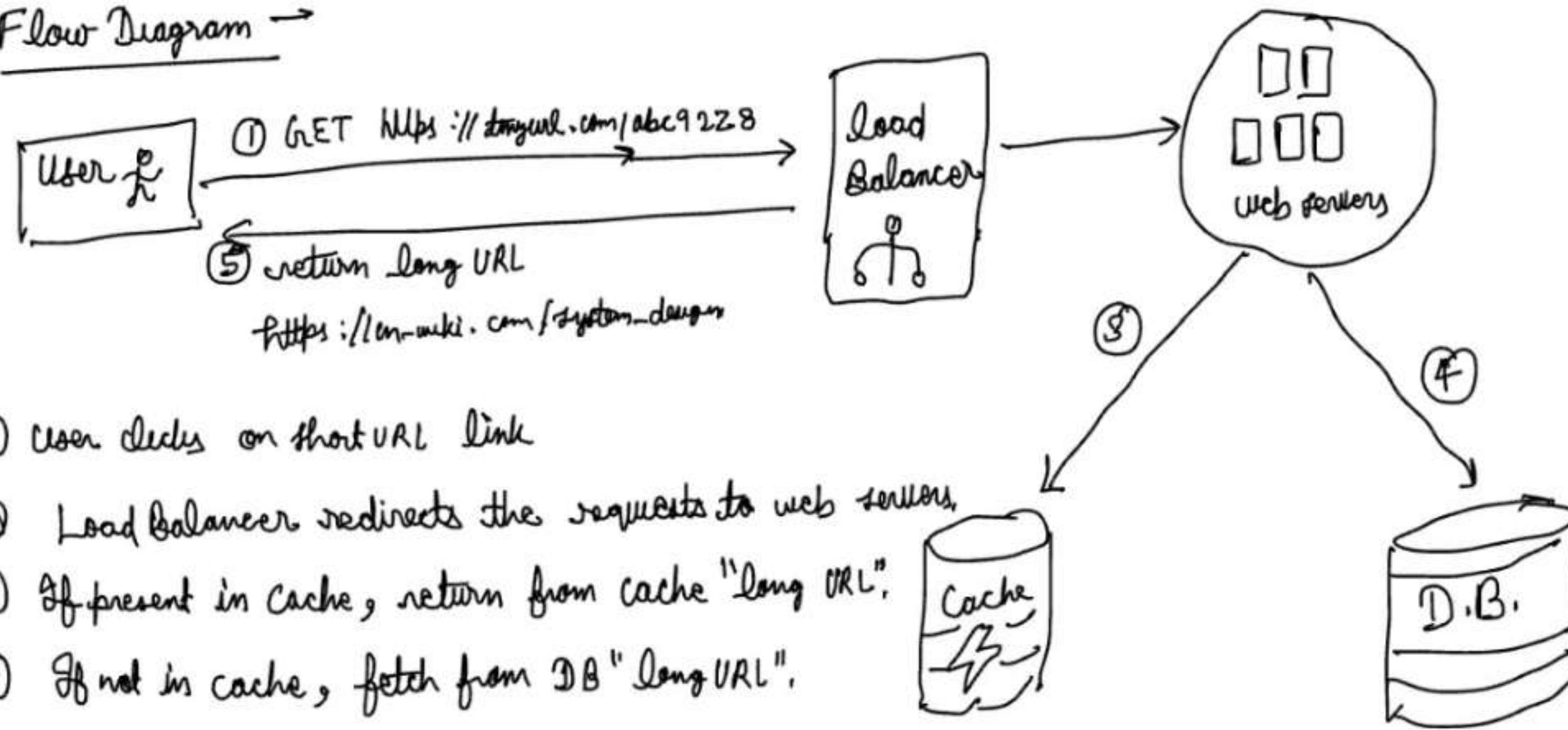
* Generally \rightarrow More Reads than Writes

$(\text{No of Read Requests}) > (\text{No of write Request})$

$\langle \text{Short URL}, \text{Long URL} \rangle \rightarrow \text{store in Cache Memory}$

\hookrightarrow Improved performance.

Flow Diagram \rightarrow



① User clicks on short URL link

② Load Balancer redirects the requests to web servers,

③ If present in Cache, return from cache "long URL".

④ If not in cache, fetch from DB "long URL".

⑤ Return long URL to user.

Step ④ WRAP UP \Rightarrow

If there is extra time in the end you can discuss \rightarrow

\rightarrow Rate limiter \rightarrow malicious user \rightarrow Bait & switch req they can't hit

\rightarrow Web server scaling \rightarrow Easy to scale up and scale down \downarrow

\rightarrow Data Base scaling \rightarrow DataBase Replication / Sharding

54) Design a Web Crawler \Rightarrow (Google Amazon Meta) \rightarrow

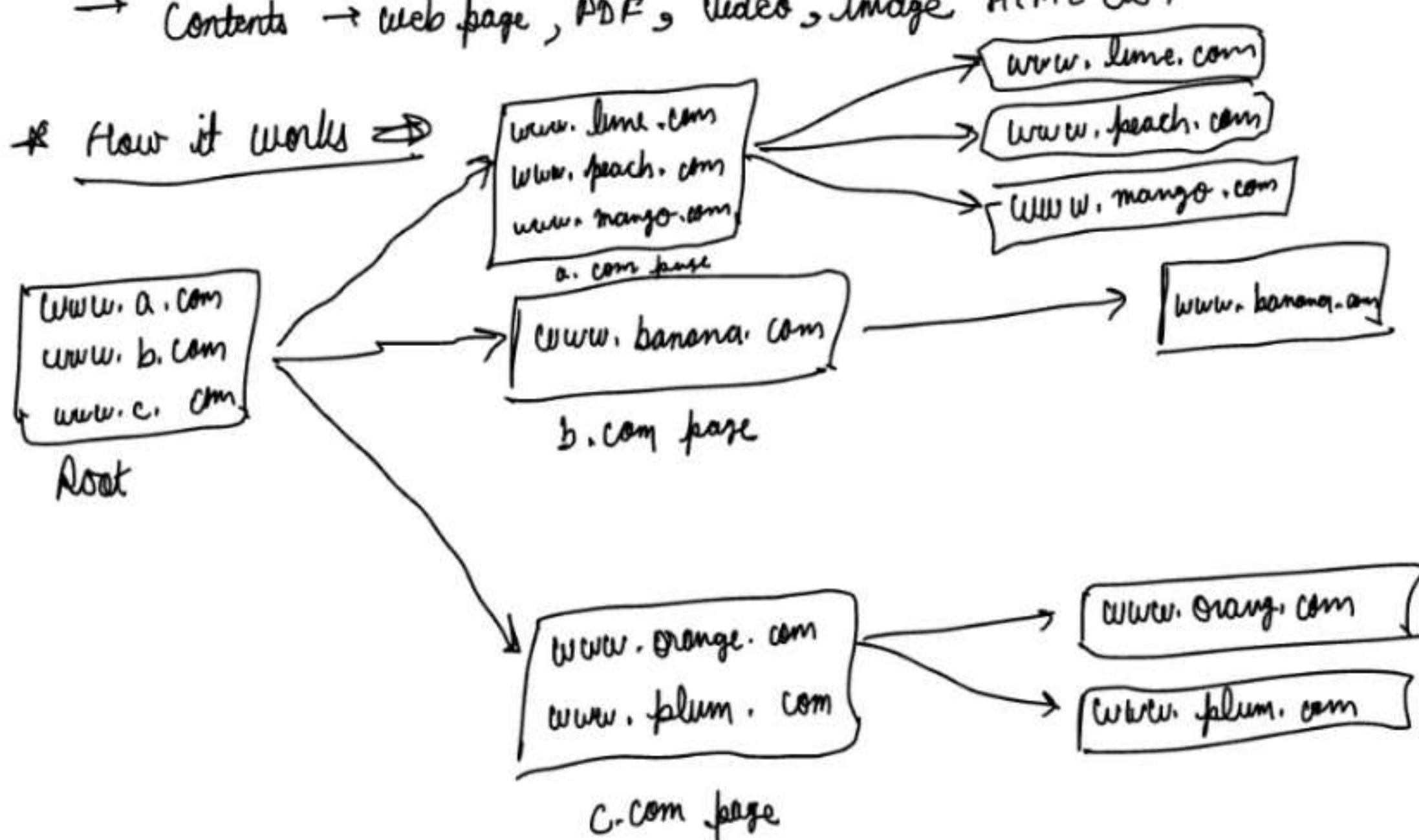
Web Crawler \Rightarrow

\rightarrow Robot OR Spider

\rightarrow Used by search engines to find contents on Web.

\rightarrow Contents \rightarrow web page, PDF, video, image, HTML etc.

* How it works \Rightarrow



Step ① Understand the problem and Establish Design Scope \rightarrow

The Basic algorithm \rightarrow

① set of URLs \rightarrow Download all web pages addressed by URLs.

