

SYSTEMS DESIGN

SYSTEMS DESIGN CONCEPTS

- i. Introduction
- ii. Client – Server Model
- iii. Network Protocols
- iv. Storage
- v. Latency and Throughput
- vi. Availability
- vii. Caching
- viii. Proxies
- ix. Load Balancers
- x. Hashing
- xi. Relational Databases
- xii. Key – Value Stores
- xiii. Replication and Sharding
- xiv. Leader Election
- xv. Peer to Peer Networks
- xvi. Polling and Streaming
- xvii. Configuration
- xviii. Rate Limiting
- xix. Logging and Monitoring
- xx. Publish/Subscribe Pattern
- xxi. MapReduce

SYSTEMS DESIGN MORE CONCEPTS

- i. Hoizontal and Vertical Scaling

SYSTEMS DESIGN BASIC QUESTIONS

- i. Design AlgoExpert
- ii. Design a Code-Deployment System
- iii. Design a Stockbroker
- iv. Design Facebook News Feed
- v. Design Google Drive
- vi. Design the Reddit API
- vii. Design Netflix
- viii. Design the Uber API
- ix. Designing WhatsApp
- x. Designing Tinder
- xi. Designing Instagram News Feed
- xii. Google Maps : Location based database

SYS DESIGN FUNDAS

4 categories → intertwined

- 1] Underlying, fundamental knowledge → Client-server model, network protocols. Gotta atleast understand all these to have a knowledge of sys. des.
- 2] Key characteristics of sys → Things you might want a sys to have, things that you might be trading off while making design decisions.
eg:- availability, ~~failover~~ wait & see, throughput, redundancy, consistency
- 3] Actual sys. components → Tangible things that you hr/ implement in a sys. Bread & butter of sys.
eg:- load balancer, proxy, cache, rate limiting, leader elec"
- 4] Tech → Real, existing products or services that you can use in sys either as actual components or to achieve a certain charac. in a sys. These are real tools
eg:- Zookeeper, XCP, Engine X, Reddits, Amazon S3, Google Cloud storage

DATA STRUCTURES & ALGORITHMS

I

SYSTEMS EXPERT

Design fundamentals → founda"al knowledge reqd. for systems design & interviews

Data structures → founda"al knowledge reqd for coding interviews

e.g. of design fundamentals: SQL, servers, cache, polling, load balancer, HTTP, database, hashing, replication, client, processes, Nginx, availability, leader elec", P2P

IA

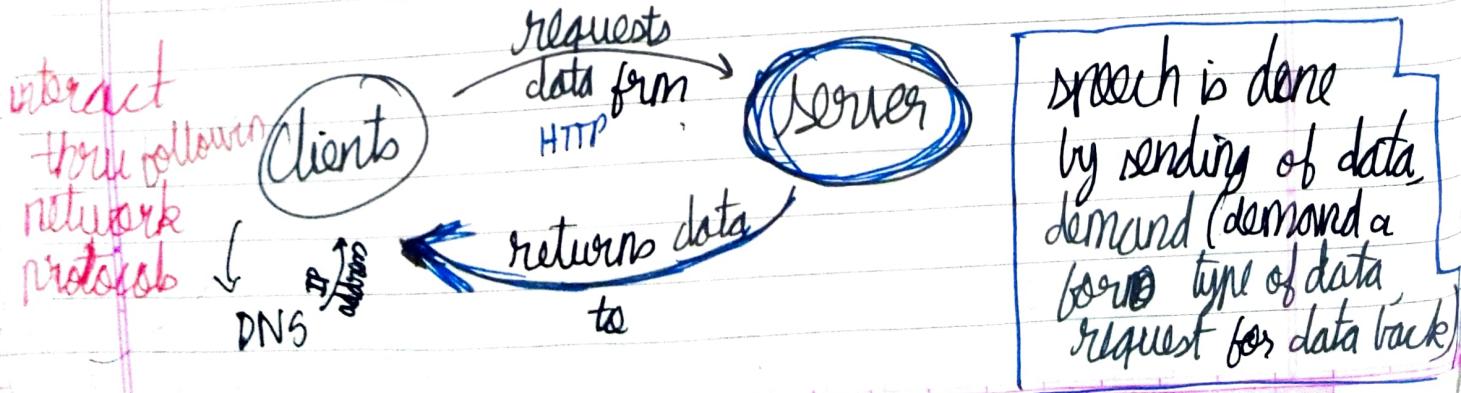
Client-Server Model/Architecture/Paradigm

founda" of modern internet, helps us understand how computers speak to one another. Paradigm consists of client requesting service, data from servers and servers providing it all

What we do? → Type URL in browser and hit 'Enter'

What happens? ^{or software} single machine can act as both client, server. e.g.: Node and Postman on PC

- (i) Client: Is something, a machine that speaks to the server.
- (ii) Server: something else, a machine that listens to the client, listens for clients to speak and then speaks back to the clients



(iii) Browser : Client
Website : Server
(ones which have server)

⇒ NOTE: A website (same) may have diff. servers in diff. loca's. Eg:- Netflix India has a diff. server from Netflix USA.

(iv) Browser doesn't really know what the server is. All that it knows is that it can communicate with server. It doesn't really know what the server represents. It just requests info from it and does stuff based on info recd. from server.

(v) To start communicating with a server, a browser 1st sends a DNS (Domain Name Server) query to find out IP address of server.

DNS query → special request going to a ~~set~~ predetermined set of servers asking for IP address of a server

IP address → unique identifier for a machine. All comps. connected to internet have ways to find out these IP addresses or choose routes to those addresses. They can send packets of data/info in form of bytes to IP address.

Analogy :- IP address is like a mailbox that some entity has granted to a machine

Eg:- Algo Expert's IP address has been granted to it by ~~Google~~ its cloud provider - Google Cloud Platform. It reserved an IP address for Algo Expert.

Eg:- Thus, browser makes DNS query for algoexpert.i
receives back an IP address and thus starts communica
with the server.

(Vi) HTTP: way to send info that comp. can understand.

Exercise : type "dig algoexpert.io" → does a
DNS query and returns IP address of
algoexpert.io

→ when browser sends HTTP req. to the server, it basically sends a bunch of bytes / char
that are gonna get packed into some special
format and sent to servers. This request also
contains IP address of your PC (a.k.a. source
address).

When server receives source address, it knows that
on which IP address it needs to send a
response to.

(Vii) Ports - Servers usually listen for requests on specific
ports. Any machine having a distinct IP address
has 16000 ports that progs. on the machine can
listen to.
On communicating with machine, you gotta specify
what port you wanna communicate on ^(client)

Analogy → IP address : Mailbox to an apartment complex
Ports : actual apartment no. but the
mail (~~you~~ client req.) arriving
at mailbox has to route to.

Most clients know the port that they should use

depending on the protocol that they're trying to speak to server with.

Protocol
HTTP
HTTPS

Port used by client

exercise: netcat - allows you to read from/write to network connection
using various protocols

nc-l 8081

says that
seen to stuff happening on this part

READING DATA

Anything that you type in terminal window 2 appears.

special IP address that always
points to your local machine
(local machine's IP address)

WRITING DATA

This terminal is entering a
communi. channel with
the machine at this IP
address at port 8081

(viii) Thus, once server receives req., it is able to read it
cuz it understands the HTTP format
Server understands that when you go to ae.jo,
you're trying to see the HTML of "ae" and so it returns
② a response viz. the HTML of "ae" to browser
which receives response and renders the HTML
on the page for you
↓
thus HTML is sent

turns HTML code into text, imgs, multimedia, interactivity effects, etc.

Add'l info:

IP address: IP addresses consist of 4 nos. separated by dots a.b.c.d where all 4 nos $\in [0, 255]$

- i) 127.0.0.1 \Rightarrow Your own local machine a.k.a localhost
- (ii) 192.168.x.y \Rightarrow Your private network.

Ex:- Your machine and all machines on your net. wifi ~~are~~ network have the 192.168 prefix

NETWORK PROTOCOL

msgs sent over wire or network = msgs sent over the internet from 1 machine to another machine

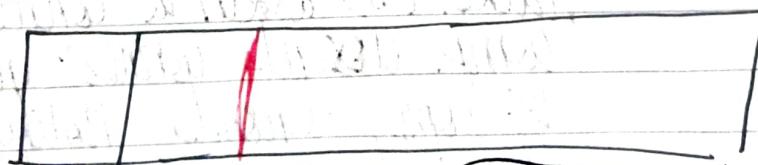
- types of msgs
- struc.
- order
- response to a msg, what shud it look like if present
- rules on when msgs can be sent to 1 another

Network protocol msg of

3 main protocols:

1) IP → Internet Protocol

- i) Modern internet effectively runs on IP when 1 machine interacts with another and (client) (server)
- sends data to it the "data" is sent in form of IP packet (fundamental block of data sent b/w machines i.e. communicaⁿ)
- ii) IP packet → made up of bytes



→ IP packet

IP header ↑
TCP header

data aka 'payload'

- iii) IP headers - At beginning of packet. Contains imp info → source IP address, destination IP address (machine sender data) (machines data will go to)

40

destinⁿ = dstn

total packet size, version of IP that packet is operating by (current used \rightarrow IPv4, IPv6). Based on IP version, packet might look diff header size : 20 to 60 bytes

- iv) Size of IP packet $\approx 2^{16} B \approx 65 KB$
Info sent as multi-IP packets across machines
Order in assurance that packets will be sent is kinda fixed.
That's why, TCP is reqd.

2) TCP - Transmission Control Protocol, that they can go to stream data thru open connecⁿ

- i) Implemented in the kernel, exposes sockets to apps, ^{which} order assuranc^e
- ii) Built on top of IP. Solves order assurance issues of IP packets in an error-free, uncorrupted way
- iii) Used in web apps mostly
- iv) TCP header \rightarrow present in data part of packet
contains TCP info like order of packets
eg:- browser wants to connect to a destn server such as Google. Browser is first gonna create a TCP connecⁿ with destn comp. / server thru a handshake

v) Handshake: Special TCP interactⁿ where 1 comp communicates with other by sending packet(s) asking to connect. Receiving comp. accepts connecⁿ invite
1st comp, responds notifying establishment of an open connecⁿ b/c 2 of em.

vi) If 1 of the machines doesn't send data in a given time period, connecⁿ can be timed out

devpr = developer

- vi) A machine can end comm. by sendin a msg "notifyin its inter" and then "TCP connect" is done
- vii) Thus, TCP is a more powerful, func'l wrapper around IP.
It lacks robust framework that developers use to define comm. channels b/w client-server in a sys. cuz its just data that fits into the IP packets underlying.

3) HTTP - Hypertext Transfer Protocol

- i) Built on top of TCP
- ii) Introduces higher level abstracⁿ (req-respons paradigm) above TCP
- iii) 1 machine req. 1 machine responds. Makes it easy for devpr to create easy to use, robust sys.
- iv) Req.: Machine that wants to interact with other machine sends this. Nr lots of persons defining HTTP
- v) Resp.: Other machine's reply to a req
- vi) Req., Resp → can be thought to be similar to obj. with imp fields. props. that describe em.

can be visualized (not exactly like this in reality) as

```
const httpReq = {
```

```
  devide { host: 'localhost',  
    port: 8080,  
    method: 'POST', // GET, PUT, DELETE, OPTIONS, PATCH  
    path: '/payments',  
    headers: {  
      'content-type': 'application/json',  
      'content-length': 51  
    },  
    body: {  
      'data': 'This is a low JSON format data piece' }  
  }  
}
```

provides data → server
retrieve data
delete data
Client sends data to server
which has various paths which dictate what kinda logic will occur per req. They're guidelines subj to your server as how we em.

```
const httpResponse = {
```

← status code: 200,

describes type headers:

of resp.

Status codes are
also like guidelines
that you can alter as
per requirement

e.g.: - 404 status

code = requested
data piece not found

Path $\xrightarrow{\text{contains}}$ logic gets
 $\xrightarrow{\text{executed}}$ as per path provided in req

Headers \rightarrow collection of k-v pairs contains imp. meta
data or info abt req.

e.g.: 403 status code \rightarrow data you're req. in is
forbidden

NOTE status code \rightarrow describes type of ~~status~~ response
status codes are like guidelines
that you can alter as per requirement

e.g.: - 404 status code = requested piece of
data not found.

Typical HTTP req. schema

host : string (eg: facebook.com)

port : integer (eg:- 3000, 80, 43)

method: string (eg:- GET, PUT, POST, DELETE, OPTIONS, PATCH)

headers: pair-list (eg:- "Content-Type" \Rightarrow application/json)

body: opaque sequence of bytes

Typical HTTP resp. schema

error
↑

→ all ok

status code: integer (eg:- 404, 200)

headers: pair list (eg:- "Content-Length" \Rightarrow 1238)

body: opaque sequence of bytes

```
JS server.js × JS http_request_example.js ●
```

```
JS server.js > ...
1 const express = require('express');
2 const app = express();
3
4 app.use(express.json());
5
6 app.listen(3000, () => console.log('Listening on port 3000.'));
7
8 app.get('/hello', (req, res) => {
9   console.log('Headers:', req.headers);
10  console.log('Method:', req.method);
11  res.send('Received GET request!\n');
12 });
13
14 app.post('/hello', (req, res) => {
15   console.log('Headers:', req.headers);
16   console.log('Method:', req.method);
17   console.log('Body:', req.body);
18   res.send('Received POST request!\n');
19});
```

...

```
Clements-MBP:network_protocols clementmihaiilescu$ node server.js
Listening on port 3000.
Headers: { host: 'localhost:3000', 'user-agent': 'curl/7.54.0', accept: '*' }
Method: GET
Headers: {
  host: 'localhost:3000',
  'user-agent': 'curl/7.54.0',
  accept: '*/*',
  'content-type': 'application/json',
  'content-length': '14'
}
Method: POST
Body: { foo: 'bar' }
```

```
Clements-MBP:network_protocols clementmihaiilescu$ curl localhost:3000/hello
Received GET request!
Clements-MBP:network_protocols clementmihaiilescu$ curl --header 'content-type: application/json' localhost:3000/hello --data '{"foo": "bar"}'
Received POST request!
Clements-MBP:network_protocols clementmihaiilescu$
```

Database → servers with long life that interact with the rest of your app thru network calls, with protocols on top of TCP or even HTTP

STORAGE

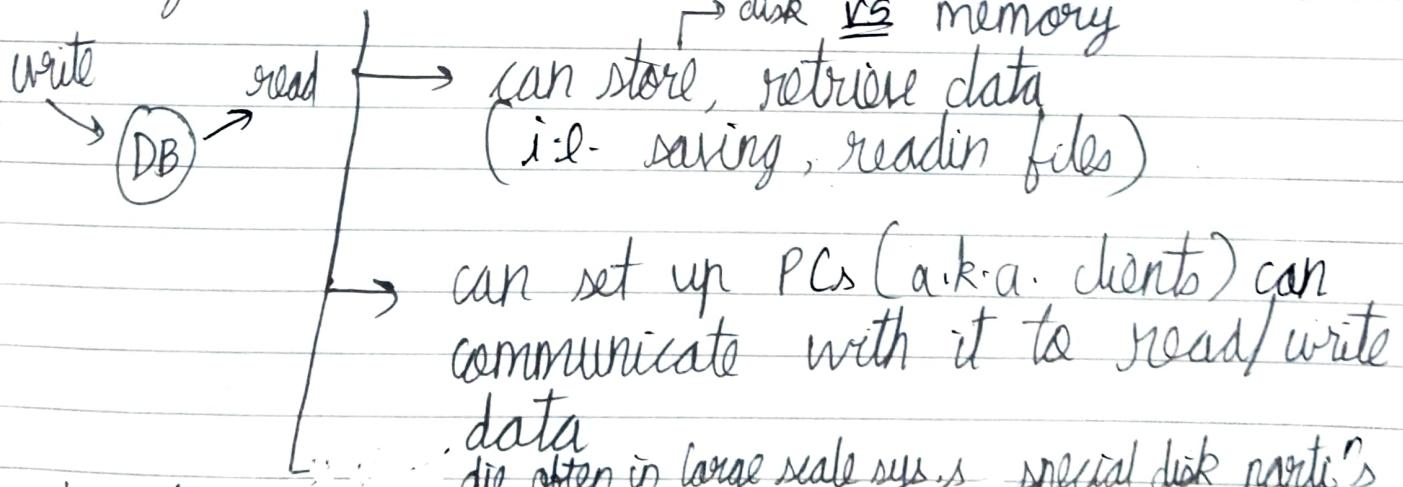
→ 1) Every system requires storage
eg:- for storin user info, metrics

→ 2) Database = a machine that helps storing and retrieving data

^{synonyms} setting ≡ writing ≡ recording ≡ storing
^{word} { pairs getting ≡ reading ≡ querying ≡ retrieving
^{w.r.t.} data.

DB

In most cases database, is actually just a server
eg:- even a PC can be made to act as a DB



→ 3) Persistence of data in a DB
a.k.a volumes are used by DB to access volumes of data even if machine die often in large scale sys.s special disk part's dies (crashes permanently)

→ non-volatile storage → can retrieve info after being power cycled

Disk - writing data to disk ~~persists~~, persists even if the DB itself crashes / faces some outage

eg:- ~~not~~ excluding extreme issues, any file you save on PC, persists even after PC ~~shuts~~ shuts down or crashes.

→ needs const. power to retain data

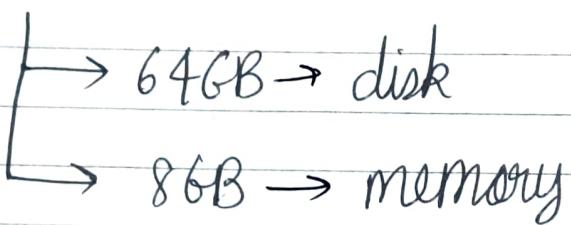
→ RAM - volatile storage

Date: _____
M T W T F S S

Memory - Data stored in this does not persist if DB goes down

Eg:- any data stored in ^{variable} your server's DB does not persist if DB server goes down

Eg:- In a mobile, say 64 GB storage, 8 GB RAM



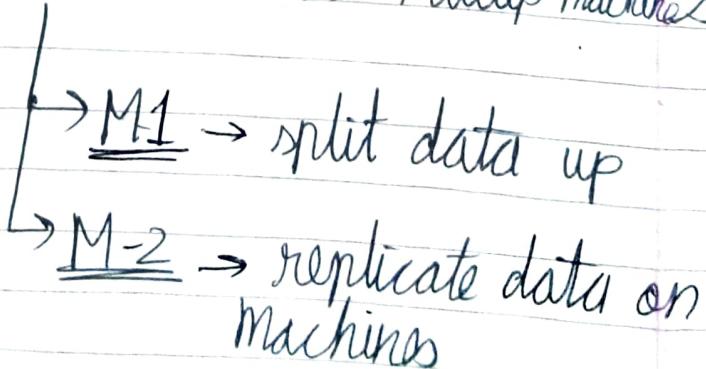
NOTE

Reading data from memory is faster than reading data from disc.