② Extract URLs \rightarrow from these web pages.

③ Add new URLs to the list of URLs then, Repeat from step ①

Is it simple & fast? \Rightarrow NO.

Let's start asking Questions \rightarrow

Candidate \rightarrow what is the main purpose of web crawler?

\rightarrow search engine indexing?

\rightarrow Data mining?

\rightarrow OR something else?

Interviewer \rightarrow Let go for 'Search Engine Indexing'

Candidate \rightarrow How many web pages does the web crawler collect/month?

Interviewer \rightarrow 1 billion page.

Candidate \rightarrow What content types are included? HTML / PDF / images etc.

Interviewer \rightarrow HTML

Candidate → Do we also need to store the downloaded web pages →

Interviewer → Yes store up to 5 years.

Candidate → What if we get web pages which are duplicates?

Interviewer → Ignore the Duplicate web pages.

Besides functionalities you should note down characteristics of a

→ Good Web Crawler ⇒

→ Scalability →

→ web is huge

→ Billions of web pages

extremely efficient तरीका पड़ेगा

using parallelization

→ Robustness →

→ BAD HTML

→ Unresponsive server

→ Malicious links etc.

→ Handle all these.

→ Politeness →

→ Crawler should not make too many requests within a short time interval.

→ Extensibility →

→ system flexible होना चाहिए

→ New feature add करने के लिए Minimal changes करने पड़े।

→ example → Crawler image files in future.

55 Design Web Crawler →

Part ② : Step ① : Back of the Envelope Estimation → ,

* Back of the envelope Estimation →

→ Assume → 1 Billion web pages download / month.

→ Query per second (QPS) → $1,000,000,000 / 30 / 24 / 3600 \approx 400 \text{ pages/second}$ download करने वाली एक वेब क्रॉलर.

→ Peak Query per second (Peak QPS) = $2 * \text{QPS} = 2 * 400 \approx 800 \text{ pages/second.}$

- Assume → web pages average size = 500 KB
- $(1,000,000,000) \times (500 \text{ KB}) = 500 \text{ TB. storage / month.}$
- Assume data will be stored for 5 years
- Total storage requirement = $500 \text{ TB} \times 12^5$
 $= 30,000 \text{ TB} = \boxed{30 \text{ PB}}$. storage will be needed for 5 years.

Q56 Design a Web Crawler (Part 3) →

Step ②: propose a high level design & Get Buy-in →

→ starting में URL स्टार्टिंग होगा जा तो get started with seed URLs.

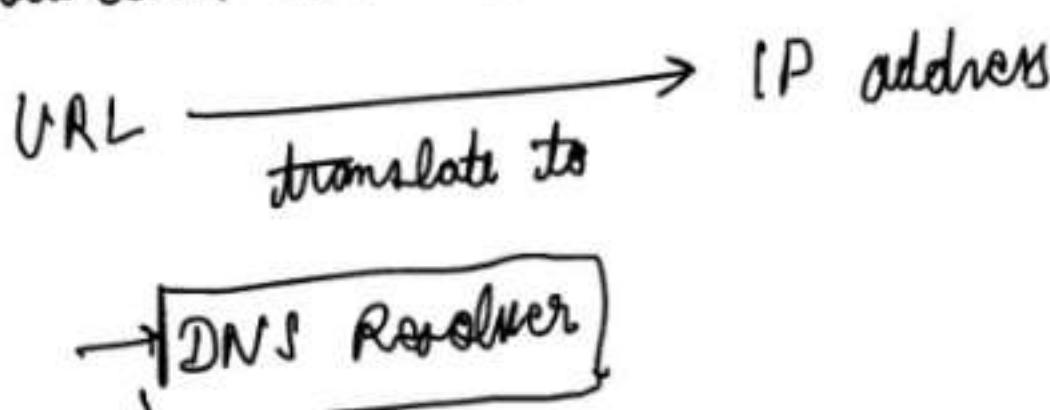
choose Good seed URLs

- Shopping URLs
- Sports URLs
- Health Care URLs etc.

→ जिन URLs के download करना है वो कहीं पर store करेंगे जा!!! → **URL Frontier**

→ URL frontier से URL मिलेगा तो download करेगा फिर web pages - HTML downloader

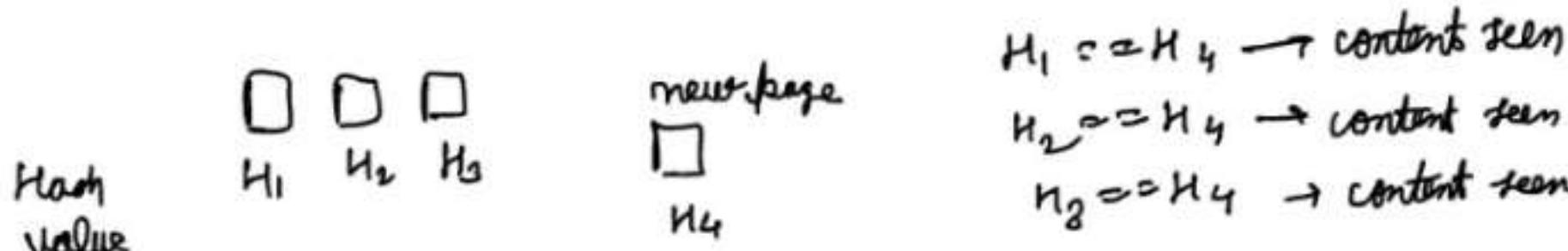
→ web page को download करने के लिये



→ Download होने के बाद parse, validate

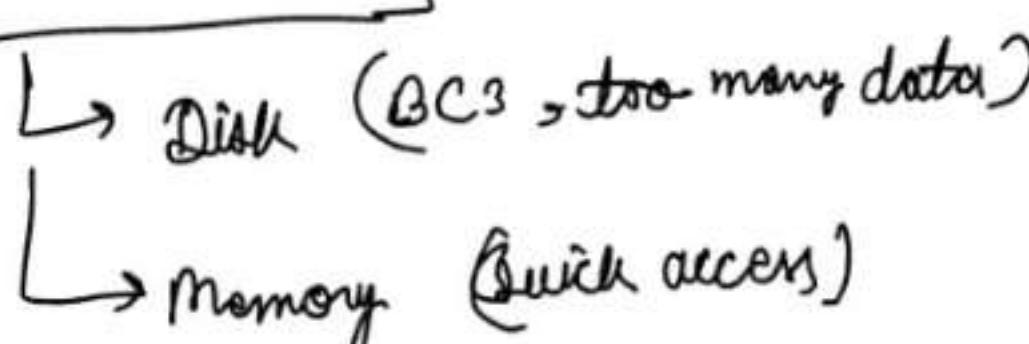
- **Content Parser**

→ Content seen? भी तो check करोगे ??



→ content को store करने की भी बात है शी

content Storage

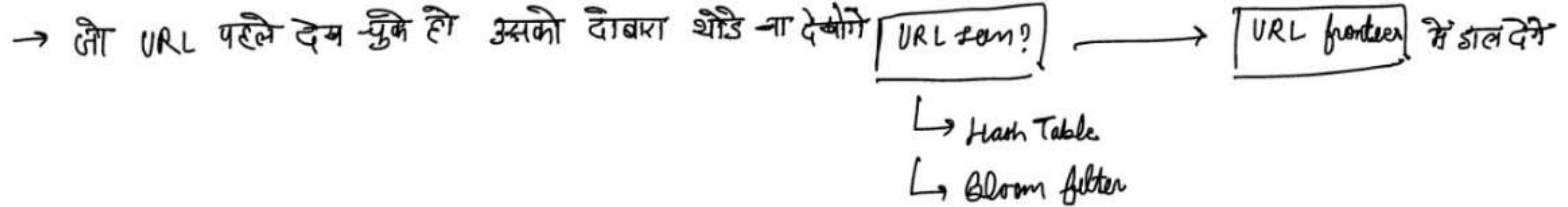


→ HTML pages से भी तो URL निकालना है।

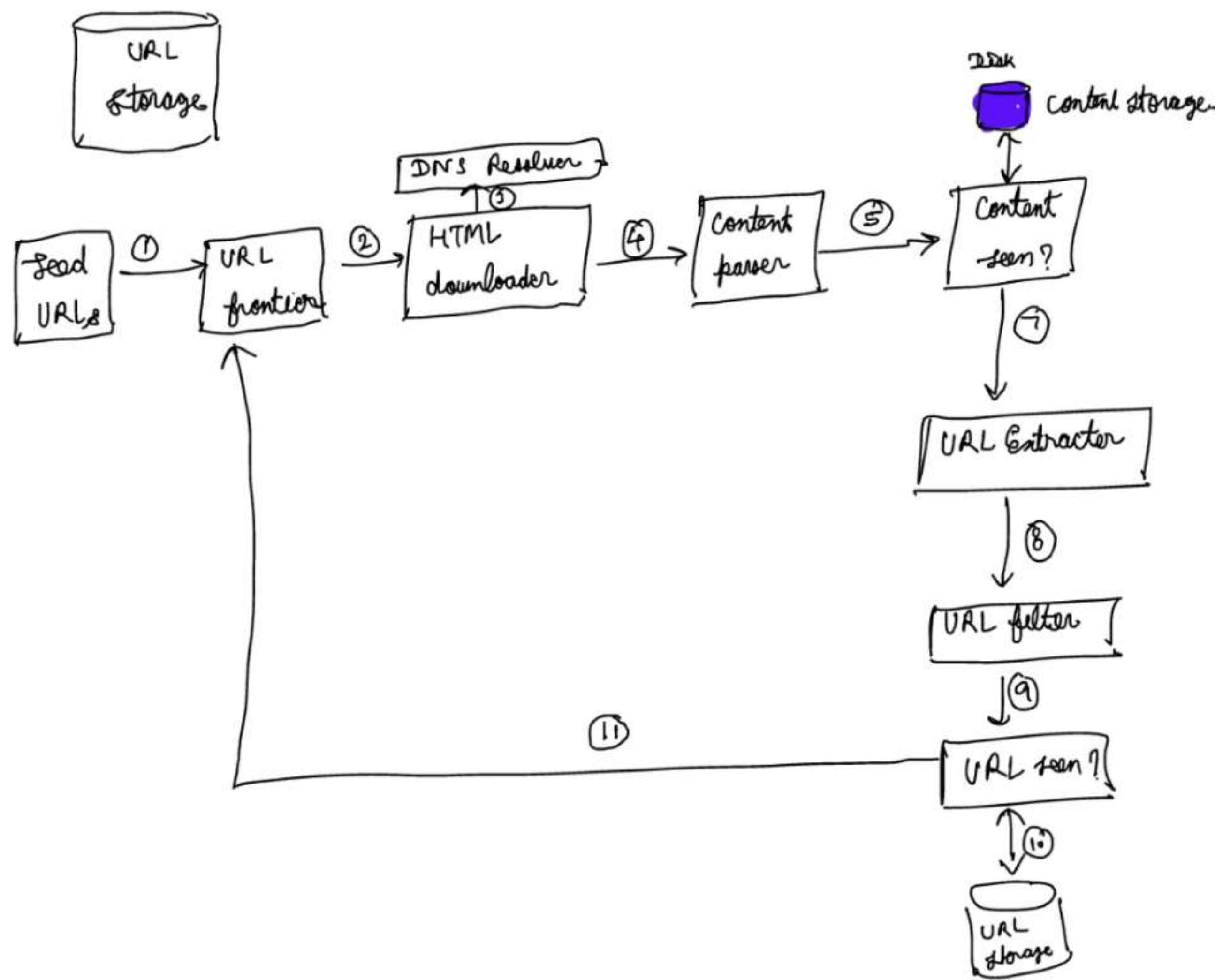
URL Extractor

→ बेकार Broken links को filter जा तो करोगे

URL Filter



→ जो URLs visit होगे उसको करी store कर लें।



57) Design Web Crawler Part ④

Step ③) Design Deep zinc (BFS v/s DFS) →

→ DFS v/s BFS.

→ URL frontier

→ HTML Downloader

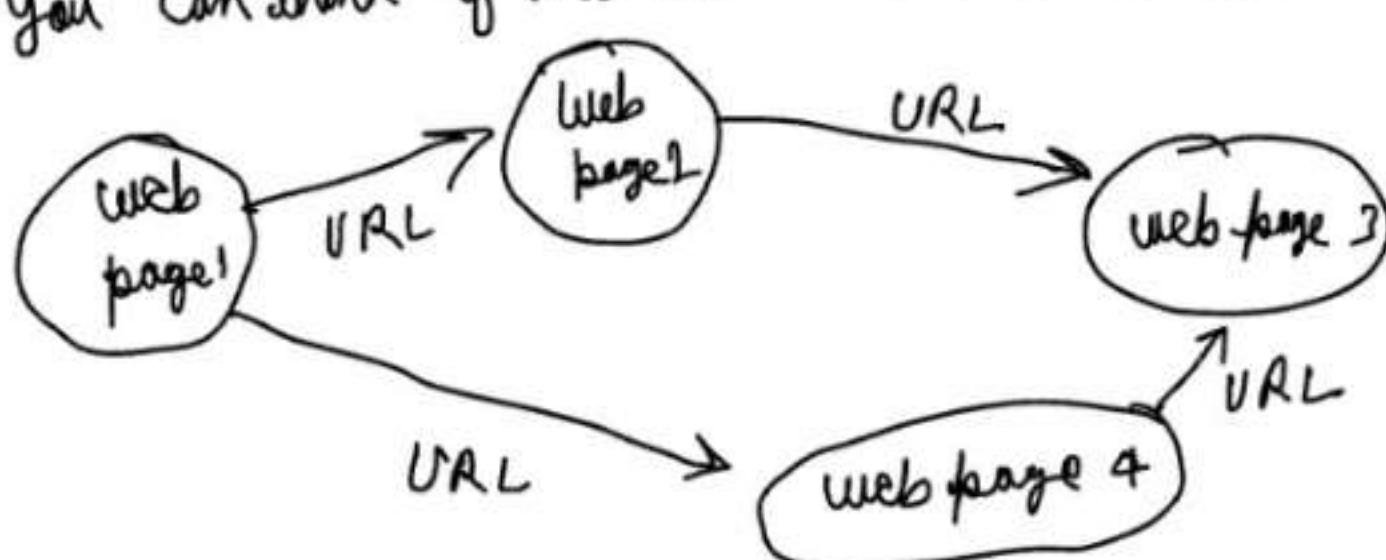
→ Robustness

→ Extensibility

→ Detect and avoid problematic Content.