There are a lot of DB offerings to give options based on performance, data security

When DB goes down, ~~because~~ depending on how critical a part of the entire sys the DB is, does the system crash or persist through

Distributed storage → storing data on multiple machines



```
JS server.js × JS http_request_example.js ●
```

```
JS server.js > ...
1 const express = require('express');
2 const app = express();
3
4 app.use(express.json());
5
6 app.listen(3000, () => console.log('Listening on port 3000.'));
7
8 app.get('/hello', (req, res) => {
9   console.log('Headers:', req.headers);
10  console.log('Method:', req.method);
11  res.send('Received GET request!\n');
12 });
13
14 app.post('/hello', (req, res) => {
15   console.log('Headers:', req.headers);
16   console.log('Method:', req.method);
17   console.log('Body:', req.body);
18   res.send('Received POST request!\n');
19});
```

...

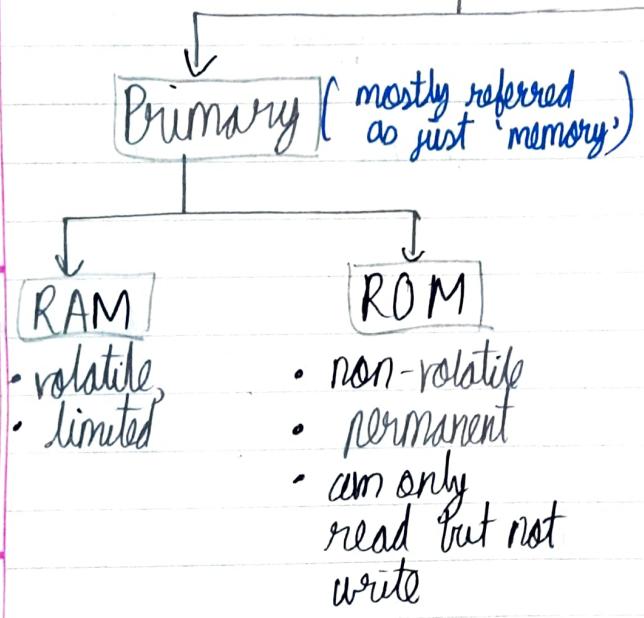
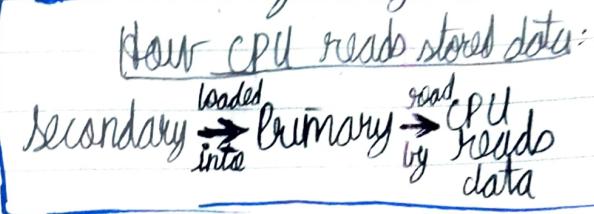
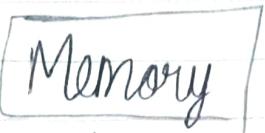
```
Clements-MBP:network_protocols clementmihaiilescu$ node server.js
Listening on port 3000.
Headers: { host: 'localhost:3000', 'user-agent': 'curl/7.54.0', accept: '*' }
Method: GET
Headers: {
  host: 'localhost:3000',
  'user-agent': 'curl/7.54.0',
  accept: '*/*',
  'content-type': 'application/json',
  'content-length': '14'
}
Method: POST
Body: { foo: 'bar' }
```

```
Clements-MBP:network_protocols clementmihaiilescu$ curl localhost:3000/hello
Received GET request!
Clements-MBP:network_protocols clementmihaiilescu$ curl --header 'content-type: application/json' localhost:3000/hello --data '{"foo": "bar"}'
Received POST request!
Clements-MBP:network_protocols clementmihaiilescu$
```

Consistency - A concept in storage referring to the staleness or up-to-dateness of data

e.g. - on accessing data from a DB, how fresh/up-to-date is the data that you get

NOTE :



Secondary (a.k.a **STORAGE**)
data stored permanently unless deleted



HDD
data is recorded magnetically on surface. accessed by spinning disc
slow, cheap

SSD
stored same as HDD but accessed using RAM module
fast, costly

- ROM contains a prog. called BIOS (Basic I/O sys) which microprocessors (i.e. computer CPU) to load OS from HDD into RAM whenever PC is turned on. Newer motherboards use UEFI (Unified Extensible Firmware Interface)

- Analogy :-
 - (i) MEMORY \rightarrow a work desk
 - (ii) storage \rightarrow Secondary memory \rightarrow desk drawers where files are stored
 - (iii) memory \rightarrow Primary memory \rightarrow table top of desk
 - (iv) Prog in memory \rightarrow tools, things on desk (easy, fast to access)
 - (v) Prog in storage \rightarrow files in desk drawers that you gotta put in table top to access (slow to access)
 - (vi) loading a prog. \rightarrow open a drawer, remove read file and put it on table top to use it

Latency & Throughput

A)

Latency and Throughput → 2 most imp measures of the performance of system

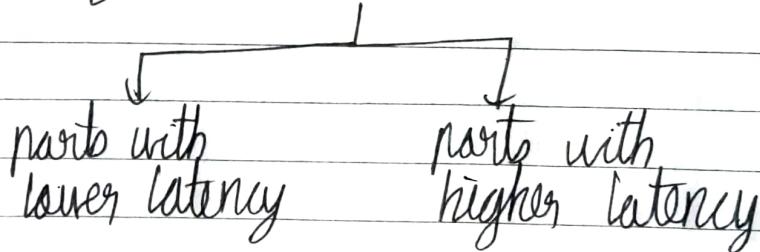
LATENCY:

- i) Latency is basically how long it takes for data to traverse a sys. i.e. get from 1 pt. in sys to another pt. in sys.

eg:- latency of network req : how long it takes for 1 req. to go from client to a server and then the processed response to go from server to client

eg 2:- time reqd. by a reader to read data from disk

- ii) Diff parts of sys hr diff. latencies
:- While designin a sys we'll face a trade-off while optimizin it cuz we'll have



- iii) Some comparisons

Part of sys + its f ⁿ	Latency
Readin 1MB from memory (RAM)	250 μS
Readin 1MB from SSD	1000 μS
Sendin 1MB over 1Gbps network	$10^4 \mu S$ → (sendin to comp- uter need to be considered) i.e. dist. is not considered
Readin 1MB from HDD	$2 \times 10^4 \mu S$
Sendin packet ($\approx 1055 B$) over network from California to Netherlands and then back to California	$15 \times 10^4 \mu S$

eg: of sending data over a network \rightarrow API
 eg of reading data from memory \rightarrow reading a variable
 in code

Takeaways:

- i) Dependin on network ^{and your PC hardware}, sometimes ~~sometimes~~ sendin, receivin data over network is faster than reading data from HDD
- ii) Sendin data around the world takes a lot longer than any other meth.
 Reason: req. gonna ~~be travel~~ get converted into ^{small} packets (into binary data) \rightarrow converted into freq modulated radiowaves \rightarrow sent to cell towers thru cables \rightarrow bounced to satellite \rightarrow passed around the world thru satellite comm.
 \rightarrow passed to destinaⁿ \leftarrow passed to all towers and reconvereted to destinaⁿ \leftarrow back to original form
- iii) Optimizin a sys \rightarrow \downarrow its overall latency
 eg:- Video games gotta have ^{to have} really low latency and lag \rightarrow delay in acⁿ passed from 1 user \rightarrow server \rightarrow receiving user. These acⁿ's are passed as network req's to server so if you're pretty far away from server, your PC will take more time to make network req/receive responses from server

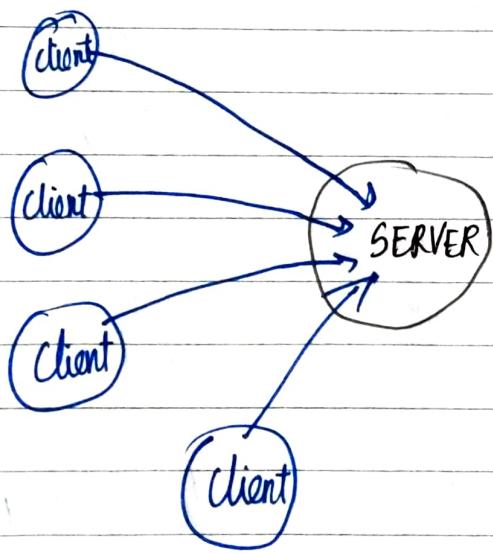
eg:- Websites → It's peace if they have low latency because in most cases, their priority is getting the info displayed to be accurate, uptime to be cont. 24x7

THROUGHPUT (TP)

Throughput is basically how much work a machine can perform in a given amt. of time → how much data can be transferred from 1 pt. in sys to another pt, in given time

unit :- bytes / sec

eg:- 1 Gbps network → network can support 1 Gb per sec.



TP is how many reqs (each hvin some data) can this server handle in given time → how much data it can let thru per sec.

To optimize sys. → pay to ↑ ~~latency~~ TP

| Case(i)
in case just ↑ TP doesn't solve prob. as you might have a server hvin 10^3 or 10^6 req. issued to it per sec so no matter how much we ↑ TP, we'll still hr a bottleneck (only some data is let thru server at server)

Case (ii) soln : have diff. servers for req.s so they don't clog at bottleneck

IMP : latency and throughput are not correlated

eg:- you might have parts of sys with really low latency (fast data transfer)

→ but then if you hr a part with < less TP, then our advan. due to low latency in other parts gets cancelled out as req.s gotta wait (clog) at this part cuz of low TP

∴ You can't make assumptions on latency or TP based on each other.

They affect each other but don't determine each other

AVAILABILITY

Fault tolerance → How ~~tolerant~~ resistant is a sys to failures. What happens if a server in sys fails? DB fails? → Is sys gonna still be operational

Availability → % of a ^{given period of} time (eg:- month, year) for which the sys. is atleast operational enough to get its primary fns satisfied

Less availability → less uptime of a sys. during given time period

- lose customers (existing)
- lose the prospect of getting new customers
- ↓
bad publi.
- ↓
money lost

sys.s that gotta hr. high availability

- DUE TO THEIR FN
eg:- sys. supporting airplane software

any amt. of downtime costs a lot

- DUE TO THE MAGNITUDE OF PEOPLE ACCESSIN
eg:- YouTube, Facebook

- CLOUD PROVIDERS

All services, platforms (dependent on cloud providers) businesses get heavily affect if cloud provider sys. fails

Service = sys.

Date: _____
M T W T F S S

eg of cloud providers : Azure, AWS, Google Cloud, Oracle Platform

Measuring availability (albt)

Availability is measured in terms of % of sys.'s uptime in a given year. (90% albt \rightarrow sys. is up 90% of the time in a yr)

Most powerful services, sys.s have ~~are~~ to be really HA. i.e. 99.9%, 99.99% etc.

Nines \rightarrow Percentages with the no. 9

eg:- 99% albt \rightarrow 2 nines of albt

Albt.	Downtime / yr
2 nines (99%)	87.7 hrs
3 nines (99.9%)	8.8 hrs
4 nines (99.99%)	52.6 mins
5 nines (99.999%)	5.3 mins

Highly available sys \rightarrow albt $>$ 5 nines albt.

SLAs and SLOs (albt. guaranteed explicitly)

SLA \rightarrow Service Level Agreement, \rightarrow Agreement b/w service provider and customers/end users of service that guarantee customers amount of albt, error free utility and some other objectives (known as SLOs \rightarrow service level objective), failing which service the service provider pays back the customer some % of their fees.

S L A

Date: _____
M T W T F S S

eg:-

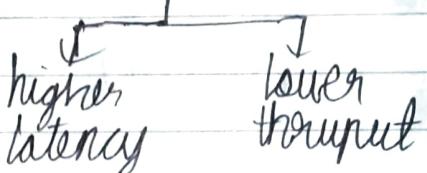
AWS (provider)

Netflix (user)

- SLO 1 → 1) Netflix get 99.999% abt.
 SLO 2 → 2) Netflix get x errors while using my service

If AWS fails to meet SLOs, it returns 10% of the service fees Netflix pays to it every month

NOTE: Achieving higher abt. is really difficult and may accompany tradeoffs like



As a sys. des. gr., you gotta decide which part of sys you want high abt. and which part you don't really need HA.

eg:- Stripe payment service for business

→ Core services → handling payments, charging customers
 These gotta have high HA else Stripe, as well as its client businesses, platforms lose lotsa money

→ Business monitor dashboard → used by client platform to use business stats

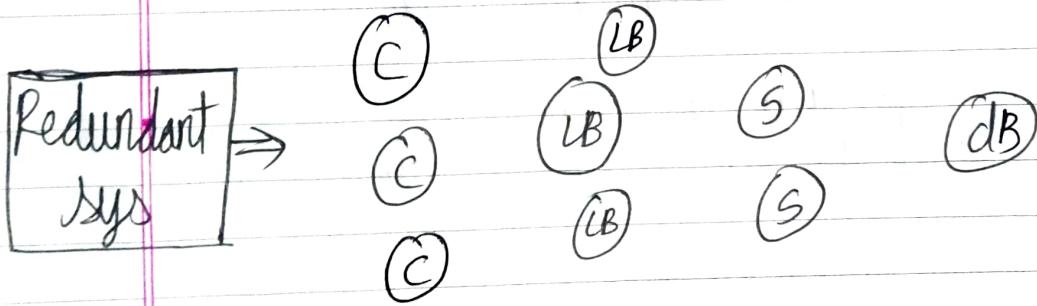
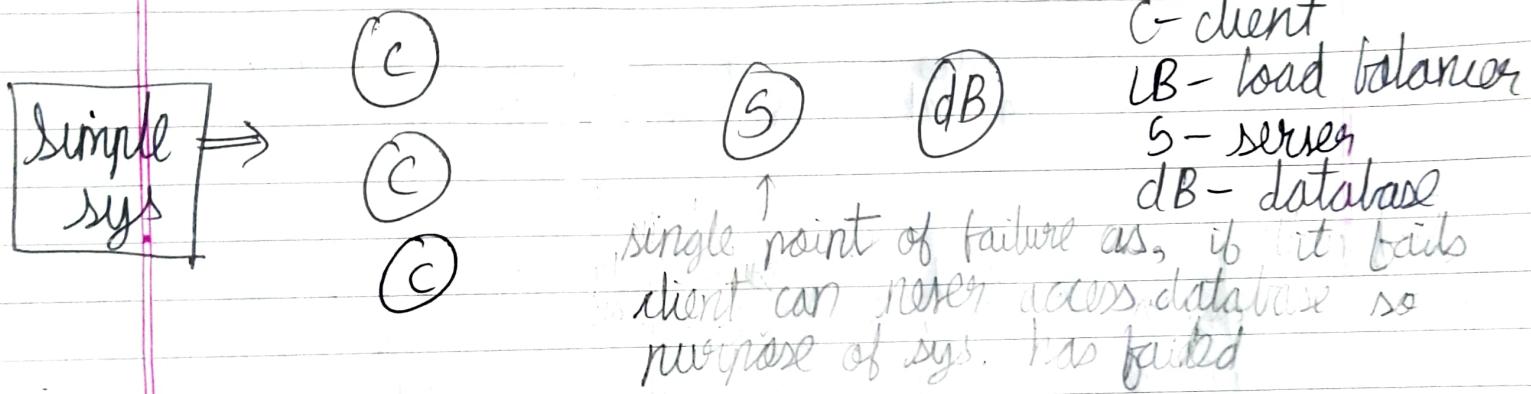
Doesn't really require HA, something that, if fails, doesn't cause catastrophic losses

How a sys. is made "AVAILABLE".

1

Single pts of failure (parts of the sys, of which, on failing cause entire sys. to fail) should be removed → for this we use **redundancy** while ~~sys~~ designing the sys.

Passive Redundancy → act of duplicating or triplicating or multiplying even more, certain parts of sys.



Load balancer → to ~~reduce~~ equally distri. load across all servers (not 1 server from becoming too overloaded & acting as single pt. of failure)

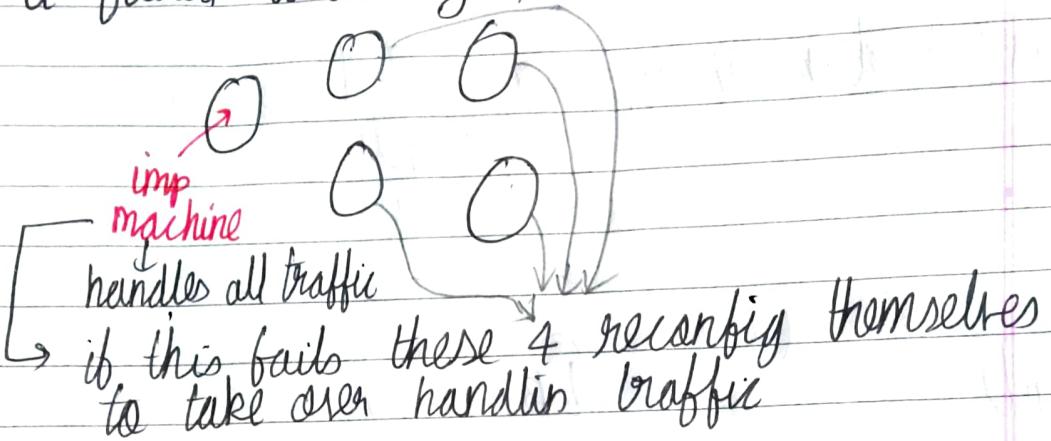
Multipled load balancers in this sys to avoid 1 LB from getting too overloaded

→ duplicated

Eg:- twin engine sys. in airplane

Active redundancy - Multip machines that work together in such a way that only 1 or few of the machines are gonna be typically handlin traffic or doin work (imp machine). If an 'imp machine' fails, the others machine know it failed and they'll take over.

eg:-



2]

Hare rigorous processes in place to handle to handle sys failures. cuz its possib. that these failures might req. human interver"

∴ Have these backup op's in mind while sys. des. to get a crashed sys., up in runnin in proper timeframe.

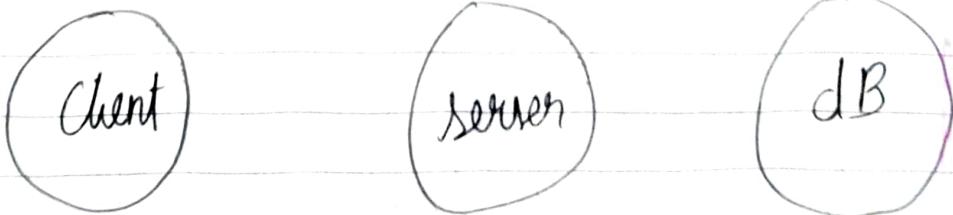
CACHING

- 1] Caching in algos → To avoid doing some ops especially computationally complex ops. that take a lotta time multiply times
∴ caching used to improve time complexity (T.C.)
- 2] Caching in sys. → Used to reduce (improve) the latency of a sys. (speeds up the sys.)
- 3] CACHING → Storing data in a locⁿ that's diff from the original data locⁿ s.t. it's faster to access this data from memory.
It's a way ~~to~~ to design a sys. s.t. if we were originally using ~~some~~ ops. or data transfers that take a lotta time (e.g.: Network requests) the sys. is designed is s.t. we do diff types of ops. or data transfers that are faster.

I

Caching in diff places in a sys.

e.g:-



Regular op. → client makes net. req. to server to get some data from dB and then data transferred from dB → server → client

1) Caching at client level - Client caches a data value so it no longer needs to go to server to retrieve it

2) Caching at server level - You want client to always interact with server but maybe the server doesn't need to go to dB for data retrieval each time

- Go the dB once → get data → cache in server → transfer data to client as per net. req.

3) Cache in b/w components - eg:- cache in b/w server and dB

not in our control when we des. sys. → 4) Caching at hardware level - eg:- CPU cache - make it faster to retrieve data from memory.

5) Caching occurs by default, at many other such kinda b/w. of sys.

II

Advantages / Utility of Caching

Caching speeds up sys. → 1) Avoid large magnitudes of net. reqs. → If the client wants to access the same data again and again, no new net req. (client → server → dB → server → client) done each time as we cache reqd. data so client doesn't need to go to dB each time.

Caching done at client or server lvl.

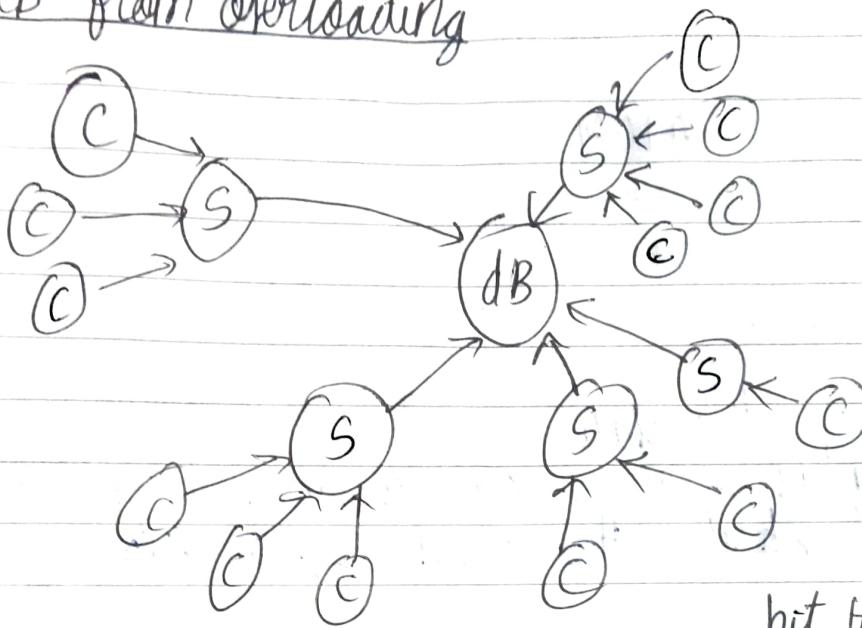
Caching speeds up sys. → 2) Avoid redoing a computationally, long, bad T.C., complex operatn each time
eg:- say on net. req. from client → server

everytime, this time - consuming algo runs.
To avoid performin this op. multiple times,
caching is done.

3) Prev. dB from overloading

dB
optimized

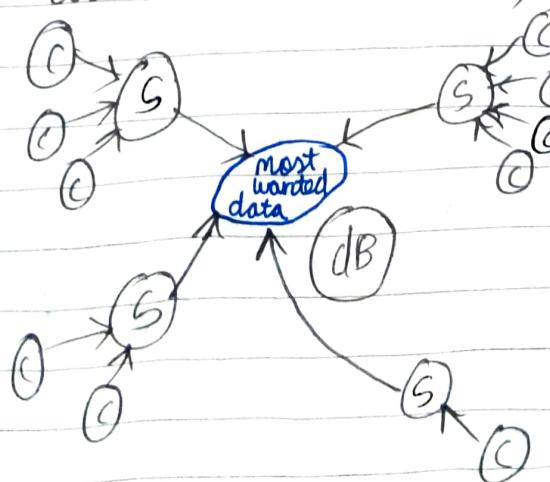
e.g:-



hit the dB

Multip clients across multip server, with the same net. req. e.g:- Each client makes ^{net.} req. to server to get Instagram profile of a popular celeb.

Indir. req. are fast, so speed is not the concern.
We ~~use~~ use caching to prev. same data being
read from dB 10^6 times (or somethin like that)
 \therefore a cache is created, outside dB, for that data



If we don't cache,
too many req.s
to dB for this data +
in gen. other data
accumulate and
overload dB

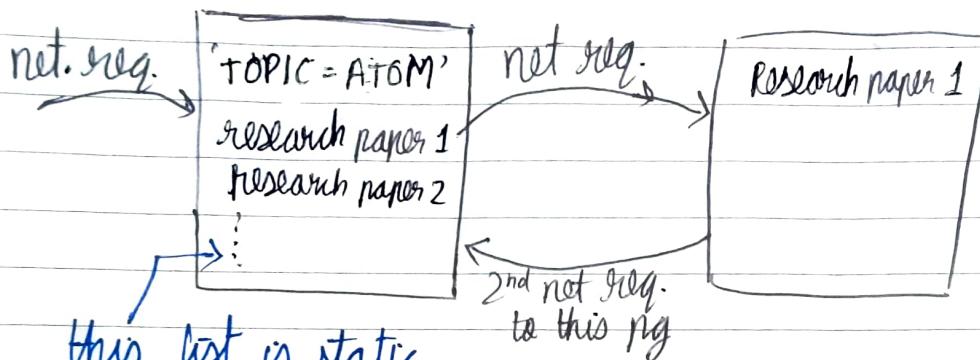
$k-v$ = key value

Date: _____
M T W T F S S

III

CACHING in acⁿ

- 1) Reasons state components ^{immutable} ^{or content} in a website
 e.g. - ResearchGate : every once a yr papers published on site



this list is static content so if we cache it after 1st net req. to this pg of website, the 2nd net req. to list pg will be a lot faster

2) Runnng code on algoswp. (AE)

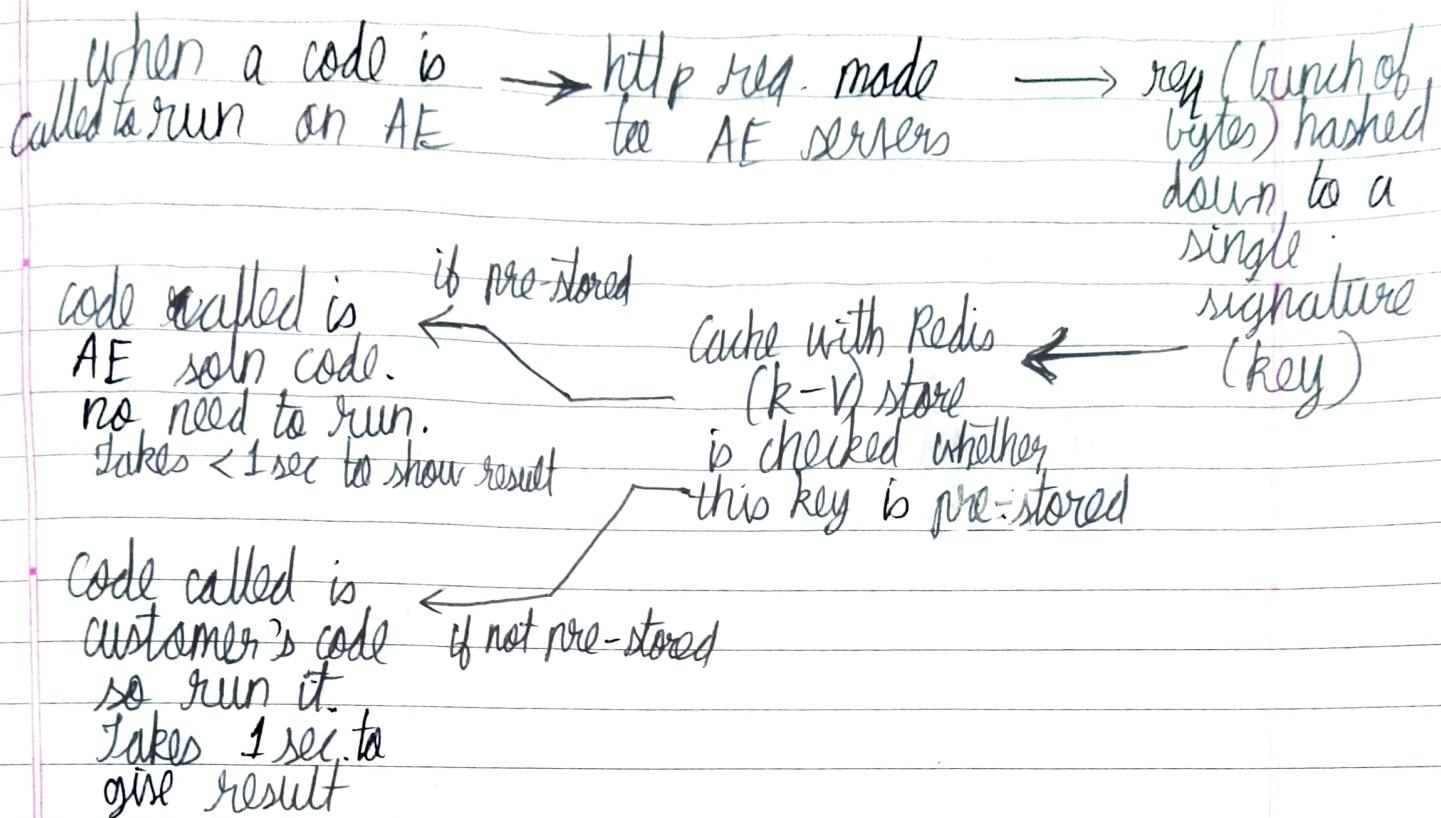
- i) Runnng, testin code takes ≈ 1 sec.
- ii) for customer code \rightarrow runnng, testin code reqd.
 for AE soln code \rightarrow w.r.t this code is correct code i.e. it'll pass all test cases so no need to run, test

Cache stored $\xrightarrow{\text{server level}}$ OR $\xrightarrow{\text{detached from server is an indep. comp.}}$

Redis (popular in-memory dB) \Rightarrow a k-v store used for caching.

All AE solns. codes pre-stored with a unique 'key' ~~in~~ in this Redis cache

Caching in acⁿ



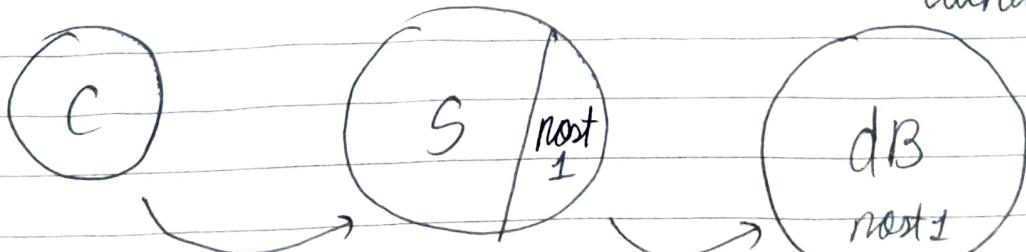
Until now we discussed caching for a sys. that needs to read data.

IV Caching and Sys. where users can read, write, edit data.

e.g.: - Readin, writin, editin posts of FB or linkedIN

1) When you write a post to FB, say that using not reqs. that post is stored in dB

→ cached in server



2 sources of truth for given post → cache in server → dB

2) When you now wanna update your post, depending on type of caching in server post stored in the

```
js server.js > app.get('/nocache/index.html') callback > database.get("index.html")
1  const database = require('./database');
2  const express = require('express');
3
4  const app = express();
5  const cache = {};
6
7  app.get('/nocache/index.html', (req, res) => {
8    database.get('index.html', page => {
9      res.send(page);
10   });
11 });
12
13 app.get('/withcache/index.html', (req, res) => {
14   if ('index.html' in cache) {
15     res.send(cache['index.html']);
16     return;
17   }
18
19   database.get('index.html', page => {
20     cache['index.html'] = page;
21     res.send(page);
22   });
23 });
24
25 app.listen(3001, function() {
26   console.log('Listening on port 3001!');
27 });
```

```
js database.js > ...
1  const database = {
2    ['index.html']: '<html>Hello World!</html>',
3  };
4
5  module.exports.get = (key, callback) => {
6    setTimeout(() => {
7      callback(database[key]);
8    }, 3000);
9};
```

DB → updated instantly
DB → updated async. ly.

Date: _____
M T W T H S S

Write-thru cache

- 1) dB and cache both updated instantly when you edit post

Write-back cache

- 1st cache is updated instantly
Then dB is updated async.

→ updated after a period of time e.g. 5 hrs, 30 days etc.

→ once cache is full and you gotta evict stuff outta cache

- 2) Takes time as you still gotta go to dB and come back for an edit

Takes less time as only cache update at the instant of edit

- 3) Edit stored in 2 sources of truth at the instant of edit

Edit stored only in cache at instant of edit
∴ if somethin happens (e.g.: delete) to our cache before dB is async. ly updated, we lose ^{the} edit.

V

Staleness

eg:- YouTube comments see?

You don't wanna update data asynch as it might happen that:

- someone posts a comment (SERVER 1)
 - other person reads it from diff. server (SERVER 2)
 - 1st person edits comment (but dB not updated) and his server's cache updated asynch yet
 - other person still sees original comment (unedited) and replies to unedited comment (STALE DATA)
- as SERVER 1's cache ↑
 by SERVER 2

Cache in server 2 is STALE : it hasn't been updated properly

Soln: Move cache outside servers into a common comp. for them all outside dB
 ∵ Each time cache is updated from server 1, server 2 accesses this latest, updated cached data

NOTE: While sys. design just see if that part of sys. is ok with staleness or not, and accordingly choose how you gonna cache the data

eg:- YouTube → view count → lit bit of staleness on video → is fine
 → comment sec'r → even a lit bit of staleness is NOT fine

VI

Date: _____
M T W T F S S

NOTE:

WHEN TO USE CACHING?

- ✓ 1) Static (immutable) data → Caching is easy to implement, fast while using
eg:- Research papers website
- ✗ 2) Dynamic (mutable) data → Caching is tricky to implement as data stored at 2 diff loca's and you gotta make sure data in both loca's is in sync.
- eg:- ^(FB) Facebook wall → not cached
- 1st net req. → my FB wall → friend's profile
net req.
2nd net req.
- On Facebook website, the data you see on your wall after 1st and 2nd net req. to FB wall page is diff.
- ✓ You're doing only 1 of reading or writing data
single thing eg: either
- ✗ You're doing multiple things → reading, writing, editing data
- ✓ You don't care abt staleness of data or if you can invalidate (get rid of) stale data in caches in a distributed manner when dealing with a distributed sys.

VII

CACHING - EVICTION POLICIES

We might wanna evict (delete) data from cache for several reasons. Eg : data is stale, cache's storage is full.

To evict data from cache, we follow certain rules a.k.a "evic" policies

- 1) LRU policy → Least Recently Used (LRU)
pieces of data in a cache are removed if ~~you~~ you have some way of tracking what pieces of data are LRU. cuz we ~~care~~ care least abt em.
- 2) LFU policy → Least Freq. ly Used (LFU) pieces of data evicted
- 3) LIFO or FIFO → Data ejected on a Last In First Out or First In First Out basis
- 4) other ways.

Choice of evic" policy depends on kind of product you wanna do sys. des. for

PROXIES

PROXY:

Just a machine-server set up to hide the IP address and other details of an interacting machine
 eg:- client server etc.)

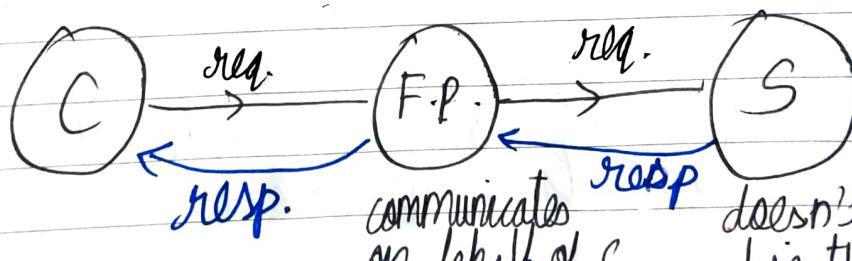
2 types (primary) :-

- i) Forward proxy usually referred to as - 'proxy'
- ii) Reverse proxy

Forward Proxy and Reverse Proxy (F.P)

1] Forward Proxy (F.P) → on client's team

- i) It is a server located b/w client / set of clients & servers / set of servers
- ii) It acts on behalf of clients / set of clients by hiding clients IP address from server and instead supplying its own IP address.
- iii) When a client sends a req. to the server and this ~~req.~~ F.P. has been properly configured by client, the req. goes as



Server has no idea abt the client and that an doesn't get req. F.P. directly from c. It gets it from F.P.