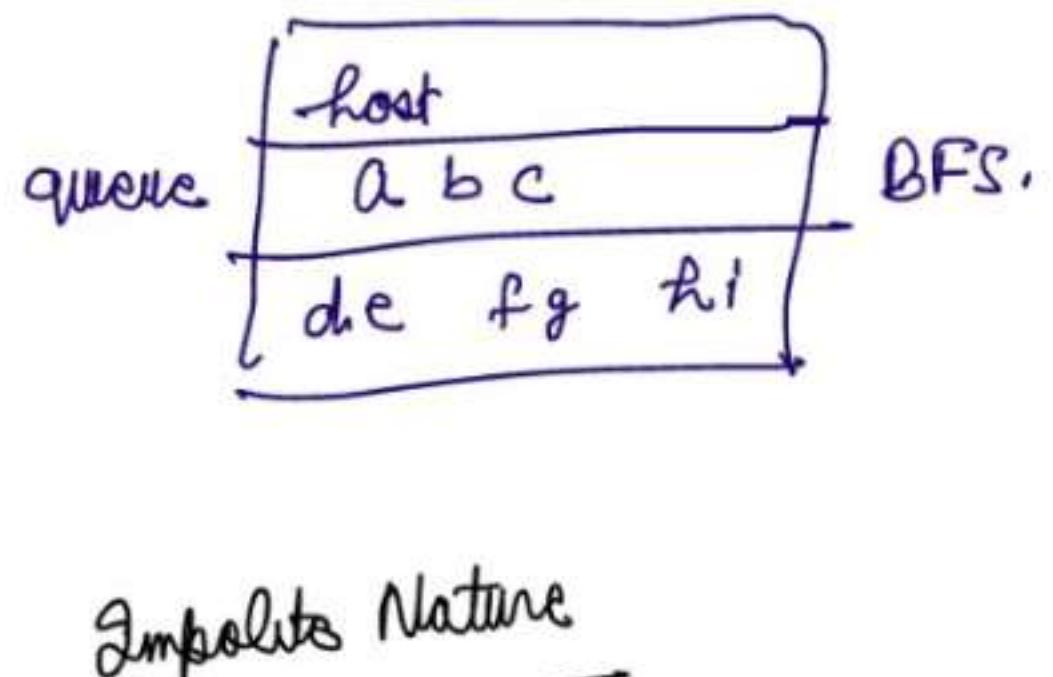
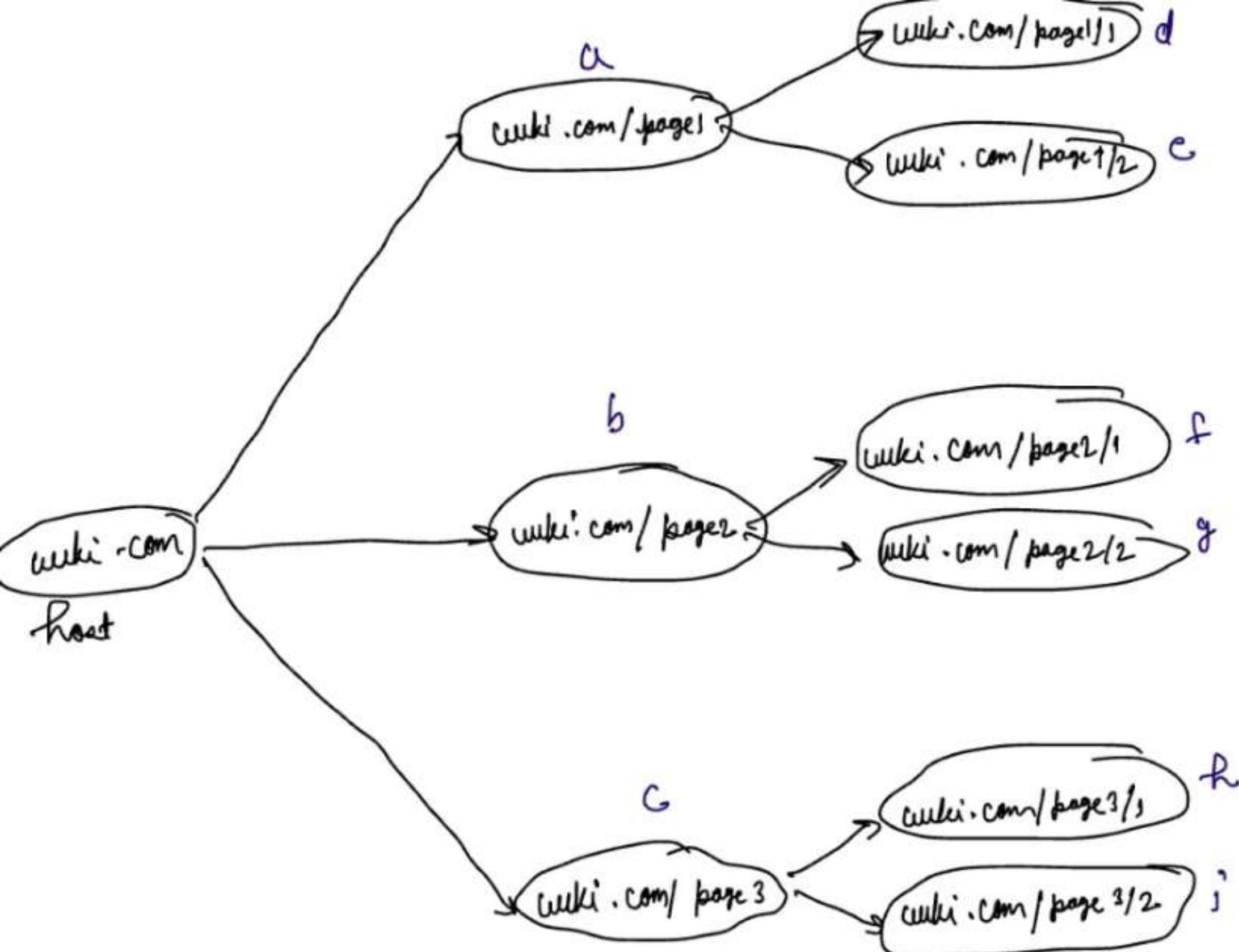
DFS v/s BFS →

→ you can think of web as → Directed graph.



CRAWLING → "Traversing a directed graph from one web page to other."

BFS & DFS →



Impolite Nature

→ इसकी Host को use करने का रोटर
↓
wiki.com

BFS →

Rank

(प्रारंभिक) Rank या priority व्याप्त है उसको पहले queue में insert करें।

webpage 1	1
webpage 2	2
webpage 3	3

⑤8 Design a Web CRAWLER part 5 →

step ③ Design Deep Dive (URL Frontier) →

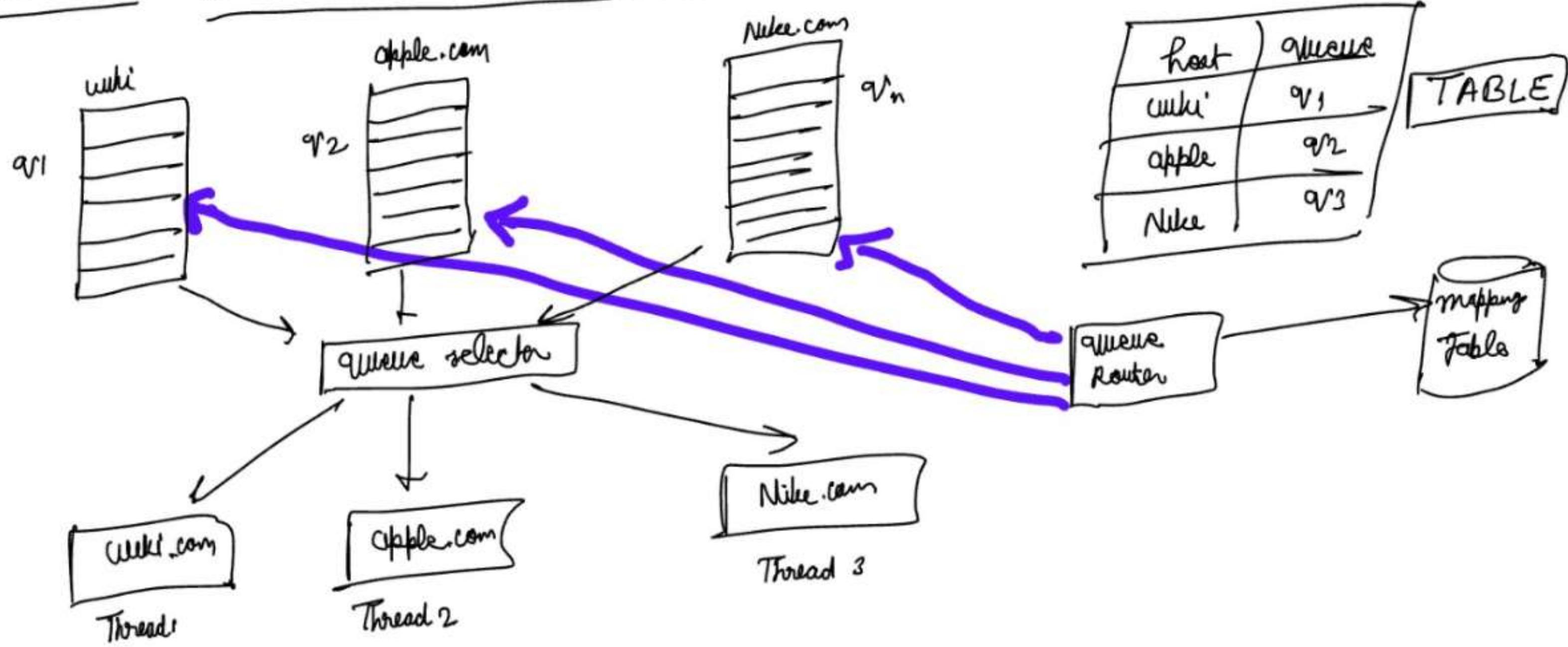
→ BFS will problem को solve करेगा

→ ensure politeness

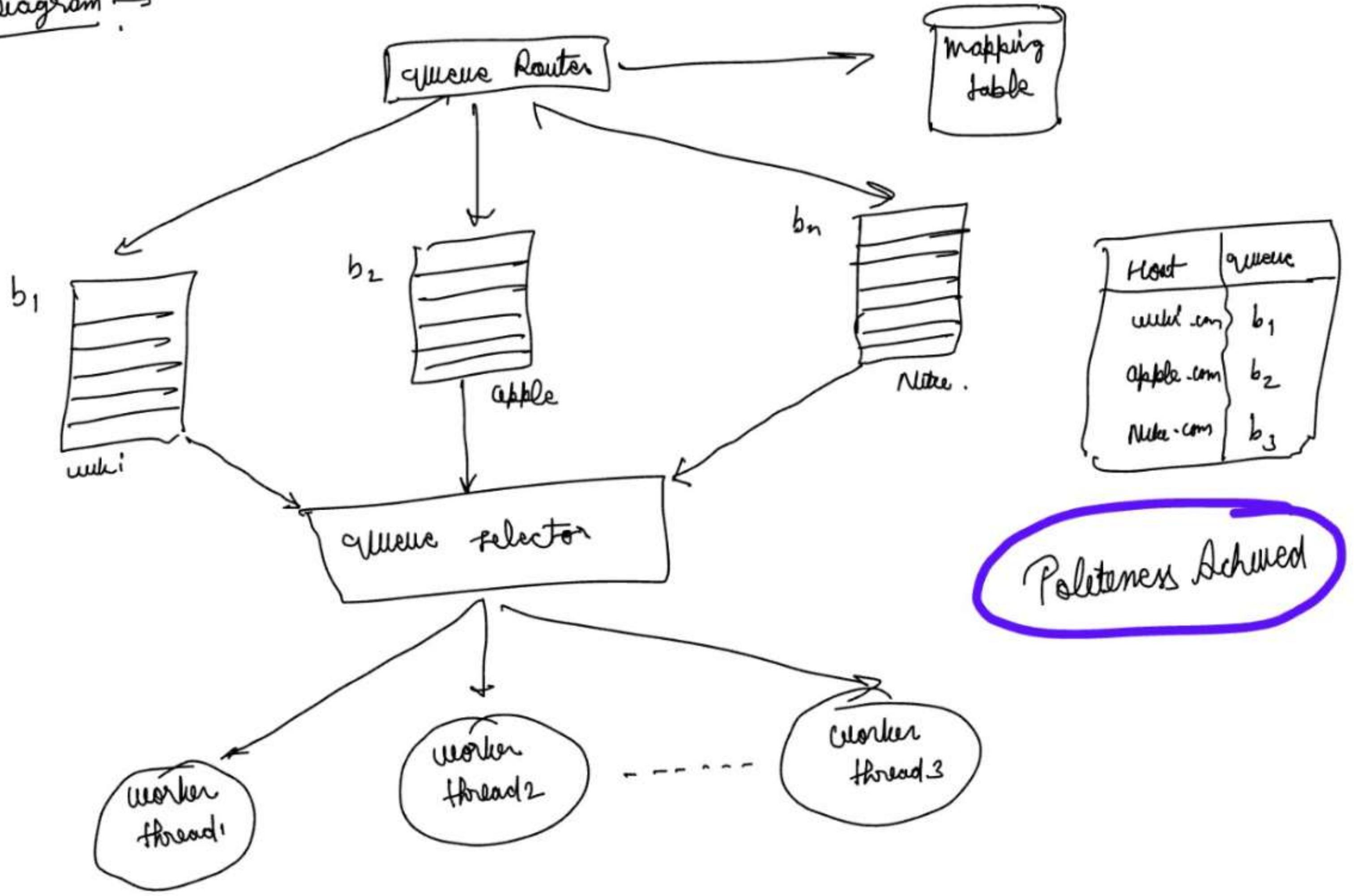
→ URL prioritization

→ Freshness.

Politeness → Avoid sending too many request to same host within a short period of time.