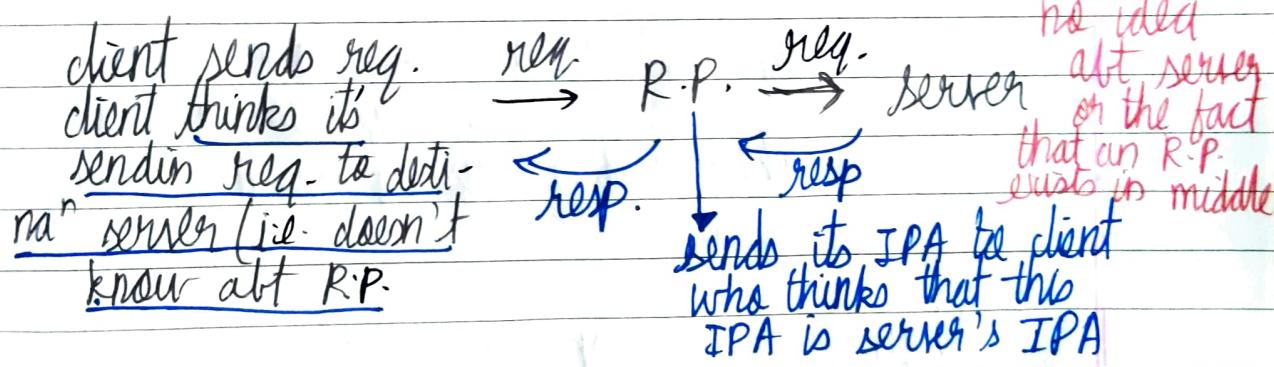
- iv) F.P. serves as a way to hide the identity of client ~~by issuing req. from server~~ by changing source IP address sent in req. to its own IP address (- IPA) instead of client's IPA
 Eg:- VPNs

v) Some types of ~~proxies~~ F.P.s make client IPA visit to server in some way, but, typically original source IPA is replaced

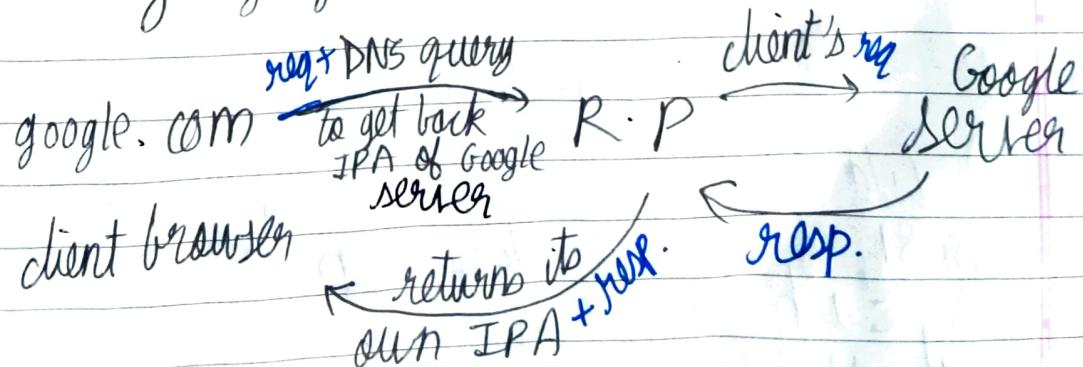
NOTE: VPN → a fraud proxy b/w client and server that hides the client's identity
 i) client can access servers restricted to it.
 Thus, client can, for eg, access a website not avail. in his/her country but avail. in other.

2] Reverse Proxy (R.P.) → on server's team.

- i) A server b/w client / set of clients and server / set of servers
- ii) When a client interacts with a server. eg:- sending a req.



eg:- May google.com sets up an R.P.



Use of R.P.s → Powerful tool for sys. des.

- 1) eg:- Config. R.P. to filter out reqs. you want your sys to ignore
- 2) eg:- Log stuff, gather metrics → can be done by R.P.
- 3) eg:- R.P. can also cache stuff (e.g.: HTML pages)
Thus, server doesn't get bothered a lot
- 4) eg 4:- As a load balancer → a server that can distri req. load amongst a bunch of servers

NOTE: 1] ~~Malicious~~ Malicious client → sends a fuckload of reqs. to a server to bring it down.

Now R.P. acts as load balancer will distri these reqs across all servers, thereby safeguarding sys. from malicious clients, viruses etc.

2] NginX is a popular web server that can be used as an R.P.

eg:- Code on NginX page:

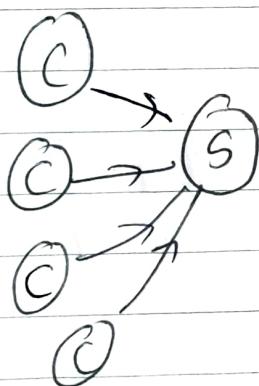
- 1) We set up an R.P. for any req. comin to port 8081 of our ~~web~~ service
- 2) Each time a req. is directed to the endpoint '/', at port 8081, the req. header on the req. known as 'systemexpert-tutorial' and set it to 'true'
- 3) Then we gonna bind this req. to the server that pts. to localhost:3000 (its name is Nodejs-backend)

Load Balancers

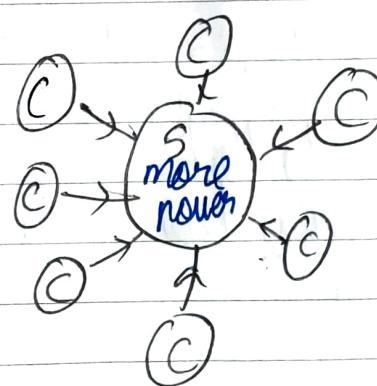
Load Balancer → (i) A server that sits b/w a set of clients and set of servers which distributes load / routes traffic evenly across all servers in our set.

- (ii) This prevents 1 server from getting too overloaded by too many client reqs.
- (iii) Also clients do not know abt load balancers' existence so load balancer = reverse proxy

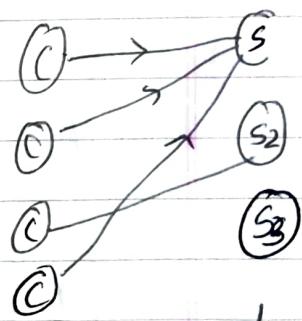
Case 1: Server overload due to lots of reqs.



Case 2: Vertical scaling: power of servers ↑ but still lots of reqs, so overlod

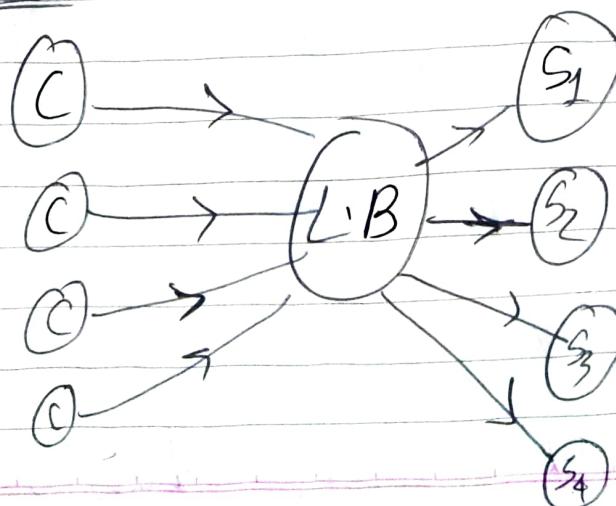


Case 3: Horizontal scaling
Multiple servers but still most client reqs to 1 server.



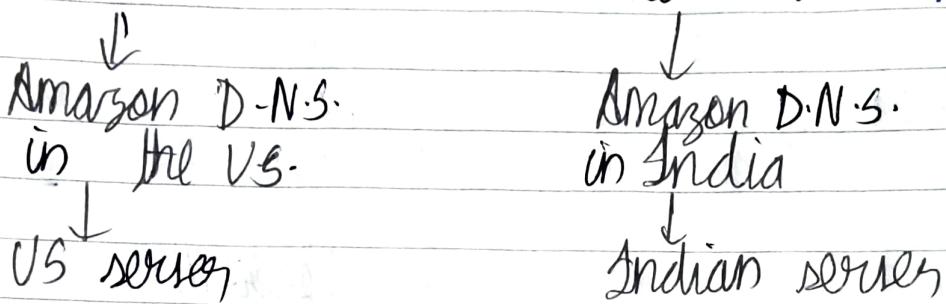
S1: overloaded

Case 4: Load Balancer



- ① Req's evenly distri. → No overload of servers
- ② throughput ↑
- ③ Latency ↓ (no server clogged)

Other places where load balancing can be done:



Q2: Netflix VS ~~VS~~ ~~Notif~~ Netflix India

e.g. On opp page \rightarrow we are accessing diff. IP addresses of google.com i.e. diff. servers and we get allotted the server related to our specific IP address (abbr. IPA) by the load balancer.

NOTE

Hardware LB → physical machine dedicated to load balancing

Software LB → ~~more~~ not physical machines
(e.g. - a cloud services LB → ~~cloud server~~)
use: More power, customization

SOFTWARE LB.

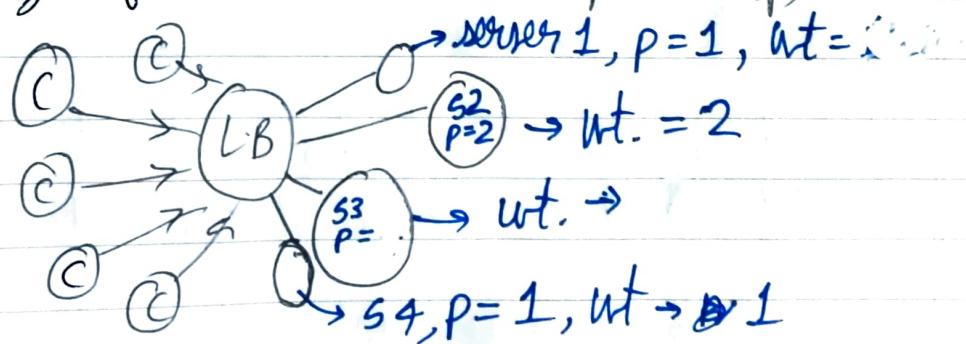
- 1) How does LB know which servers it can access?

i) sys. designer config. LB and servers to know abt each other

i) sys. designers config. L-B and servers to know abt each other

- ii) On adding / removing old servers, it registers / deregisters itself w.r.t. LB
- 2) Ways to select servers to direct load to
- ① Random redirect → ~~No logic used. Just random selecⁿ of server from available servers~~
by chance, one server may get overloaded
 - ② Round Robin → LB goes thru all servers in order and allot a batch of reqs. to a server and so on until all servers have equal load so, for next batch of reqs., it'll again start from 1st server
- ↓
server's power not optimally used.

Let size of server \otimes & its power (P)



- ③ Weighted Round Robin : LB goes as per round robin approach but if a weight is added on 2nd server, LB will send it a couple more reqs.
 \therefore Power of each server \Rightarrow optimized

Request Batch Round Robin(RR) Wtd. R.R.

1	S1	S1
2	S2	S2 } avg wt = 2
3	S3	S2 }
4	S4	S3 } wt = 3
5	S1	S3 }
6	S2	
7	S3	S4

(3) Performance/Load Based: LB will do performance/health checks on servers

- how much traffic a server is handling at any given time
- how long a server is taking to respond to traffic
- how many resources server uses
- etc.

Based on these checks, it allocates reqs. (more reqs. to low : high performance servers)

(4) IP based server selec: On getting requests from clients, L.B. hashes the IP address of the client and depending on this hashed value, the client's req. gets directed to a specific server

Use :- Caching

e.g.: - You're caching results of certain reqs.

in your servers, it's helpful to redirect the client to a series in which, the ~~req~~ response to the request ~~is~~ has already been cached

- eg:-
- i) Client 1 searches for "cheese pizza"
 - ii) ~~LB~~ knows that, the response to this request is ~~is~~ in ~~server 1~~
 - iii) LB hashes down IPA of client to say "1" → basically, a VID
 - iv) Now, it sees that in server S3 a result for ~~req~~ a prior req. made with VID = 1 is cached and ready
 - v) LB redirects client 1's req. to server S3
 - v) Client 1 gets his/her response super fast.

⑤ Path-based → LB selects ^{server or set of servers} based on path of ~~the reqs~~.

eg:- AlgoExpert has servers S1, S2, S3, ..., S25

All reqs related to payments } → S1, S2, ..., S9

All reqs related to running code → S10 - S18

All reqs related to watching vids → S19 - S21, S22, S24, S25

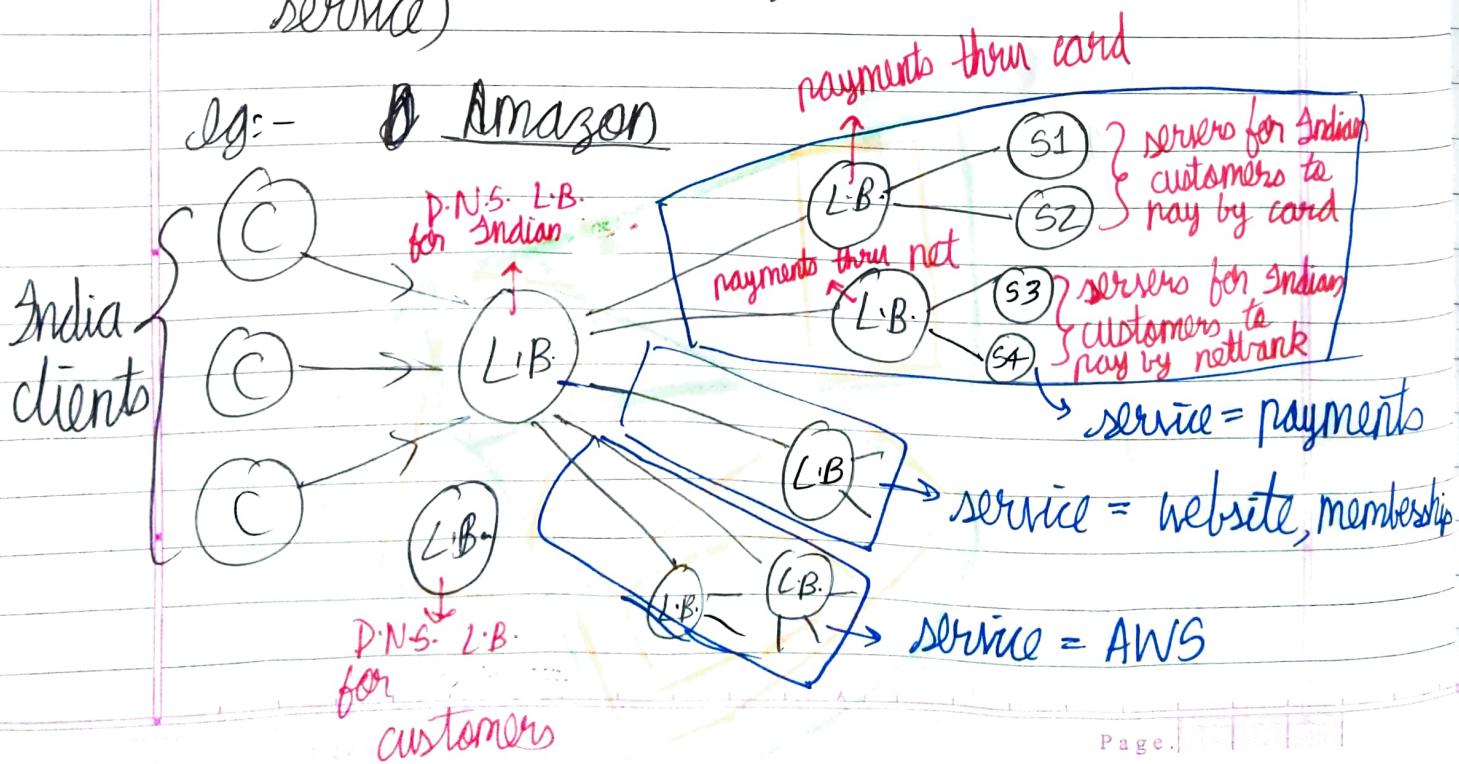
All reqs related to "about" page → S23

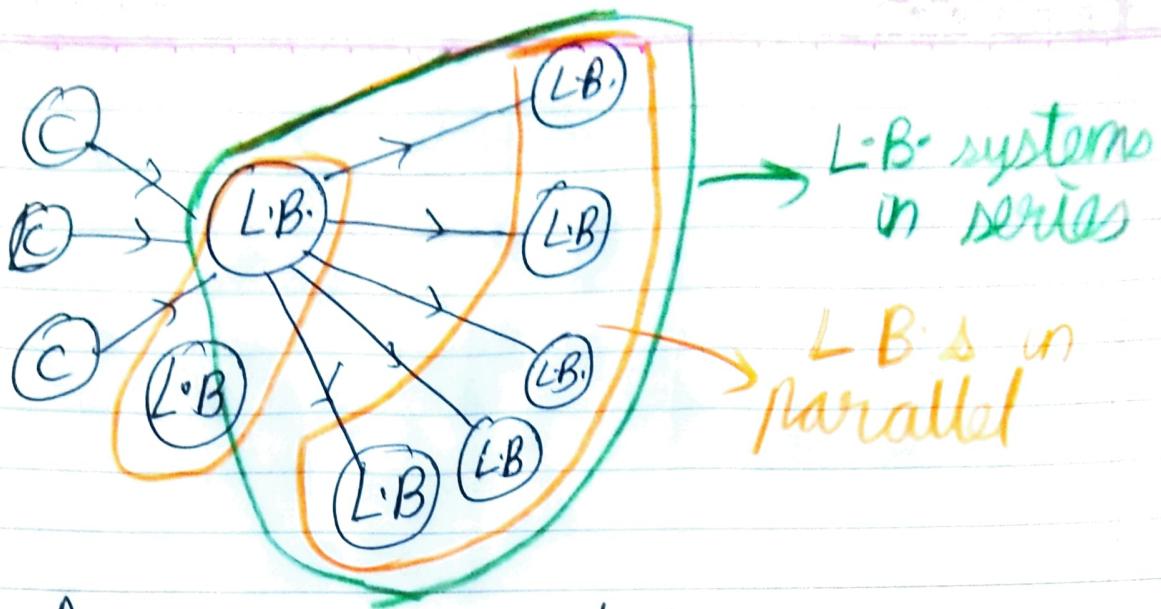
Use:

- 1) Any deployment of big change to any service (eg:- running code service) affects only the servers allotted to this service.
- 2) If one set of servers for a service start failing, atleast other services are still up and running. (unaffected)

MULTIPLE LOAD BALANCERS

- 1) In series → To keep divide a batch of reqs into smaller batches by based on some criteria (path eg:- path) on the req.
- 2) In parallel → To ① just prevent 1 L.B. itself from getting overloaded or ② Redirect to diff. servers based on req. criteria or server criterial (eg:- performance, diff kinda service)





Amazon eg:- Just L-B-S

code on opp. page → shows used R.R. selectⁿ meth.



Hashing: An ~~op~~ opⁿ ie. performed to transform an arbitrary piece of data (str., arr, or other data type or struct) into a fixed size value, typically an int.

e.g. in sys. des:- Arbitrary Data → IP address, HTTP request, username,

1

Need for Hashing

Clients

(C1)

(C2)

(C3)

(C4)

L-B.

LB.

Servers

(S1)

(S2)

(S3)

(S4)

① While servers select strategies such as round robin or random get your work done in normal scenarios, they fail to optimize cache hits in a in-memory cache scenario.

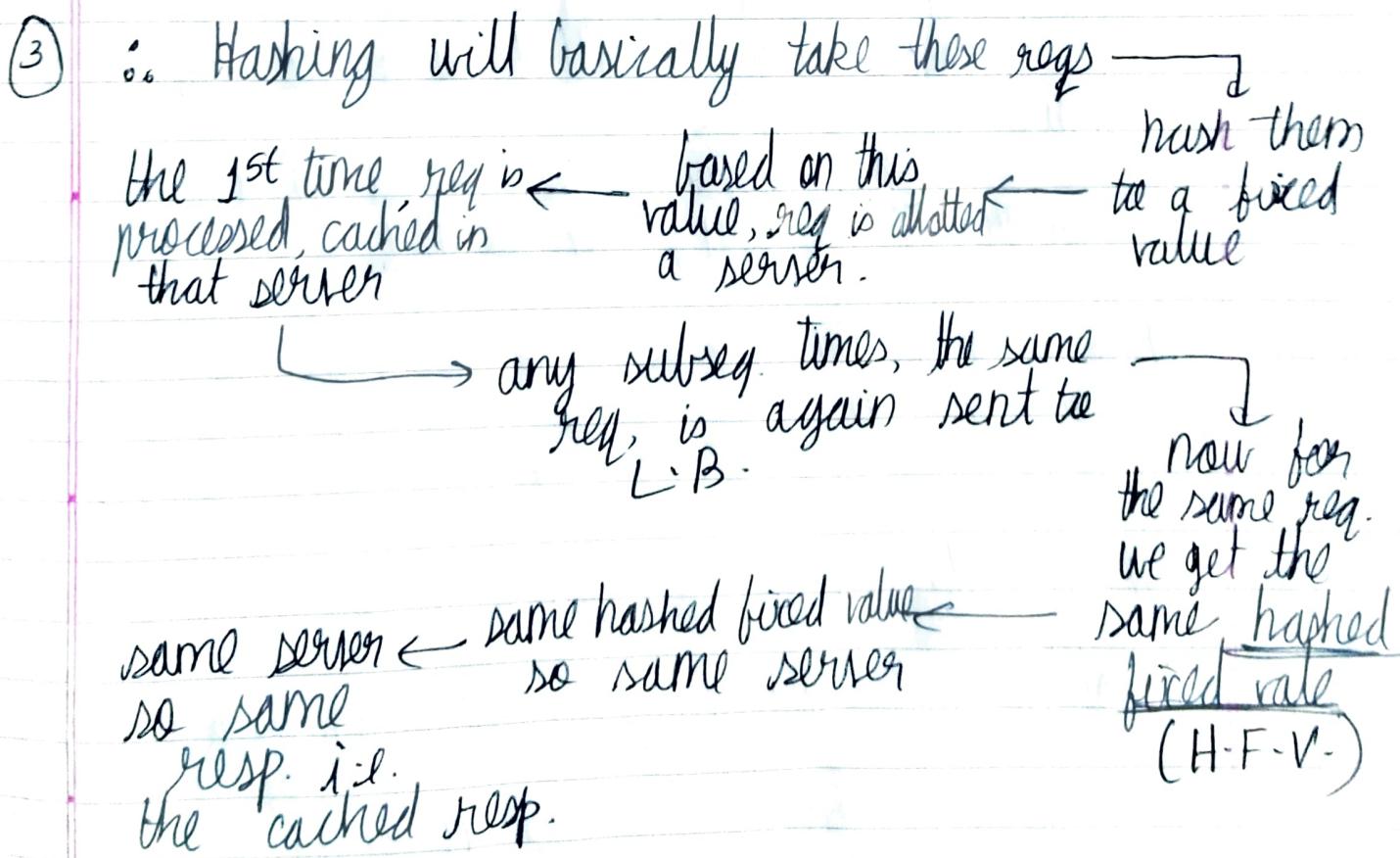
② In Memory Cache

a) Cache hit → no. of times a req. received a cached res. from our sys.
(faster)

(SSS)

- 6) For this, round robin or random server select strategy don't work as they do not guarantee you that if the 1st time a client makes a req. and it gets ~~processed~~ a processed and cached in S1, the 2nd time, when client makes the same req., he/she will be redirected to S1 only
∴ chances of cache hits ↓

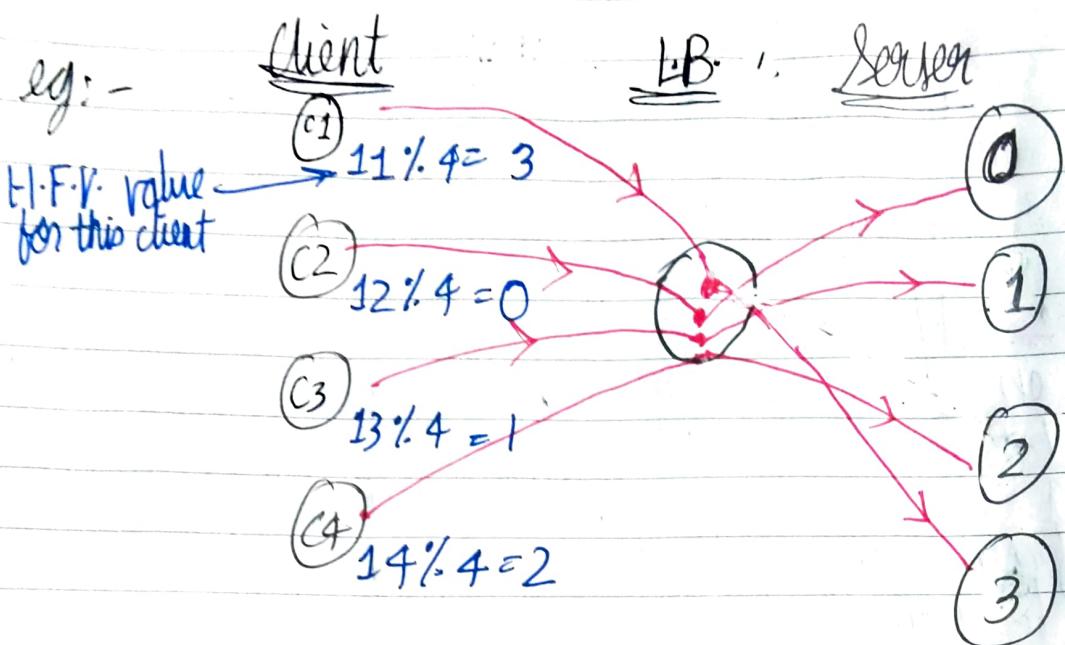
- c) ∴ If your sys. des is st., it relies heavily on in-memory cache in its servers, but your L.B. is receiving req.s from clients following a SSS that doesn't guarantee that reqs. from some client OR same reqs. will be routed to the same server every time, then, your in-memory caching sys. falls apart



2] Simple Hashing \rightarrow

$$\text{server allotted to client (SATC)} = (\text{H.F.V.}) \% \text{ (no. of servers)}$$

e.g.:-



- ① Uniformity in Hashing Function (H.F.) \rightarrow No. of collisions
- i.e. same server allotted to a no. of clients should be minimized.
 - i.e. clients should be distributed across servers equally, uniformly.

NOTE: Typically ~~you~~ you never write your own H.F., You use ~~the~~ pre-made industry-grade H.F.s such as MD-5 or SHA-1 or Bcrypt

- ② Cache hits \uparrow for same reqs

- ③ Problem :- Adding or removing ~~one~~ (or any) of servers completely fucks up cache routes (C.R. \rightarrow a route which a req. must follow to get its matching cache a.k.a. route for cache hits)

$\text{MEM} \rightarrow \text{Memory}$

Date _____
M-T-W-T-F-S

e.g.: - Adding a server to cache eg. after in-mem. caches have been created for reqs which already in their resp. SATC. Say server S5 is added now

Client	H·F·V	$n_i = 4$	$n_f = 5$	SATC _i	SATC _f
C1	11	4	5	3	1
C2	12	4	5	0	2
C3	13	4	5	1	3
C4	14	4	5	2	4

n_i = initial no. of servers

SATC_i = initial order (the one with in mem. cache) allotted to a req. from this client

- As visible, no new reqs (i.e. reqs after adding of S5) are going on their cache route.
- ∵ They won't reach server with cached res.
- ∵ Cache hits ↓

3] Consistent Hashing

SAP → Server Allotment Procedure

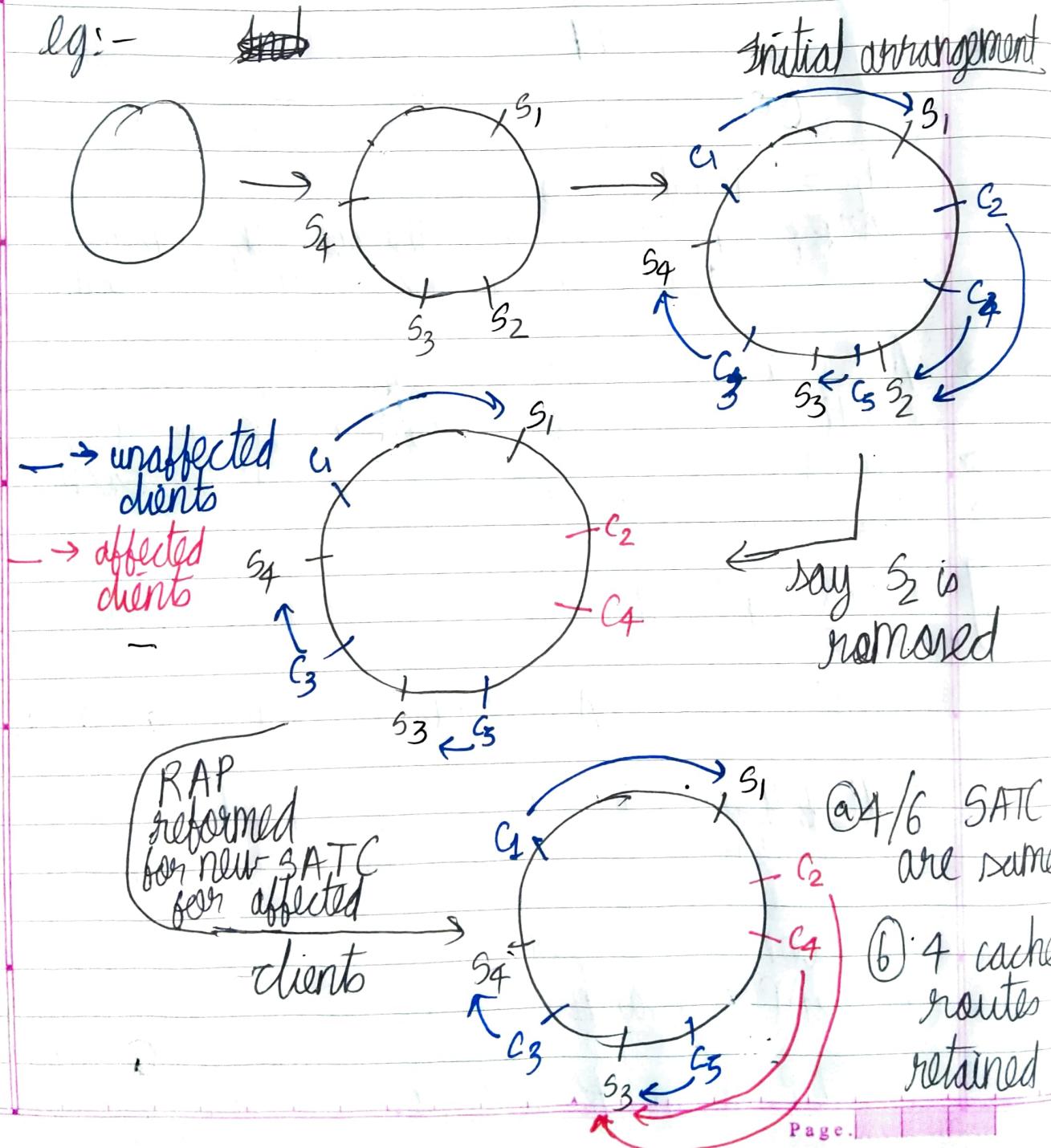
Visualisation: ① Servers are placed on a circle, preferably at equal dist. covering entire circle in equal parts.

- ② Clients are now placed on the same circle
 - ③ SAP → Matching client to its closest server if we traverse circle in $\frac{1}{2}$ dirⁿ
- SAP is performed to get SATC for each client

④ Why circle? :- Cuz it's a shape that is closed

∴ It is easy to just say that on adding or deleting a server, we know the next closest server just for the clients affected and so when we re-perform SAP to get SATC for each client we retain cache routes for all unaffected clients

e.g:- ~~Diagram~~



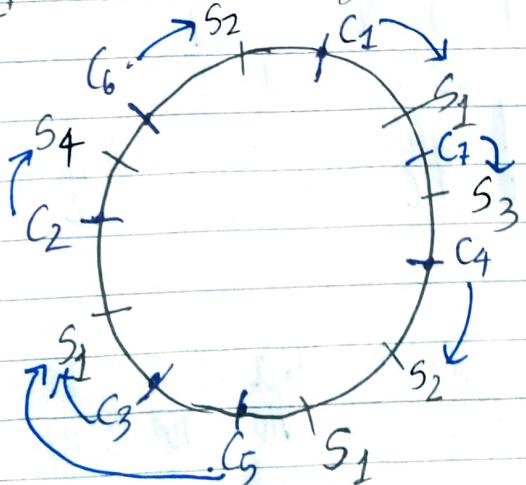
Thus, we retain at least some cache routes irresp. of addin' / removal of server.

high latency

- NOTE:
- ① Also, as u see, C_1 is very far from S_1 . This prob. can be solved by arranging servers ~~at~~ at equal dist.s, coverin entire area of wide ~~real world~~ set up servers in geographical loca's s.t. you minimize dist.s b/w ^{all} servers and clients mappin to em

② We can traverse circle in P or G dir.
~~P~~ Dirⁿ doesn't really matter as long as it's fixed.

- ③ Say if you have a more powerful server i.e. server scaled vertically, just place it on more pts. in circle to ensure more clients are directed to it ~~real world~~ place powerful server/ scale a server vertically, in geographical loca's where no. of clients is high.



$$\text{Power} \rightarrow S_1 > S_2 > S_3 = S_4 = S_5$$

S_1 serves 3 clients

S_2 serves 2 clients

S_3, S_4 serve 1 client

Rendezvous Hashing:

- i) For every client, H.F. calculates and allocates score for each server (i.e. ranks the servers) (H.R.S.)
- ii) The highest ranking server is chosen as the destination server for our client. A uniform H.F. will ensure that for each client, we're a diff. H.R.S. i.e., a diff. destination server.
- iii) Thus, each server acts as an H.R.S. for same client.
- iv) ~~Removing server~~ → If a server is removed then we go to the client associated with it and we then allot the 2nd highest ranking servers for that client as the new destination server for that client.
- v) ~~Adding a server~~ → Re-ranking carried out for clients but new server will be pushed as H.R.S. only in 1 ranking so 1 client's destin. server will change but for other clients' ~~destn.~~ servers won't change so cache routes won't change so reqs. will still keep getting to the server with the in-mem. - cache.

Note: SHA → Secure Hash Algos → "older" of cryptographic H.F.s in the industry. These days, SHA-3 is a popular choice to use in a sys.

```
const serverSet1 = [ 'server1', 'server2', 'server3', 'server4', 'server5', ]  
const serverSet2 = [ 'server1', 'server2', 'server3', 'server4', ]  
  
const username = [ 'username0', 'username1', 'username2', 'username3', ..... 'username9' ]  
  
//SIMPLE HASHING  
function hashString(string) {  
    let hash = 0;  
    if (string.length === 0) return hash;  
    for (let i = 0; i < string.length; i++) {  
        charCode = string.charCodeAt(i);  
        hash = (hash << 5) - hash + charCode  
        hash |= 0  
    }  
    return hash}  
  
//RENDEZVOUS HASHING  
function computeScore(username, server) {  
    const userNameHash = hashString(username);  
    const serverHash = hashString(server);  
    return (userNameHash * 13 + serverHash * 11) % 67}  
  
//SIMPLE HASHING  
function pickServerSimple(username, servers) {  
    const hash = utils.hashString(username);  
    return servers[hash % servers.length]}  
  
//RENDEZVOUS HASHING  
function pickServerRendezvous(username, servers) {  
    let maxServer = null;  
    let maxScore = null;
```

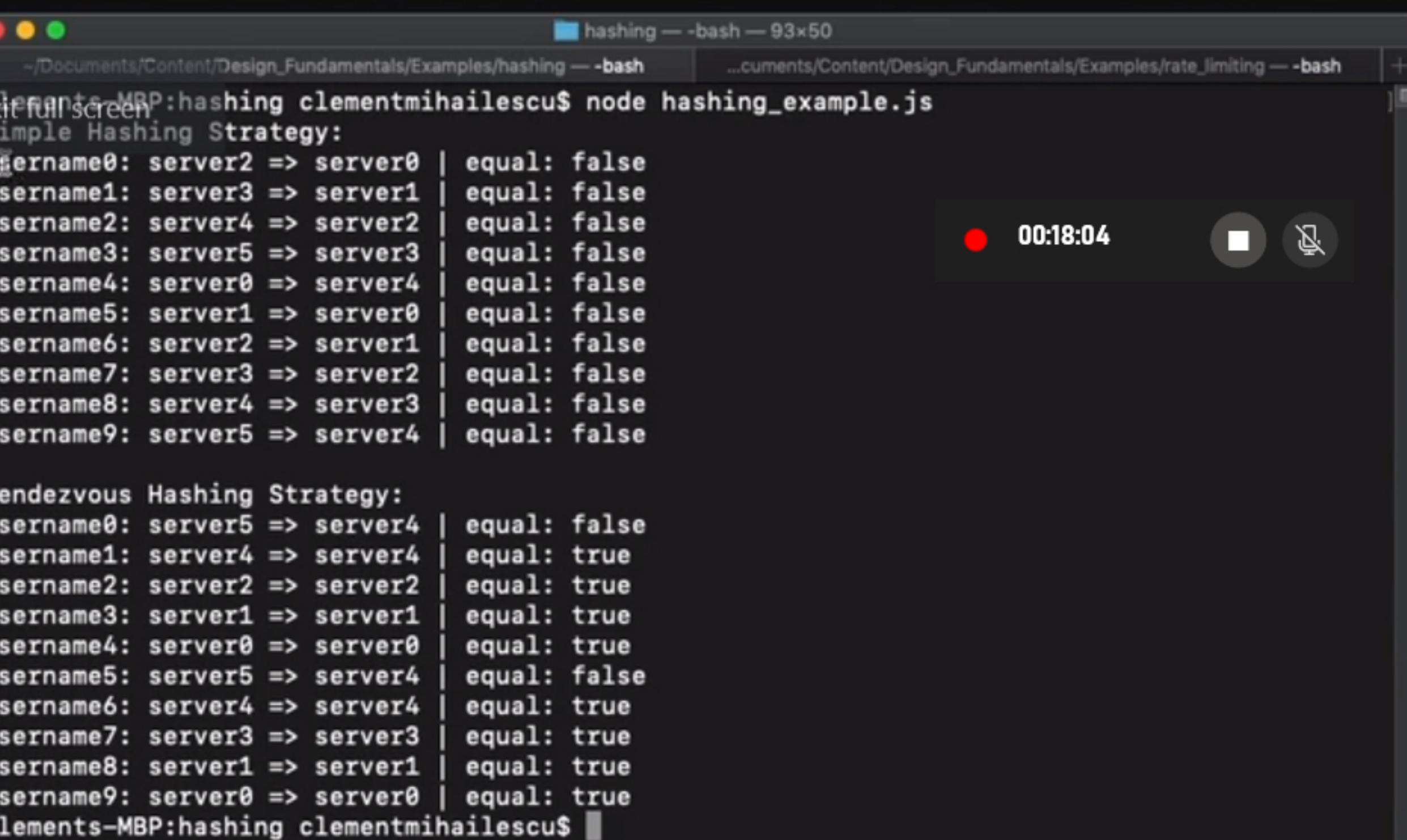
```
for (const server of servers) {  
    const score = utils.computeSCore(username, server);  
    if (maxScore === null || score > maxScore) {  
        maxScore = score;  
        maxServer = server  
    }  
}  
  
return maxServer  
}
```