Nice Diagram →



Priority →

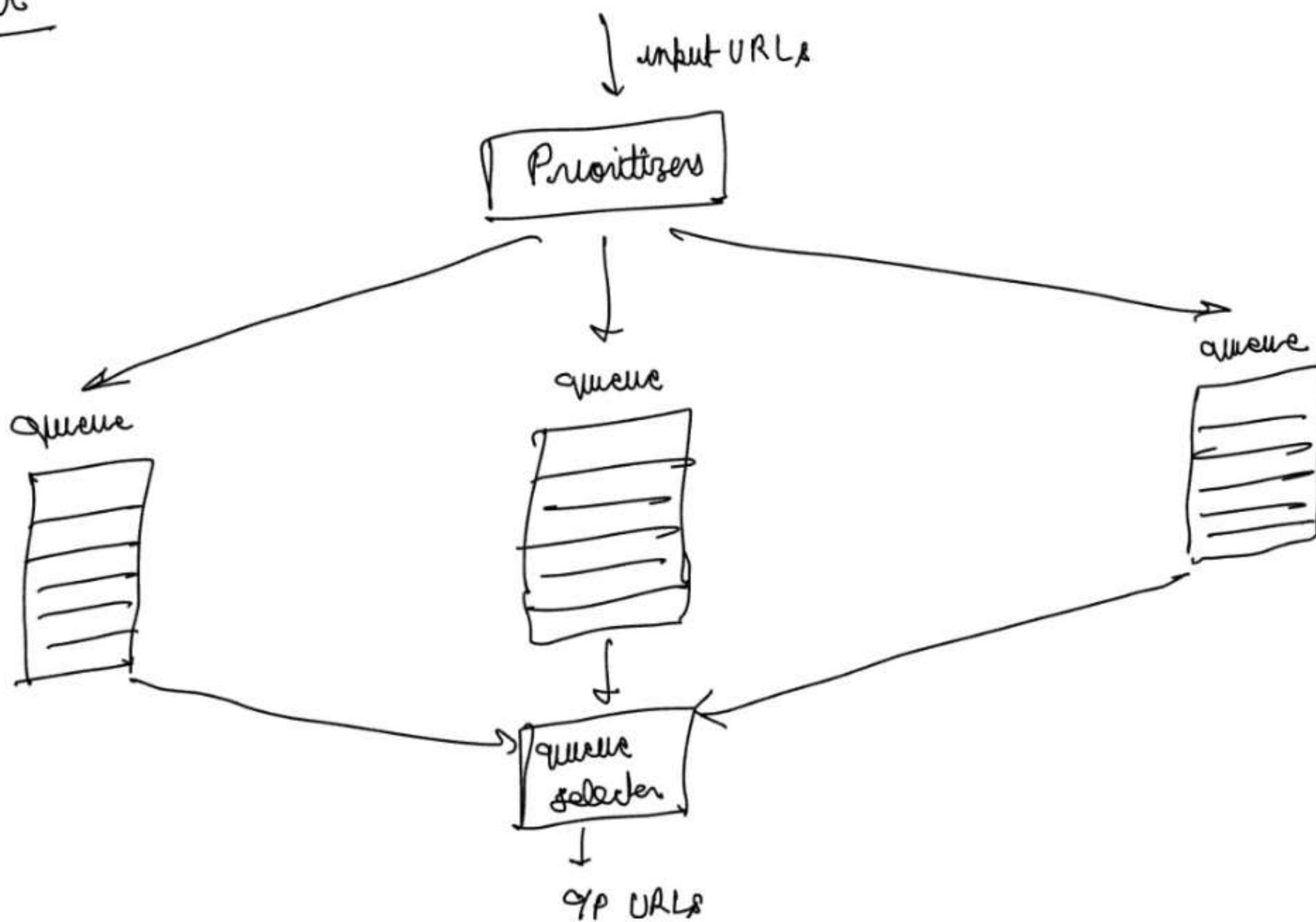
Random Discussion Forum
--- "Apple" ---

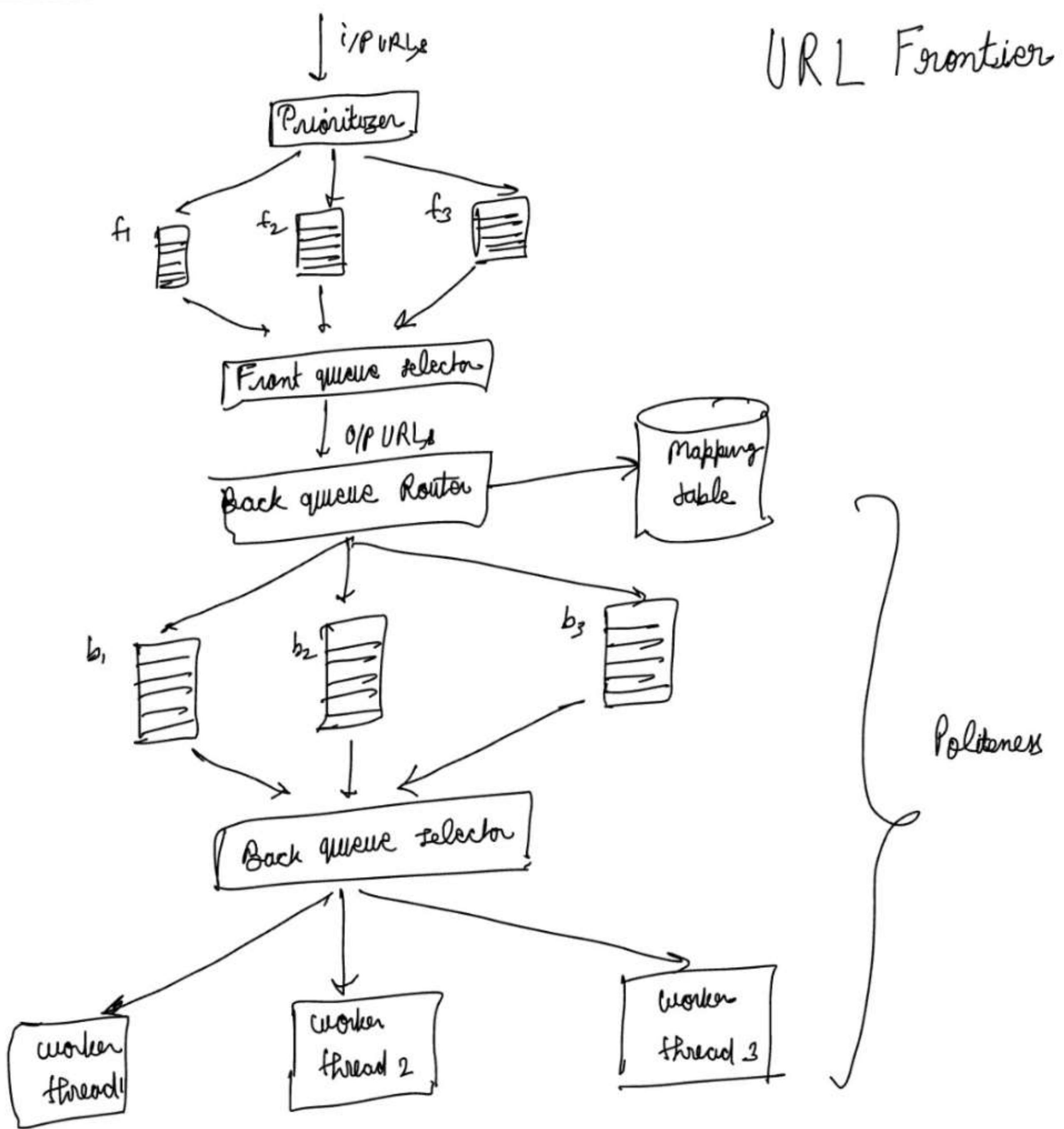
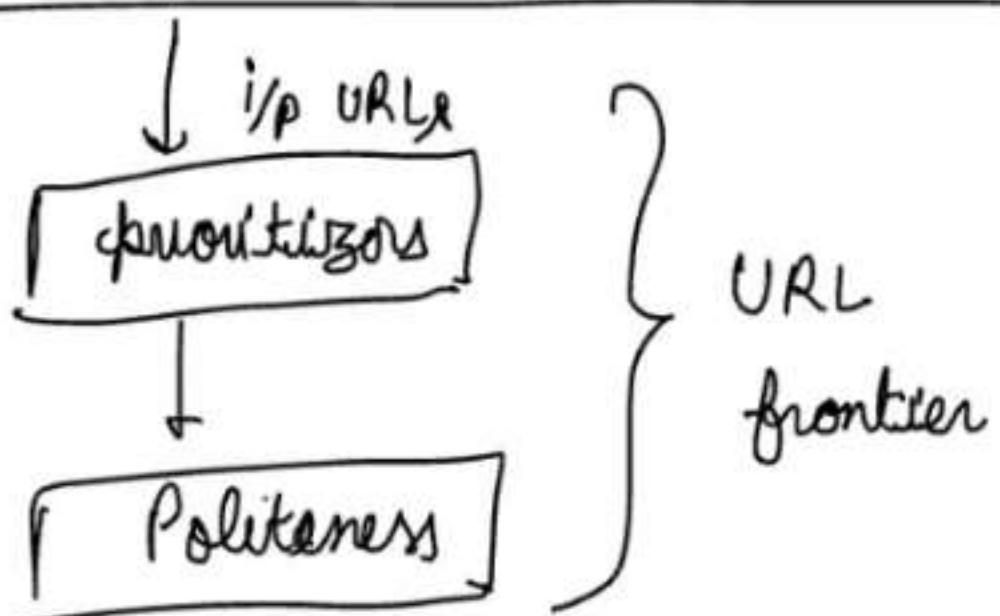
Apple Home Page

--- "Apple" --- ←
most important
most Priority.

usefulness → prioritise
↳ More web traffic
↳ update frequency

Prioritizer →





Freshness →

Web Pages →

- ↳ Edited
- ↳ Deleted
- ↳ Added.

→ Last update history.

→ Prioritize → web pages.

Storage for URL Frontier

→ Real World → Millions of URLs

→ Putting everything on Memory

↳ Neither Durable

↳ Nor Scalable.

→ Disk में भी सब कुदू नहीं रख सकते

Because Disk is slow.

→ so we use Hybrid.

→ Disk → Majority URLs

→ Memory → Remaining URLs

To avoid read & write load on Disk.

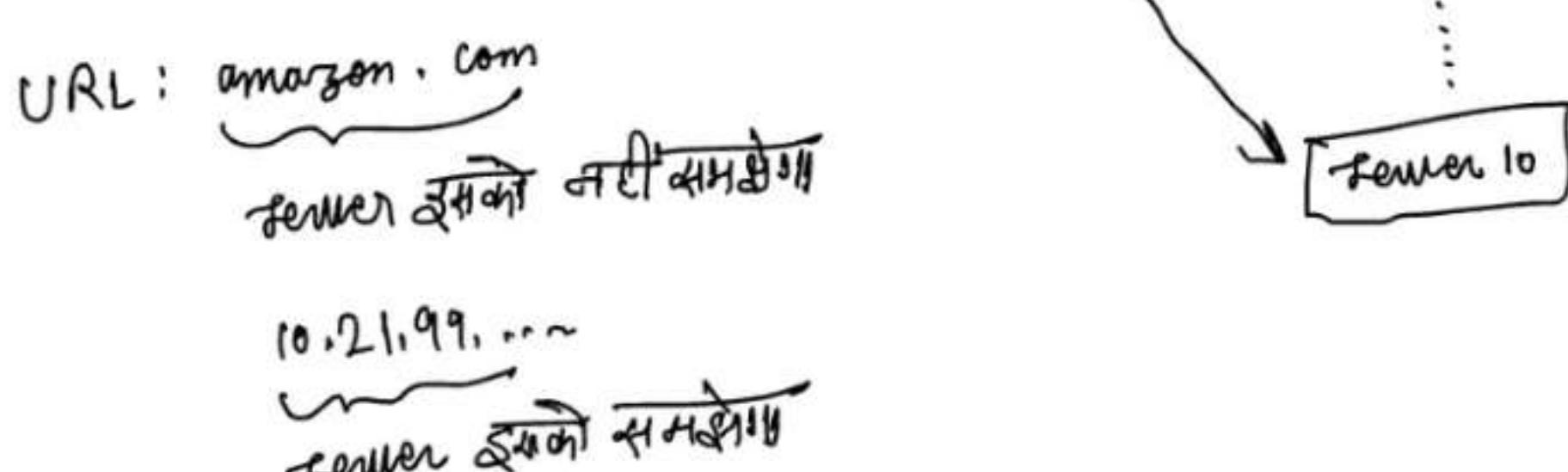
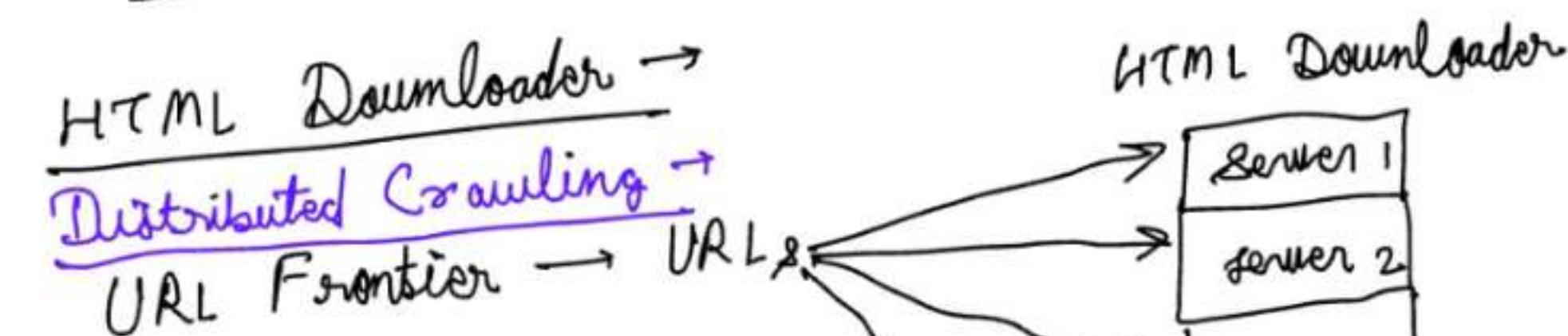
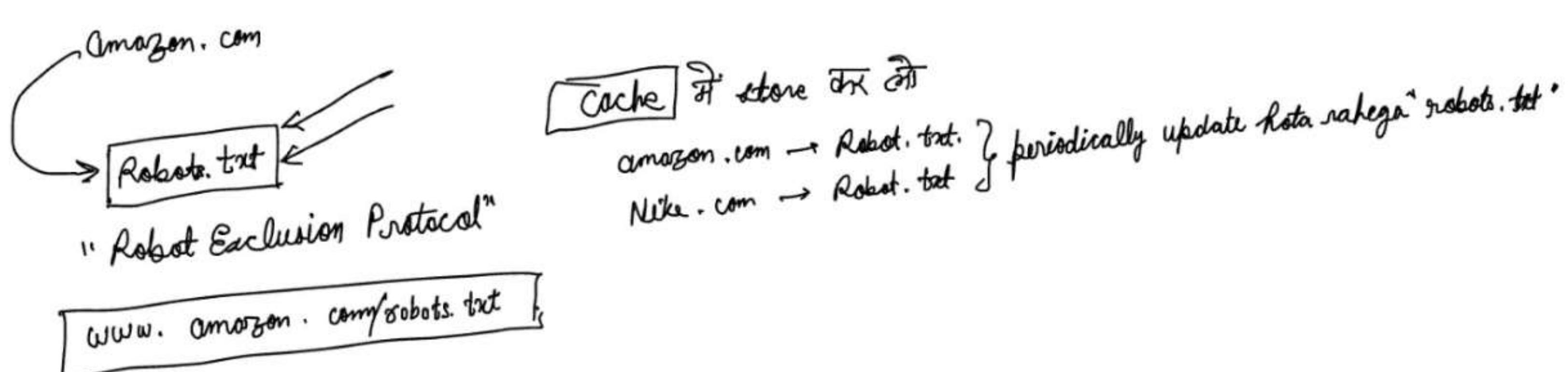
(59) Design Web Crawler (Part 6) ⇒

Step ③ Design deep dive Complete & WRAPUP →

* HTML Downloader → Download web pages from internet.