//SIMPLE HASHING

```
console.log('Simple Hashing Strategy')  
for (const username of usernames) {  
    const server1 = pickServerSimple(username, serverSet1);  
    const server2 = pickServerSimple(username, serverSet2);  
    const serversAreEqual = server1 == server2  
}
```

//RENDEZVOUS HASHING

```
console.log('Rendezvous Hashing Strategy')  
for (const username of usernames) {  
    const server1 = pickServerRendezvous(username, serverSet1);  
    const server2 = pickServerRendezvous(username, serverSet2);  
    const serversAreEqual = server1 == server2  
}
```



```
hashing_example.js — hashing
JS hashing_example.js × JS hashing_utils.js
JS hashing_example.js > [e] server1
...
36
37  function pickServerRendezvous(username, servers) {
38    let maxServer = null;
39    let maxScore = null;
40    for (const server of servers) {
41      const score = utils.computeScore(username, server);
42      if (maxScore === null || score > maxScore) {
43        maxScore = score;
44        maxServer = server;
45      }
46    }
47    return maxServer;
48  }
49
50  console.log('Simple Hashing Strategy:');
51  for (const username of usernames) {
52    const server1 = pickServerSimple(username, serverSet1);
53    const server2 = pickServerSimple(username, serverSet2);
54    const serversAreEqual = server1 === server2;
55    console.log(` ${username}: ${server1} => ${server2} | equal: ${serversAreEqual}`);
56  }
57
58  console.log('\nRendezvous Hashing Strategy:');
59  for (const username of usernames) {
60    const server1 = pickServerRendezvous(username, serverSet1);
61    const server2 = pickServerRendezvous(username, serverSet2);
62    const serversAreEqual = server1 === server2;
63    console.log(` ${username}: ${server1} => ${server2} | equal: ${serversAreEqual}`);
64  }
Press Esc to exit full screen
```

Like here, you're thinking, okay well,

Relational Databases

- 1) lots of DBs out there
- 2) 2 major categories depending on rules imposed on these DBs

① Relational a.k.a SQL DB (eg:- POSTGRESQL)

- (i) A type of DB that imposes on the data stored in it, a tabular like structure (data stored as table)
- (ii) Record → instances of the entities that the resp. tables represent
- (iii) Attribute → Attributes of the entities
- eg:- table → payments, record for a shop
 row → each customer payment details such as date, amount, transaction medium
 column → the attr. → ↑ ↑ ↑

Cust. name	Date	Amt	Transac' medium
A	17/9	10	Credit Card
B	15/9	45	Cash
C	14/9	3	Netbanking
D	13/9	87	Cash

(iii) Schema → Define the shape of the table i.e. it will basically specify what attr (col) each entity (row) will have

eg:- Schema for a book library

Each row → book name

For each row, cols → borrowed on, returned on, details of payment, condition of book

② Non relational dB aka noSQL (eg! - MongoDB
Google Cloud Datastore)

(i) dBs don't impose tabular struc. on data stored in them. Sometimes they may impose some kinda other struc. but def by not table

(ii) flexible, less rigorous.

SQL → Structured Query Lang.

- i) Most relational dB (rDB) support SQL
- ii) SQL is used to perform complex queries in rDBs
- iii) Because rDBs support SQL, they ~~support~~ come with powerful querying capabilities and this is also the reason why most people choose rDBs over non-rDBs for part of their sys.
- iv) Why SQL > Python, JS

Cuz for Python, JS you gotta load data in mem to perform these kinda queries

So when we talk abt large scale distri. sys, we got TBs of data and we can't do this trivially

SQL dBs aka RDBs

1] Must use ACID transactions → transaction supports ACID props.

Atomicity Consistency Isolation Durability

i) Atomicity - If a transaction consists of multiple sub-ops then these sub-ops are gonna be considered as 1 unit so they'll either all succeed or all fail
eg:- funds transfer in banks

- sub-opn 1 → Reduce money in person 1's acc
- sub-opn 2 → Increase money in person 2's acc

ii) Consistency (aka Strong Consistency)
↓
dB is never invalid and never stale

→ i) Any transaction in dB is gonna abide by all the rules in the dB → dB is never invalid

ii) Any future transaction in dB, is gonna consider any past transaction's in dB
i.e. no 'state' state in dB where 1 transaction has exec. ed but another transaction doesn't know it has exec. ed
→ dB is never stale

iii) Isolation - Multiple transactions can occur at the same time but under the hood, they will actually have been exec. ed as if they had been done sequentially (i.e. put in a queue) one-by-one

iv) Durability - When you make a transaction in a dB, the effects of that transaction in the dB are permanent

i.e. data stored in the DB is effectively stored in disc (not mem.)

Database Index (DBI)

A DBI is a data structure that improves the speed of data retrieval ops on a DB at the cost of additional writes and storage space to maintain the index data structure.

Consider this eg for bank passbook

Month	Date	Amount credited	Amount debited
January	28 th	10000	-
Feb	7 th	5000	-
March	15 th	700	-
April	27 th	2000	-

Say we wanna figure out on which day of the year was the highest amt. credited. (Ans: Jan 28th)

i) Conventional Meth. :- Traverse thru entire look thru entire DB to find that amt
 $O(N)$ T

ii) Using DBI :- Auxiliary data structure created i.e. optimized for fast searching on a specific attr. (col.) in the table.

Say, our DBI = table that has all credited amounts stored in sorted order ($O(N \log N)$) and

Key-Value Stores

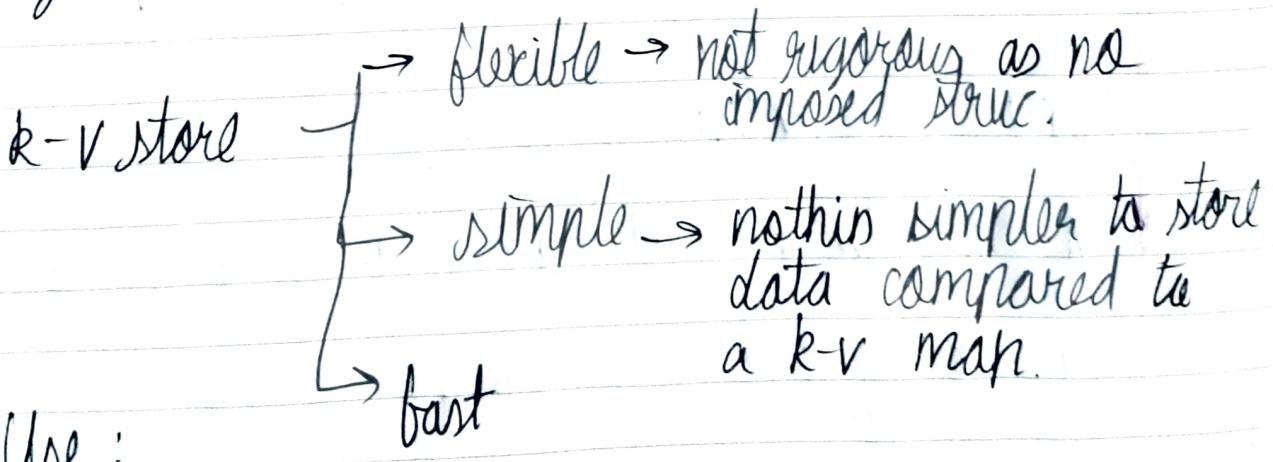
Sometimes rDBs are too cumbersome ~~in terms of~~ when weighed against the services they offer. Also, they impose a tabular struct on data viz not reqd. everytime. In those cases, we use non-rDBs for more flexibility and ease.

- Key-Value Store → i) A popular non-rDB (i.e. noSQL)
- ii) Data stored as k-v pairs
 - iii) Allows mapping of keys (strings usually strings) to ~~values~~ arbitrary values.

eg 1:

<u>Key</u>	<u>Value</u>
foo	-48.8
bleh	Apple Pie
whut	[x, y, z]

eg 2: hashTables



Use:

- i) Caching : Values → Responses to network reqs
Keys → Hash, IP address, username

Date: _____
M T W T F S S

2) Dynamic config: Basically, just storing a parameter that diff. parts of your sys. rely on in a separate place so all these parts can access it.

e.g.: - using config when you wanna switch what code you wanna run depending on what envmt. you're in → (development, prod, testing)

Now, this param. is stored in a k-v store (usually an obj. in .json or .yaml file) ~~as~~ as working with diff. values of this same param (key) is a lot easier.

NOTE: Advantage of using a k-v store → ∵ values are accessed directly thru keys, no search thru DB or any other fancy lookup reqd.
∴ latency ↓, throughput ↑ of our sys.

e.g.: - DynamoDB, Redis, Etcd, Zookeeper

NOTE: ~~store~~ 2 types of k-v stores based on where they write their data to:

i) Write data to disc → a) How to get back data from disc
b) Data persists even after k-v store server itself crashes

ii) Write data to mem → a) Faster to retrieve data from mem
b) Data lost if k-v store server crashes

Use: - Caching → As honestly, it doesn't affect the sys. much if caches crash but we do want reading from cache = fast

~~When a cache~~

When a caching in-memory k-v store goes down, we pretty much don't lose anything except some cache hits which is not that imp. a loss.

def of k-v stores:

- 1) Etcd → strong consistency
→ high availability
→ Use: implement leader elec' on a sys.
→ writes data to disc
- 2) Redis → writes data to mem. a.k.a in mem. k-v store
→ V. lit persistent storage
→ Uses: best-effort, fast caching implementation of "rate limiting"
- 3) Zookeeper → writes data to disc
→ strong consistency
→ high availability
→ Use: store imp config
implement leader elec'

A screenshot of a developer's desktop environment. On the left, a code editor window titled "server.js — key_value_stores" shows two tabs: "JS server.js" and "JS database.js". The "server.js" tab contains the following Node.js code:

```
JS server.js > app.get('/nocache/index.html') callback
1 const database = require('./database');
2 const express = require('express');
3 const redis = require('redis').createClient();
4
5 const app = express();
6
7 app.get('/nocache/index.html', (req, res) => {
8   database.get('index.html', page => {
9     res.send(page);
10  });
11 });
12
13 app.get('/withcache/index.html', (req, res) => {
14   redis.get('index.html', (err, redisRes) => {
15     if (redisRes) {
16       res.send(redisRes);
17       return;
18     }
19
20     database.get('index.html', page => {
21       redis.set('index.html', page, 'EX', 10);
22       res.send(page);
23     });
24   });
25 });
26
27 app.listen(3001, function() {
28   console.log('Listening on port 3001!');
29 });
```

On the right, a web browser window titled "Purchase | AlgoExpert" displays the URL "algoexpert.io/purchase". The page features a dark blue header with the AlgoExpert logo and tagline "Ace the Coding Interviews". Below the header, a large white banner with the text "The Ultimate Platform." and "Organized. Guided. All-encompassing. That's AlgoExpert." A modal box in the center of the page says "You've purchased all available products." At the bottom of the browser window, there is a footer with links: "Contact Us" | "FAQ" | "Reviews" | "Become An Affiliate" | "Legal Stuff" | "Privacy Policy". The status bar at the bottom of the screen shows "Ln 7, Col 47" and "Spaces: 2".

Configurations

A] Config → set of params / consts that we meet up
 Parts of our sys / app use so, we store em in
 a separate (seemingly isolated) file and access
 em everywhere

Use 1 : Storing config files for params that affect
 certain code blocks in your app in a
 certain manner depending on what
 env you're in → dev, prod, testing

FORMAT:

Usually config files are stored in .json or .yaml
 but they can be stored in any format and
 depends on kinda work people's pref., etc.
 but it pretty much doesn't differ in terms
 of overall functionality of the file

C] 2 types of config:

1) Static config : Config bundled with code so
 if you wanna change any specific param in
 config file, you gotta redeploy entire code
 (i.e. make it go thru entire testing, code
 review process etc.) cuz config is packaged
 with code

Advantage : Each time you submit code to code-base,
 code undergoes thorough review
 process. ∵ change might break the
 app, the breakage will likely be
 caught on personal tests, quality

assurance tests (tests that happens before app is deployed). Safer approach

Cons → Slower to see changes as entire app's gotta be redeployed

2) Dynamic Config: Config is completely separated from app code. ∵ Config changes are instantaneously visible in app.
 This is a bit complex as it's backed by a dB that your app is querying to see what the curr. config

Use: eg:- build a UI on top of a config. ∵ any changes in UI → super easy to implement, visualize. w/o redeployment of app code

Pros: Fast changes in app. Deploying new ~~the~~ features real fast.

Cons: Risky cuz not a lotta thorough review, testing.

Implementation: Dynamic config is used with tools built around it to make it safe

- eg:- i) builds a review sys on top of a UI built on top of ~~config~~ dynamic config
- ii) access controls → only certain ppl in org can make config changes
- iii) build deployment sys around config

eg:- app related to config" is deployed
only after specific periods of time

period of time
after which app
is deployed
automatically

eg:- twice / week

period of time after
which config" changes
are deployed

eg: twice / day

a build deployment sys.
around config"

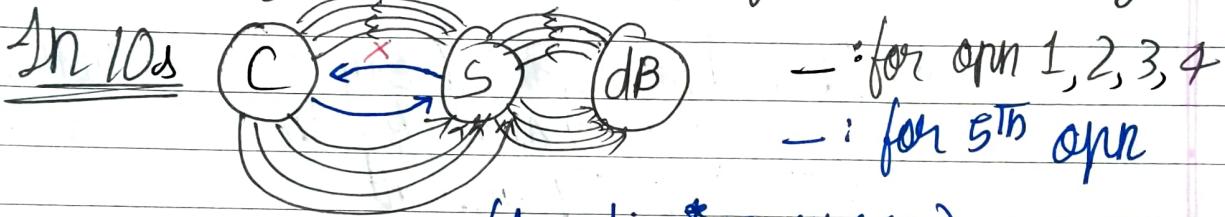
iv) Only few users (say 1%) actually see
the config" change for the 1st few hrs/days
after app deployment

RATE LIMITING

Rate limiting :

- ① Setting a threshold of sorts on certain opns. past which these opns return errs. limiting the amt. of opns. that can be performed in given time

e.g.: - Only 5 HTTP req. from client every 10 secs.



→ (1 machine * many reqs) = many reqs

- ② Denial of Service (DoS) attack ⇒ Malicious client floods the sys with too many req and effectively clogs the sys with way more traffic than it can handle and thus, sys is brought down

- ③ Rate limiting prevents clogging of servers by DoS attacks. past a certain no. of reqs, it'll throw errs. and stop taking in reqs.

e.g.: - "code exec" engine on AlgoExp. Run code opn" is rate lim. to prevent someone from spamming the opn" i.e. clogging the DB with reqs.

- ④ Rate limiting can also be dependent on other factors:

- i) Based on user: Identify user issuing req.

by lookin at req. headers or authentica's credential
 and ~~user~~ once identified we can rate limit
 no of opns. by user eg:- 3 opns/min for
 a student to access his/her grades on insti. website

- i) Based on IP address
- ii) Based on region of the world
- iii) Based on whole sys. eg:- a sys as a whole
 should never allow servers to handle $> 10^4$
 reqs per min.
- v) Based on a certain org. fallin under same
 IP address

⑤

DDoS attack \rightarrow Distributed DoS attack

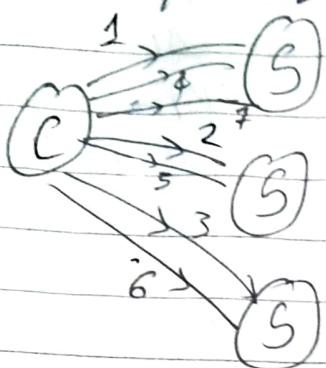
\hookrightarrow (many machines * some reqs by each machine) = (many reqs to servers)

might not be identifiable as part of same org/grp.

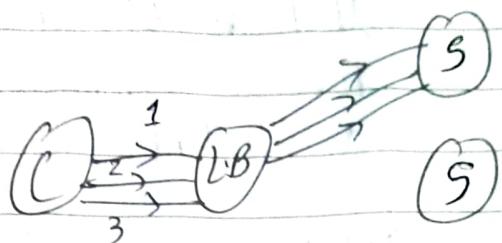
⑥

In large - scale distri sys \rightarrow we. gotta hr. v.
 powerful rate limiting ~~to prevent~~ \rightarrow i.e. we
 need to hr. really powerful load balancin
 to ensure that a client's req are always routed
 to the same server (amongst a set of servers
 cuz its a big sys) to ensure that we have ~~to~~
 some prev access records of this clients reqs. to
 the sys, stored in-memory in the ~~same~~ ^{some servers} sys
 we can compare and say whether amt. of
 reqs. this client is issuing is within a safe
 boundary or not

w/o LB in



with LB



Client can flood all servers as rate lim.
is met / server but not across servers (sys as a whole) so sys is flooded with traffic

Client always routed to same server → it has nr/r access records of this client so it'll rate lim. client if it's flooding with too many reqs.

∴ It is imp. to config load balancing keeping rate lim. in mind

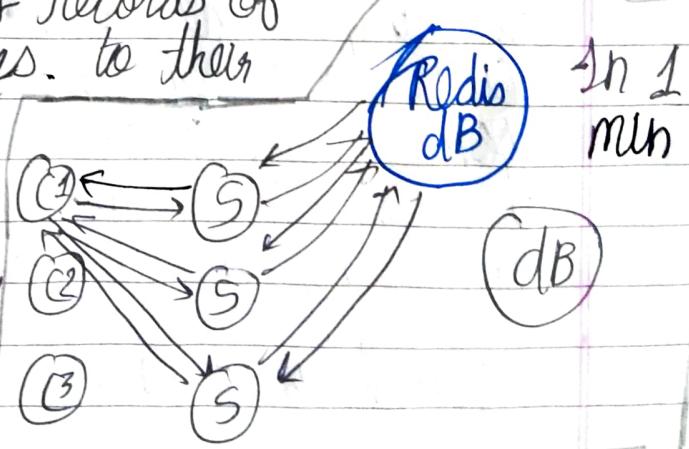
7 ∵ in most cases, large scale distri. sys's use handle rate lim. in a separated service or dB instead of in-mem. in servers. All ~~these~~ your servers speak to this ~~dB~~ and separate dB and realize if they gotta rat lim. client or not

which stores nr/r records of all clients' reqs. to their resp. servers

e.g:- Redis. → ^{eg 5 req/}
_{min allowed}

separate dB
that stores nr/r client access records and does rate lim logic and likewise,

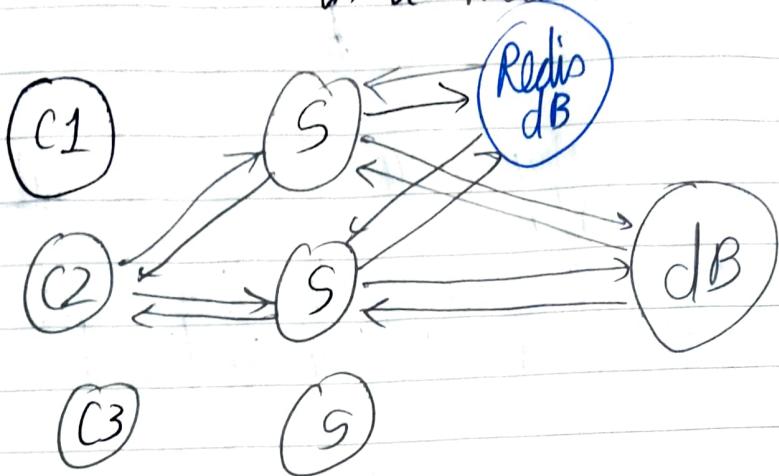
tells a server to rate lim a client or not



on its 3rd req.

Date: _____
M T W T F S S

In our ex:- C1 is blocked / rate lim from all other servers as it issued > 2 reqs. in a min.



C2 has ≤ 2 reqs per min so it gets back data from main dB both times.

⑧ tier based Rate Limiting

ex - HTTP reqs:

1 req./0.5s but only 3 reqs/10sec but only 10 reqs/min

Thus, we ~~still~~ give users bit more freedom whilst still maintaining a strong grasp over their power to keep req. in

Complex to code this out tho

server.js — rate_limiting

JS server.js X JS database.js

JS server.js > app.get('/index.html') callback > [0] previousAccessTime

```
1 const database = ...;
2 const express = require('express');
3 const app = express();
4
5 app.listen(3000, () => console.log('Listening on port 3000.'));
6
7 // Keep a hash table of the previous access time for each user.
8 const accesses = {};
9
10 app.get('/index.html', function(req, res) {
11   const {user} = req.headers;
12   if (user in accesses) {
13     const previousAccessTime = accesses[user];
14
15     // Limit to 1 request every 5 seconds.
16     if (Date.now() - previousAccessTime < 5000) {
17       res.status(429).send('Too many requests.\n');
18       return;
19     }
20   }
21
22   // Serve the page and store this access time.
23   database.get('index.html', page => {
24     accesses[user] = Date.now();
25     res.send(page + '\n');
26   });
27});
```

```
rate_limiting — bash — 93x25
...ents/Content/Design_Fundamentals/Examples/rate_limiting — node server.js | ...cuments/Content/Design_Fundamentals/Examples/rate_limiting — bash
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
<html>Hello World!</html>
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
<html>Hello World!</html>
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
Too many requests.
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
Too many requests.
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
Too many requests.
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
<html>Hello World!</html>
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
<html>Hello World!</html>
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: clement' http://localhost:3000/index.html
Too many requests.
Clements-MBP:rate_limiting clementmihailescu$ █
rate_limiting — curl -H user: antoine http://localhost:3000/index.html — 93x23
~/Documents/Content/Design_Fundamentals/Examples/rate_limiting — curl -H user: antoine http://localhost:3000/index.html
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: antoine' http://localhost:3000/index.html
<html>Hello World!</html>
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: antoine' http://localhost:3000/index.html
Too many requests.
[Clements-MBP:rate_limiting clementmihailescu$ curl -H 'user: antoine' http://localhost:3000/index.html
```

SPECIALIZED STORAGE PARADIGMS

1]

BLOB STORAGE (BLOB = Binary Large Object)

i) Blob → arbitrary piece of unstructured data

e.g. - video file, img file, large txt file, large binary (compiled code) file

ii) Blob store → storage sol'n for blobs.

iii) Info abt Blob stores

① They don't count as a dB since they only allow storage and retrieval of data based on the name of the blob.

② They're usually used for storin large size blobs for small and large scale sys.s

③ Blob-store accesses blobs usually thru a key or blob name
It specializes in storin massive amts. of unstructured data i.e. blobs.

④

Blobk-v store

Blob accessed thru key

Value accessed thru key

Optimized to store massive amts. of unstructured data → high availability and durability of data

Optimized for latency

(5) Really complex prob → accessing blobs from blobstore, so generally relyin on popular blobstores soln. (eg:- Google Cloud Storage, Amazon S3) is suggested as these are big companies that have that kinda infra. to support blobstore soln.
 Money is charged based on

- how much storage is used
- how often are blobs stored, retrieved

2] Time Series Database (TSDB)

i) TSDB - special kinda dB optimized for storing and analyzing time - indexed data

ii) Time Indexed Data - Data pts. that specifically occur at a given moment in time
 Data = events that happen at a given time say every millisecond

To perform time-series like computa's on this data, eg → computing avg

- ↓
- TSDB is used
- computing local maxima, minima
 - aggregatin all data b/cr 2 specific pts in time

iii) Use cases:

① Monitoring: ~~Monitoring~~ Looking, ^{analysis}computin a bunch of events occurin in our sys. at a given timestamp.

② IoT: Lots of devices which are constantly sending or capturing telemetry or some other data in their envts.

③ Stock prices, Cryptocurrency prices → changing all the time

(iv) Eg:- Influx DB, Prometheus

3] Graph dB

(i) Used when within the data that you store, there are a lot of relationships b/w the data (i.e. indiv data pts. in dataset) / multip levels of relationship b/w data

Eg:- Facebook accounts

$\begin{cases} \text{dataset 1} \rightarrow \text{acc 1 : Harry, age: 16, eye-color: green} \\ \text{dataset 2} \rightarrow \text{acc 2 : Ron, age: 17, eye-color: brown} \end{cases}$ } RELATIONS-
 dataset } HIP: friends

(ii) In SQL table format, querying certain pieces of data that rely on a lot of diff relationships b/w data stored in tables, becomes really complicated

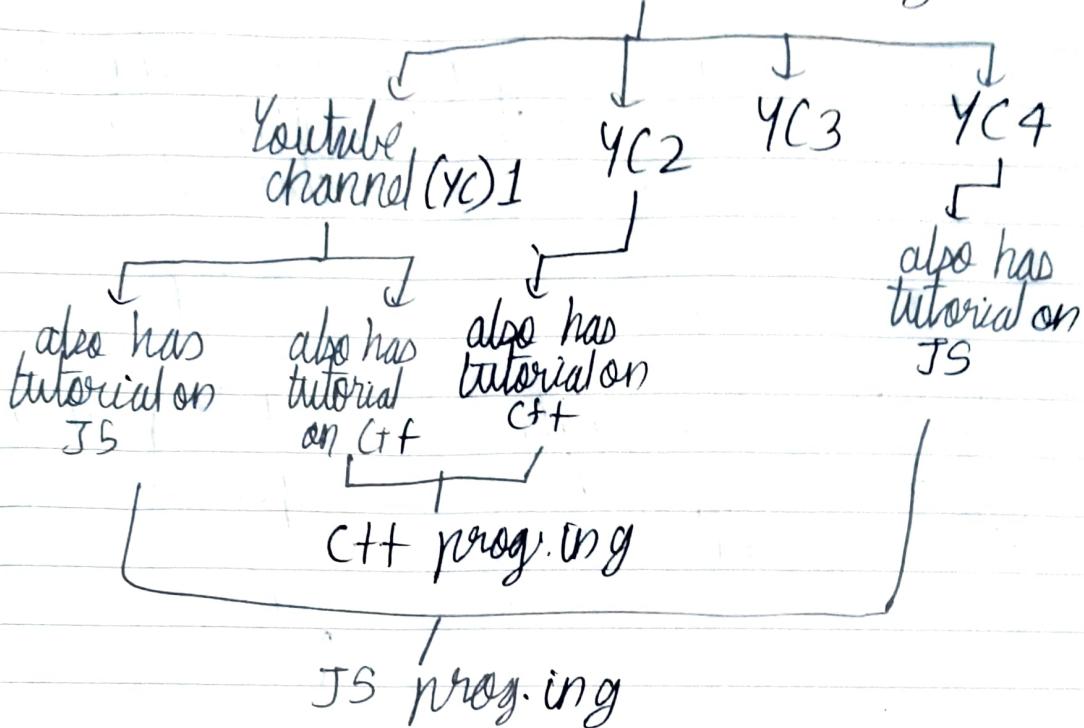
(iii) Graph dB → ① dB built on top of the graph data model

② : The concept of a relationship b/w indiv data pts. is high priority
 ③ Complex queries on deeply connected data

(iv) Use cases: done very fast

① Social Networks

② Websites with relationships b/w people and pages or content
 eg:- Content → Python Blog . org



(v) Most popular graph dB - Neo4j

(vi) Cypher - A 'graph query lang' originally developed for the Neo4j graph dB but is now used in a lot of other dBs (aka "SQL for graphs")

~~•~~ Cypher queries are lot simpler than SQL queries when dealing with deeply connected data

eg:- to find cypher query on Neo4j

MATCH (some-node: some-label)-[: SOME_RELATIONSHIP]->(some-other-node: some-label { some-prop })

4) Spatial DB

- i) Optimized for storing spatial data
- ii) Spatial Data : literally any data that deals with geometric space.
e.g. - loca's on a map, restaurants in a locality

iii) Database idx

i) Used to perform complex queries on 1 col (i.e. 1 attr say latitude) really fast.

Spatial idx

Used to perform complex queries on 2 cols (i.e 2 attrs → latitude, longitude)

Also, these queries are referred as spatial related queries

e.g.s. of spatial related queries → finding all loca's in vicinity of a specific loca
dist b/w 2 loca's

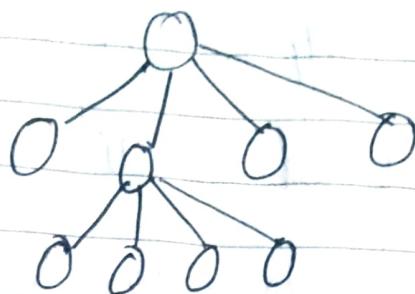
e.g.s. of a spatial idx →

- Quad tree
- R tree
- K-D tree
- M tree

iv) Quad tree

- ① A tree data structure used to index 2-D spatial data
- ② Rule : Each tree node has no. of child nodes = 4

Eg: of quad tree



#4 or nothing

- ③ Quad tree nodes → contains some form of spatial data that has a max capa.

For node value

\rightarrow < max capa.
 node is undivided i.e. leaf node condn

\rightarrow > max capa.

node is given + kids nodes
 and its data entries is split across the 4 kid nodes

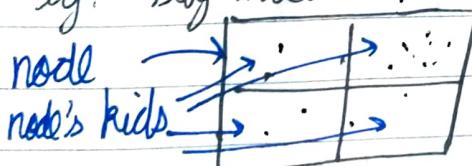
- ④ Graphical visualisation of quad tree

a) quad tree node = rect.

b) quad tree node value = dot inside rect.

c) max capa. of quad tree node value = max no. of dots a rect. can accommodate
 eg:- say max capa. = 10
 node was split before we split it into 4 kids each having some of these dots.

d) of entire quad tree = grid filled with rect.s that are recursively subdivided into sub-rect.s



eg on opp page: storing loca's in the world,
 let max-node-capa. = n

i) Root Node = world = outer most rect.

ii) If entire world has more than n -loca's,
 submost rect divided into 4 quadrants
 (each representing a region in the world)

iii)

Regions that have

$\rightarrow > n$ loca's
 further subdivided into 4 small sub-regions (sub-rects.) i.e.
 corres. quad tree node is given 4 kids nodes

$\rightarrow < n$ loca's
 undivided rects = leaf nodes in quad tree

iv)

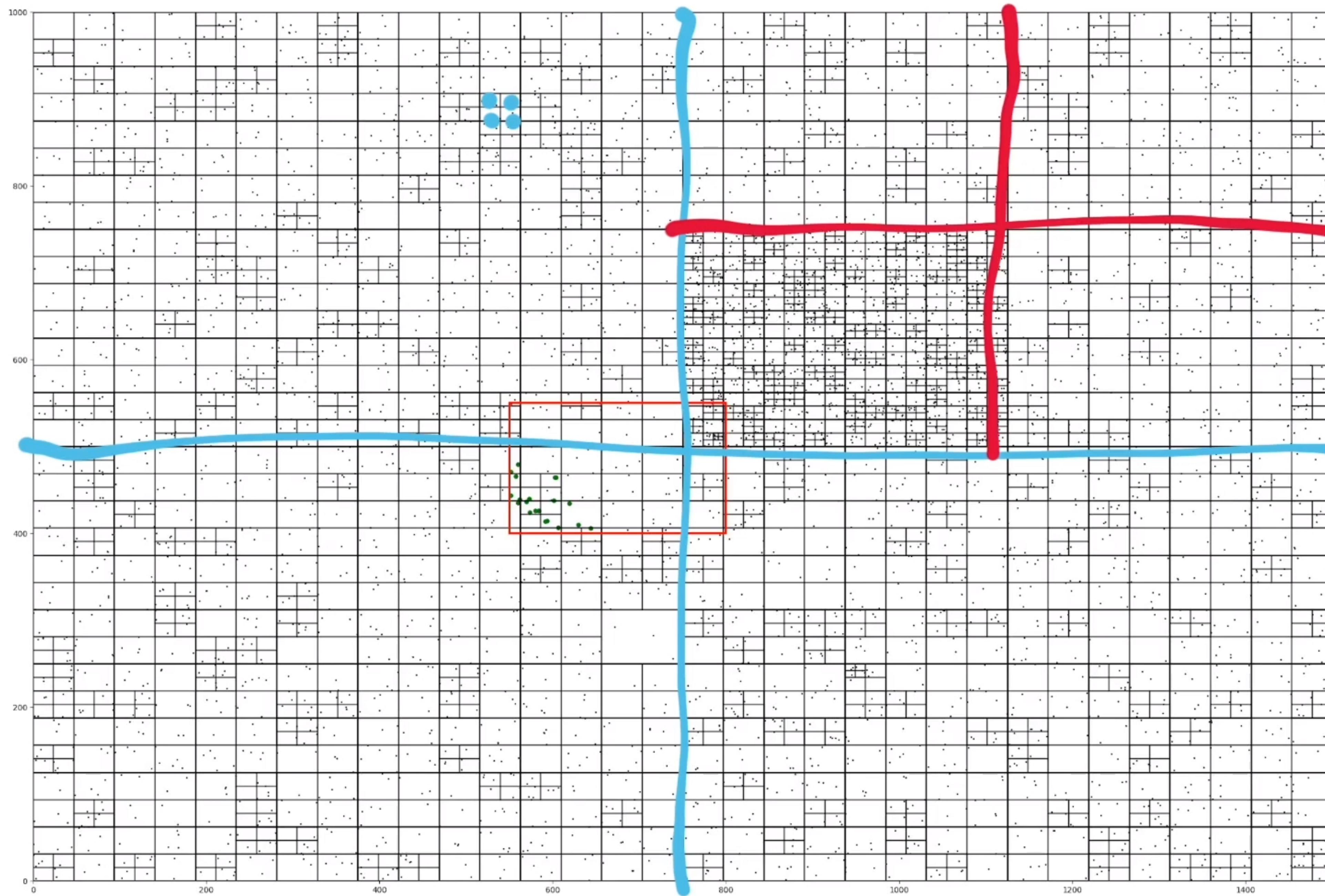
Parts of grid with

\rightarrow many subdivided rects
 represent densely populated areas
 eg: cities

\rightarrow few subdivided rects.
 represent sparsely populated areas. eg: villages

v) Finding a given loca" in quad tree: $O(\log_4 x)T$
 where $x = \text{total no. of loca}'s$

Analogy: kinda like Quadratic search versus version of Binary Search



Replica" and sharding

- 1) A dB is critical to sys.
- 2) ∵ sys's performance \propto dB's performance
- 3) $\begin{array}{c} \text{dB} \\ \text{eg 1) unavail.} \end{array} \quad \begin{array}{c} \text{sys} \\ \text{unavail.} \end{array}$
 $\begin{array}{c} \text{eg 2) low throughput} \\ \text{low throughput} \end{array}$
- 4) While des. in a sys and making it performin, we gotta ensure its dB is performant, else sys will fall apart

I]

Replica"

The act of duplicating data from 1 dB servers to others

use case (i) : To act as backup if main dB server fails

- ① ~~When~~ When main dB goes down, we can't perform any opns. (read, write, etc.) on the data ^{clients}
- ② To prev. clients from experiencing this we establish a secondary dB (backup) i.e. a replica of the main dB.

- ③ For this to happen the main dB, on receivin req.

→ performs whatever opn (read, write, etc) the req. is for

→ updates the replica s.t. replica == main dB
 i.e. in 'sync' with dB at any given pt. in time

④ we do this to ensure that replica dB can take over in case main dB fails

⑤ main dB fails → replica takes over as your main dB after some time ↓

now, they swap roles, if main dB takes over again as main dB

now, main dB is updated by replica

main dB comes back up

replica goes back to act as a standby.

⑥ write oprns to main dB take a lil longer as you gotta write data in main dB and replica to ensure replica is always 'in sync' (i.e. up-to-date) w.r.t. main dB.

Each write oprn of main dB should sync'dly be done in replica. If write oprn fails on replica due to some reason (e.g.: network part), then the write oprn. should also not be completed on main dB.

⑦ Replica should NEVER have stale data.