Interviewer → Can web crawler download everything from web page? (Because its not safe)

Candidate → No, Crawler को कुदू Rule follow करने पड़ेगे जो web page ने सेट किय है।



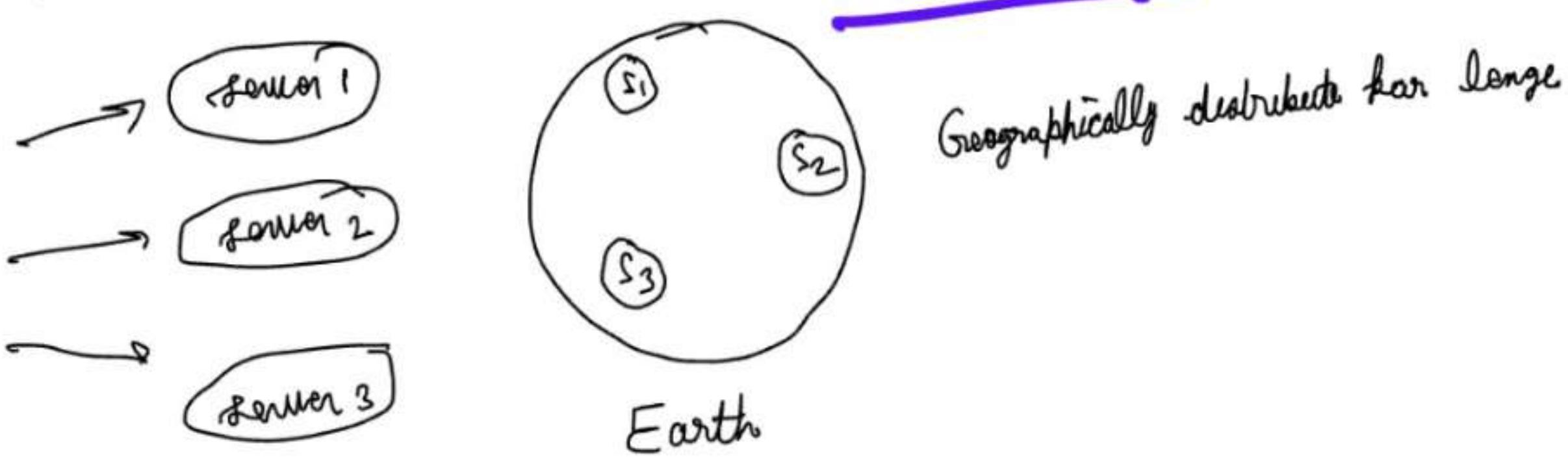
→ Cache \Rightarrow DNS result को store कर लेगा।

Amazon. com \rightarrow 10.21.99. ...

Nike. com \rightarrow 10.32.87. ...

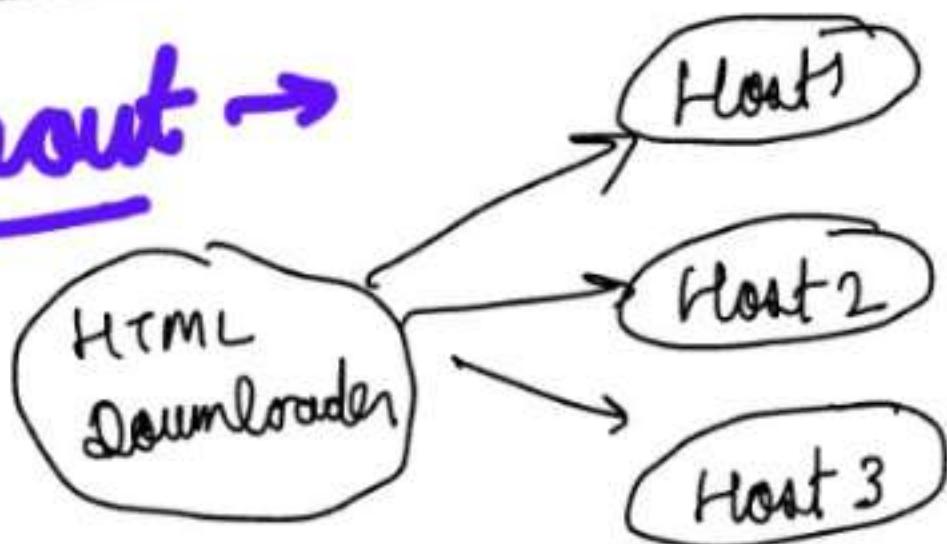
HTML Downloader \rightarrow

Locality:-



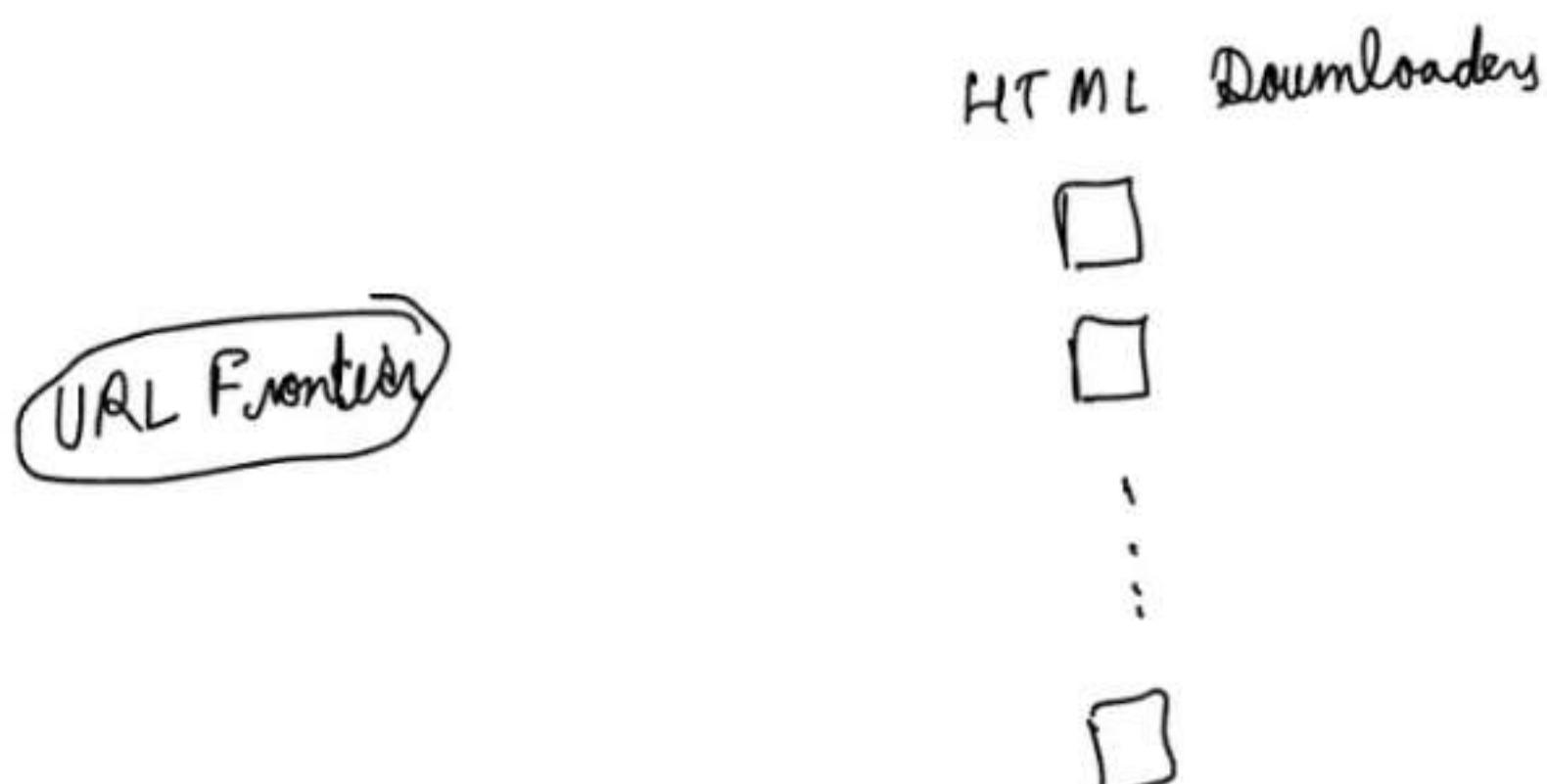
HTML Downloader \rightarrow

set Timeout \rightarrow



Robustness \rightarrow

→ किसी गड़बड़ को handle करने की ability.



→ We will use consistent Hashing.

मान लिये S_1 crash कर गया

→ Consistent Hashing

HTML downloader

↳ Add/ Remove Server.

↳ we will save state & data of CRAWL

↳ Exception Handling

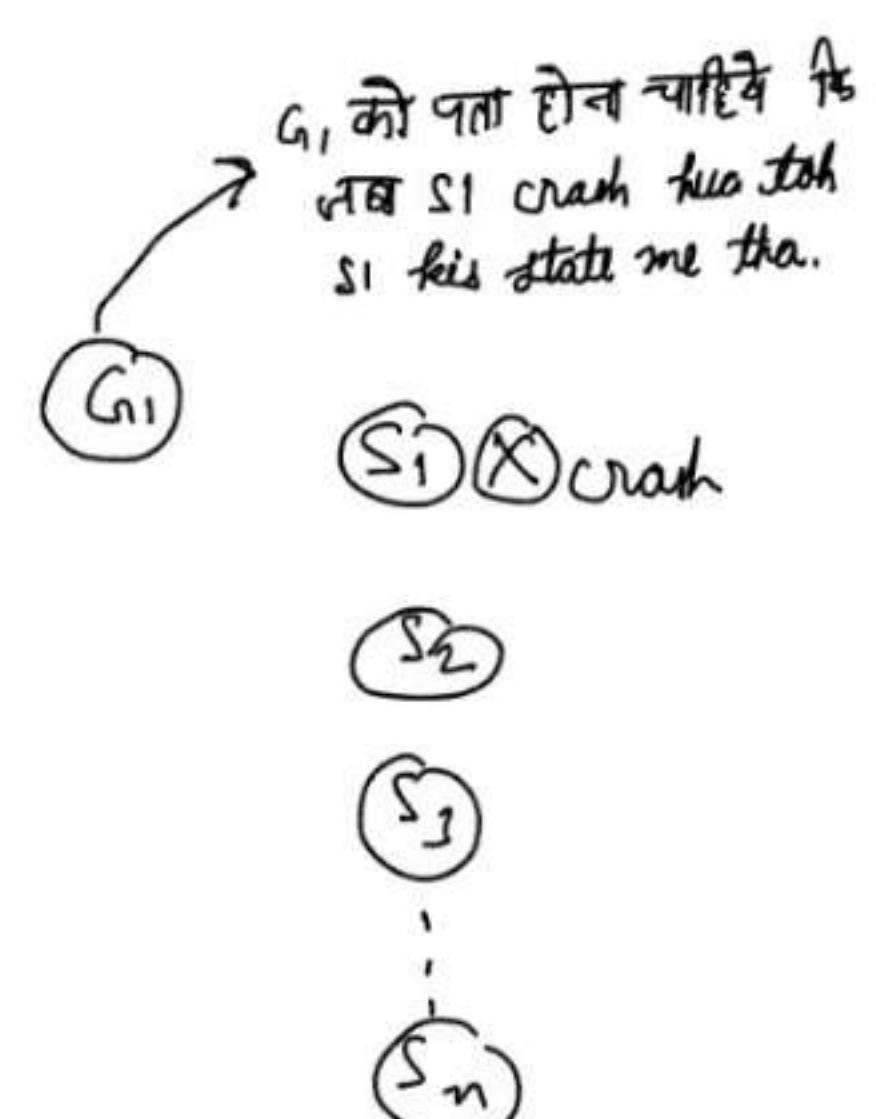
(Handle error, crash etc gracefully)

$S_1 \otimes$ CRASH

S_2

:

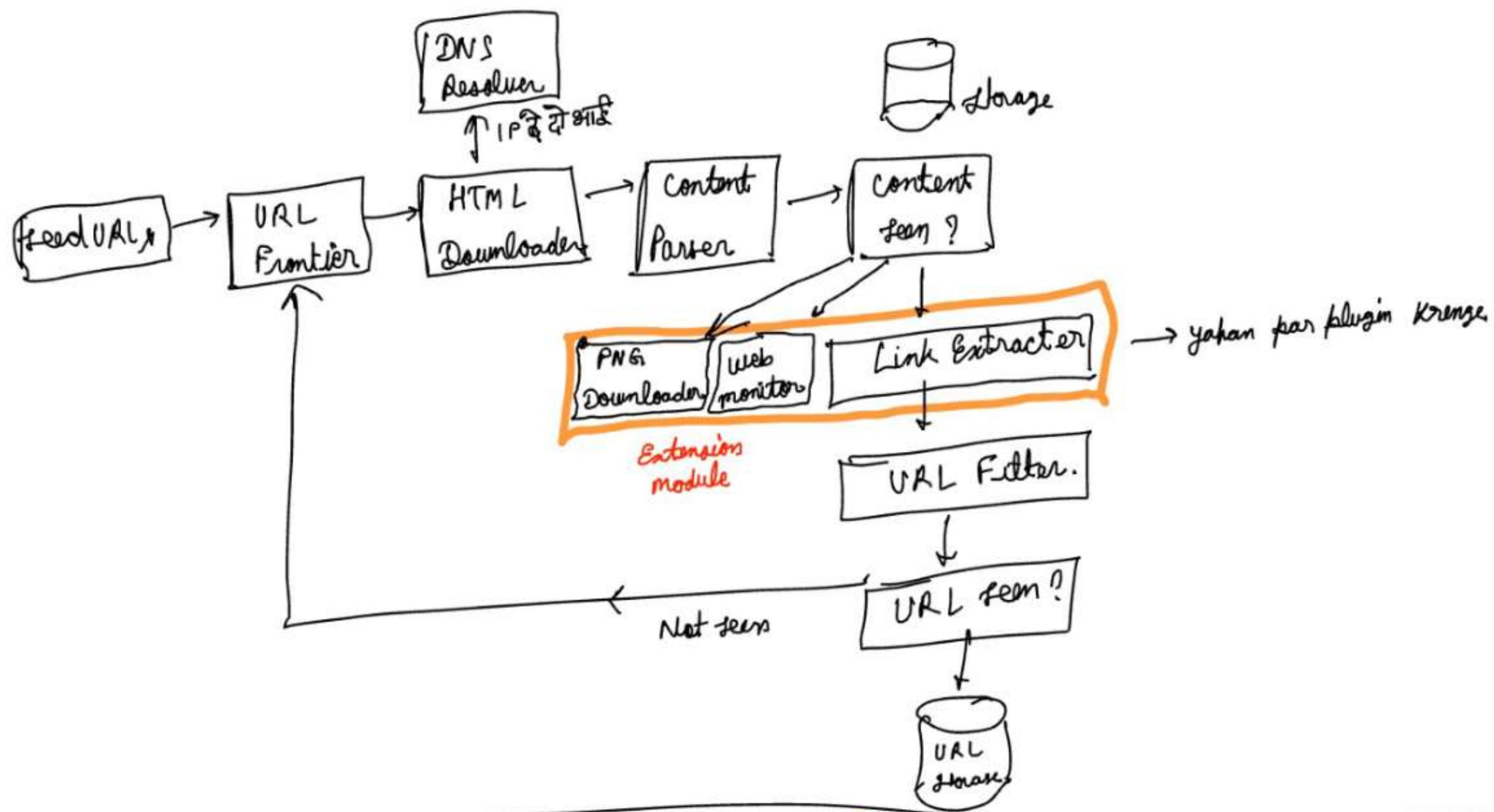
S_n



* Extensibility →

Interviewer → What if in future, I want to download PNG files also?
Where in the design you will put it.

Candidate → New Modules plug कर सकते हैं।
Let's recall our high level Diagram.



Detect & Avoid Problematic Content →

- Redundant (Duplicate) Data 30% duplicate data
- Harmful Data (Spider Traps) → www.studentstakeexample.com/foo/bar/foo/bar/ foo/bar/...
- Meaningless Data. (Data Noise)
 - e.g. → advertisement
 - code snippet
 - spam URLs

WRAP UP →

Good CRAWLER →

- ① scalability
- ② Politeness
- ③ Extensibility
- ④ Robustness

website → AJAX, JS



Web pages	
Link generated dynamically!	
" "	1
" "	2
" "	3

Server side Rendering →

- unwanted Pages → filter out.
- Horizontal scaling
- Availability consistency & Reliability .