⑧ Redundancy of sys ↑

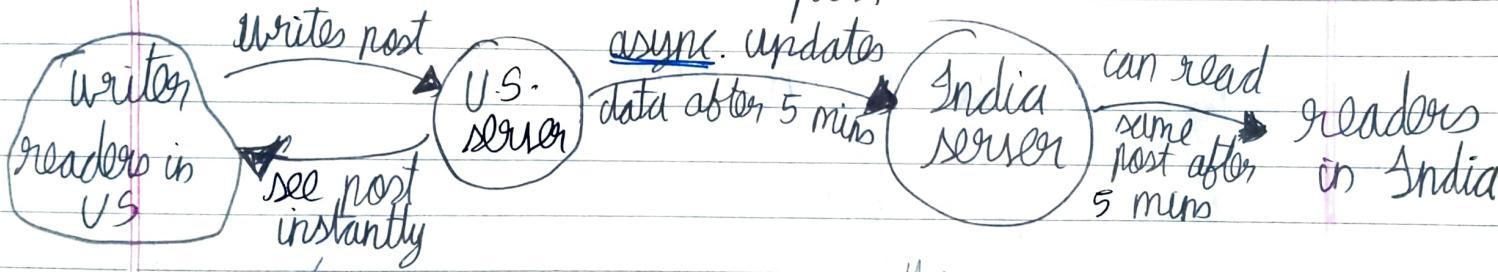
use Case (ii) : To move data closer to clients, thereby ↓ latency of accessin ~~and~~ specific data.

① Consistency : When you save data to a dB, it gets

stored but time taken \propto dist. of server from user.

- ② To bring server closer to user i.e. ↓ latency, we ~~can~~ use a replica server close to user and user stores his/her writes to it. (takes lesser time as server is closer)
- ③ We async.ly update data (say once/5min) present in all servers to ~~be~~ have consistent data everywhere
- ④ Used only in sys.s where we can afford the delay this async. update causes

e.g.: - 1 user in US writes post on LinkedIn. 1 user in India reads this post



Reason \rightarrow server is closer (not like cases where server is on other side of world so you gotta make a time consuming round trip to 1st write data (writer) and then read data (reader))

NOTE: vice versa Indian writers' posts ~~can't~~ are also read by US readers 5 mins after post is published

e.g.: - If you don't need result of a deployment of new feature in your app to be seen all over the world instantly.

II] Sharding

① It may happen that main dB is handling way more say (say, 1 billion/sec) than it should be and thus, a ~~bottleneck~~ bottleneck starts forming at main dB servers.

horizontally scaling of a server

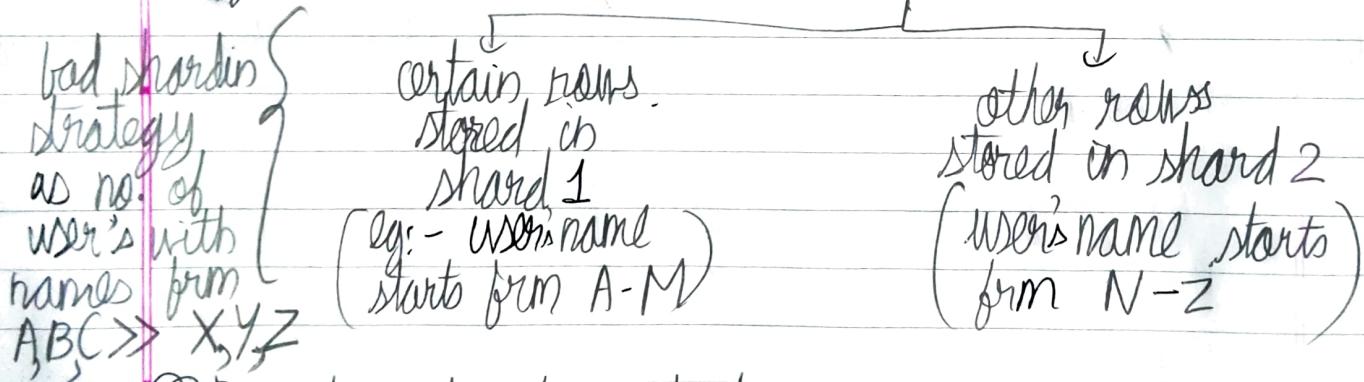
② So he ~~will split up main dB data~~ small nos' of data ~~is main~~ and store them in separate small servers (known as **shards**) and use them to ~~provide~~ provide data to diff. kinds of reqs.

③ Shards a.k.a. 'data partitions'

Sharding : Act of splitting a dB into 2 or more pieces (**shards**)

Advantage : Throughput ↑ (as no bottleneck formed)

e.g.: - say we got an rdB → table



④ Popular sharding strategies

i) Based on clients' regions

ii) Based on type of data being stored

e.g.: - data → split → shard 1 → user data

iii) Sharding based on the hash of a column (only done for structured data)

- ⑤ Hotspots :- Certain shards that get more traffic than other shards cuz of the kinda data they store
- Cause → shardin key or hashin fⁿ is sub-optimal
- Effect → worsen distri. of workload across shards

⑥ Minimizing hotspots i.e. evenly splitton up data

i) Use a hashin fⁿ (H.F.) that guarantees uniformity in determinin what data goes to which shard.

ii) Also, we need to ensure ~~the~~ H.F. doesn't change much to ensure that a particular type of data always goes to the same servers

iii) Consistent hashin is kinda useful (H.F. is const)

→ Addin a new shard → ~~the~~ consistent hashin minimizes the pieces of data we gotta migrate from existin shards to new shard

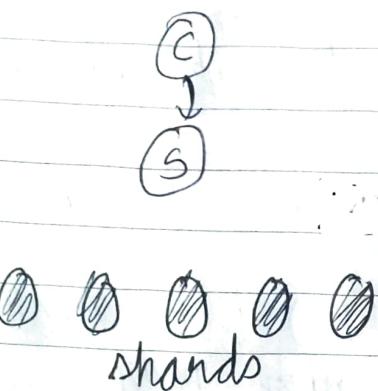
→ Shard goes down → consistent hashin can't do anything

↳ soln:- replicate each shard.

∴ ~~think~~ Give a lotta thought to comin up with a good logic to split up data while designin a sys.

(7) Storage of sharding strategy logic

M-1 → stored in server



M-2 → stored in reverse proxy b/w server & shards



Place of storage basically decides that based on clients req which shard is it → gonna req. data from OR write data to

```
//REVERSE PROXY - STORES LOGIC TO SPLIT INCOMING DATA (WRITE REQUEST) INTO SHARDS
```

```
const axios = require('axios');

const express = require('express');

const SHARD_ADDRESSES = ['http://localhost:3000', 'http://localhost:3001'];

const SHARD_COUNT = SHARD_ADDRESSES.length;

function getShardEndpoint(key) {

  const shardNumber = key.charCodeAt(0) % SHARD_COUNT;

  const shardAddress = SHARD_ADDRESSES[shardNumber];

  return `${shardAddress}/${key}`;

}

app.post('/:key', (req, res) => {

  const shardEndpoint = getShardEndpoint(req.params.key);

  console.log(`Forwarding to ${shardEndpoint}`);

  axios.post(shardEndpoint, req.body).then((innerRes) => res.send());

});

app.get('/:key', (req, res) => {

  const shardEndpoint = getShardEndpoint(req.params.key);

  console.log(`Forwarding to ${shardEndpoint}`);

  axios.get(shardEndpoint).then((innerRes) => {

    if (innerRes.data === null) { res.send('null'); return; }

    res.send(innerRes.data);

  });

});

app.listen(8000, () => { console.log('Listening on port 8000');});
```

```
//TO BASICALLY FORWARD DATA TO REVERSE PROXY

const express = require('express');
const fs = require('fs');
const app = express();
const PORT = process.env.PORT;
const DATA_DIR = process.env.DATA_DIR;

app.use(express.json());

app.post('/:key', (req, res) => {
  const { key } = req.params;
  console.log(`Storing data at key ${key}`);
  fs.writeFileSync(destinationFile, req.body.data);
  res.send();
});

app.get('/:key', (req, res) => {
  const { key } = req.params;
  console.log(`Storing data at key${key}`);
  const destinationFile = `${DATA_DIR}/${key}`;
  try { const data = fs.readFileSync(destinationFile); res.send(data); }
  catch (e) { res.send('null'); }
});

app.listen(PORT, () => { console.log(`Listening on port ${PORT}`);});
```

```
replication_and_sharding — node aedb.js — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — node aedb.js
Clements-MBP:replication_and_sharding clementmihai@escu$ DATA_DIR=aedb_data_0 PORT=3000 node aedb.js
Listening on port 3000!
```

```
replication_and_sharding — node aedb.js — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — node aedb.js
Clements-MBP:replication_and_sharding clementmihai@escu$ DATA_DIR=aedb_data_1 PORT=3001 node aedb.js
Listening on port 3001!
Storing data at key a.
```

```
replication_and_sharding — node aedb_proxy.js — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — node aedb_proxy.js
Clements-MBP:replication_and_sharding clementmihai@escu$ node aedb_proxy.js
Listening on port 8000!
Forwarding to: http://localhost:3001/a
```

```
replication_and_sharding — -bash — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — -bash
Clements-MBP:replication_and_sharding clementmihai@escu$ curl --header 'Content-Type: application/json' --data '{"data": "This is some data."}' localhost:8000/a
Clements-MBP:replication_and_sharding clementmihai@escu$ curl -w "\n" localhost:8000/a
```

```
replication_and_sharding — node aedb.js — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — node aedb.js
Clements-MBP:replication_and_sharding clementmihai$ DATA_DIR=aedb_data_0 PORT=3000 node aedb.js
Listening on port 3000!
Storing data at key b.
Retrieving data from key b.
Retrieving data from key b.
Storing data at key bar.
Storing data at key baz.
Storing data at key foobar.
Storing data at key foobaz.

replication_and_sharding — node aedb.js — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — node aedb.js
Clements-MBP:replication_and_sharding clementmihai$ DATA_DIR=aedb_data_1 PORT=3001 node aedb.js
Listening on port 3001!
Storing data at key a.
Retrieving data from key a.
Retrieving data from key a.

replication_and_sharding — node aedb_proxy.js — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — node aedb_proxy.js
Clements-MBP:replication_and_sharding clementmihai$ node aedb_proxy.js
Listening on port 8000!
Forwarding to: http://localhost:3001/a
Forwarding to: http://localhost:3001/a
Forwarding to: http://localhost:3000/b
Forwarding to: http://localhost:3000/b
Forwarding to: http://localhost:3001/a
Forwarding to: http://localhost:3000/b
Forwarding to: http://localhost:3000/bar
Forwarding to: http://localhost:3000/baz
Forwarding to: http://localhost:3000/foobar
Forwarding to: http://localhost:3000/foobaz

replication_and_sharding — -bash — 93x24
~/Documents/Content/Design_Fundamentals/Examples/replication_and_sharding — -bash
Clements-MBP:replication_and_sharding clementmihai$ curl --header 'Content-Type: application/json' --data '{"data": "This is some data."}' localhost:8000/a
Clements-MBP:replication_and_sharding clementmihai$ curl -w "\n" localhost:8000/a
This is some data.
Clements-MBP:replication_and_sharding clementmihai$ curl --header 'Content-Type: application/json' --data '{"data": "This is some data."}' localhost:8000/b
Clements-MBP:replication_and_sharding clementmihai$ curl -w "\n" localhost:8000/b
This is some data.
Clements-MBP:replication_and_sharding clementmihai$ curl -w "\n" localhost:8000/a
This is some data.
Clements-MBP:replication_and_sharding clementmihai$ curl -w "\n" localhost:8000/b
This is some data.
Clements-MBP:replication_and_sharding clementmihai$ curl --header 'Content-Type: application/json' --data '{"data": "This is some data."}' localhost:8000/bar
Clements-MBP:replication_and_sharding clementmihai$ curl --header 'Content-Type: application/json' --data '{"data": "This is some data."}' localhost:8000/baz
Clements-MBP:replication_and_sharding clementmihai$ curl --header 'Content-Type: application/json' --data '{"data": "This is some data."}' localhost:8000/foobar
Clements-MBP:replication_and_sharding clementmihai$ curl --header 'Content-Type: application/json' --data '{"data": "This is some data."}' localhost:8000/foobaz
Clements-MBP:replication_and_sharding clementmihai$ 
```

Leader Elecⁿ

Gonna understand entire concept thru ~~one~~ one of its actual ~~use~~ use cases

- ① Case → Managing subscrip's and payments in some sort of service such as Amazon Prime, Netflix (subscripⁿ on recurring basis)

- ① dB → stores data pertaining to subscripⁿ and users

- ↳ is user curr. ly subscribed?
- ↳ date for subscripⁿ renewal
- ↳ price

- ② 3rd party service → service that conducts the payment transaction (debit from user → credit to you) eg. Paypal, Stripe
 Needs to comm. with dB to know

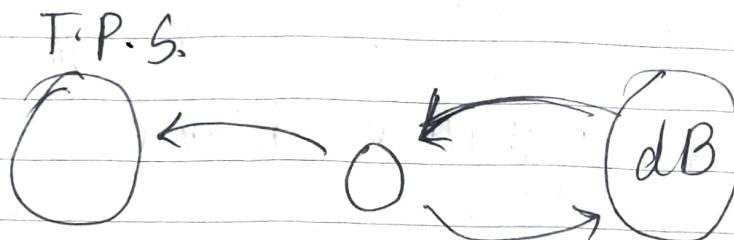
- ↳ when ~~st~~ a user shud be charged again
- ↳ what's the charge?

~~No direct connect b/w dB and 3rd party service~~

- ↳ diff. to implement

- ↳ dB is pretty sensitive part of sys. Contains imp info. and its not wise to divulge all of it by allowing a T.P.S. to connect directly

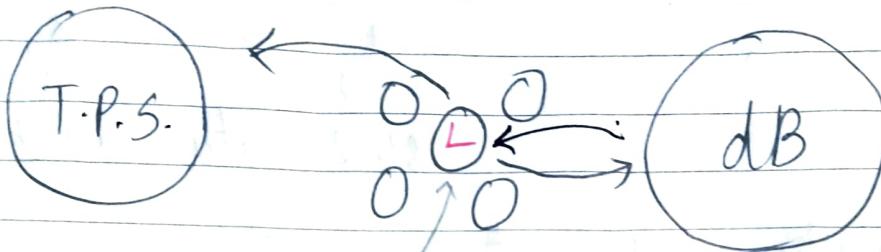
- ③ Service in middle → talks to dB periodically
 → figures out when a user's subscription is gonna renew
 → how much to charge the user
 → transmits only relevant info. to T.P.S.



Optimizing this setup.

- 1) If the server in the middle fails, entire payment sys. collapses.
 We introduce 'redundancy' (passive) in our sys. by horiz. scaling the service in the middle i.e. have 5 instead of 1 servers in the middle to do the transacⁿ logic
- 2) All 5 servers are doing the same thing i.e. asking dB for transacⁿ info and ~~req.~~ req. in T.P.S. to charge the user
 We don't wanna duplicate the req. to T.P.S. to charge the user i.e. make that req. just once
- 3) Leader elecⁿ → If you have a group of machines (is ~~over~~ case-servers) that are in charge of doing the same, instead of having them all do that thing, we elect 1 of the machines as leader

Network part^n → some network failure that makes some machines no longer be able to comm. with other machines
 and only that machine is gonna be responsib. for doing the opns. that all of the machines are ~~not~~ built to do.



say this is the leader

Other servers (followers) → just on standby to take over when leader fails i.e. if a new leader is elected from amongs these servers and it takes over.

4) Complexity of leader elec^n: Leader elec^n's logic is complex cuz:

- multiple machines are made to elect a single leader (**consensus** - agree upon something together)
- they all gotta be aware of who the leader is at any given time
- they all gotta be capable of re-electing a new leader in case curr. leader

Difficulty of logic mainly lies in

- multiple machines that are distributed are share state. Diff. cuz we dunno what might happen in a network (e.g.: network part^n)
- make multiple machines gain consensus

Here, gainin consensus → agreein who's the curr. leader

- 5) Consensus Algo. → complex, math-heavy algo that allows multiple nodes (servers) in a cluster (grp) to reach consensus i.e. agree upon some data-value.
- Eg:- Paxos, Raft

- 6) Usually, ppl just use some pre-existing 3rd party industry tool that offers them this logic whilst ~~itself~~ itself runnin a consensus algo under the hood

Eg :- Zookeeper, Etcd → help implem. leader elec'n v. easily.

7) Leader elec'n using Etcd :

- ① Etcd → k-v store
 → highly available
 → strongly consistent

Consistency → if you got ~~and~~ machine(s), readin & writin to the same k-v pair in the k-v store, you're always gonna get the same, correct value irresp of when or from which machine this k-v pair is accessed.

- ② Etcd achieves high availability, strong consistency by implementin ~~over~~ consensus algo - Raft

- ③ There are gonna be multiple machines that can read, write to the main k-v store

that Etcd supports. These machines need:

- high availability: To know which other machines are available ~~are~~ after a leader dies
- strong consistency: single source of truth for all k-v pairs in store

④ In Etcd each k-v pair can be visualised be of form:

<u>key</u>	<u>value</u>
status of machine (i.e. is it a leader) (or: a follower)	name / IP address (basically a VID (for a machine in grp i.e. node is cluster)
∴ we'll have 1 one k-v pair as follows	<u>"leader": VID for this machine</u>

→ k-v pair represents our "leader"

⑤ All machines (in our case - servers) comm. with k-v store any given pt. in time we gotta by a leader present as k-v pair in the store else we gotta elect a new one.

```

import etcd3
import time
from threading import Event

# The current leader is going to be the value with this key
LEADER_KEY = "/algoexpert/leader"

# Entry point of the program

def main(server_name):
    # Create a new client to etcd
    client = etcd3.client(host="localhost", port=2379)

    while True:
        is_leader, lease = leader_election(client, server_name)

        if is_leader:
            print("I am the leader.")
            on_leadership_gained(lease)
        else:
            print("I am a follower.")
            wait_for_next_election(client)

    # This election mechanism consists of all clients trying to put their name into a single key, but in a way that
    # only works if the key does not exist(or has expired before)

def leader_election(client, server_name):
    print("New leader election happening.")
    # Create a lease before creating a key. This way, if this client ever lets the lease expire, the keys associated
    # with that lease will all expire as well.
    Here, if the client fails to renew lease for 5 seconds (network partition or machine goes down), then the
    leader election key will expire.
    lease = client.lease(5)

    # Try to create the key with your name as the value. If it fails, then another server got there first
    is_leader = try_insert(client, LEADER_KEY, server_name, lease)
    return is_leader, lease

def on_leadership_gained(lease):
    while True:
        # As long as this process is alive and we're the leader, we try to renew the lease. We don't give up the
        # leadership unless the process/machine crashes or some exception is raised
        try:
            print("Refreshing lease; still the leader.")
            lease.refresh()
        except Exception:
            # This is where the business logic would go (eg: ask 3rd part service to charge user based on Db INFO)
            do_work()
        except KeyboardInterrupt:
            lease.revoke()

```

```

print("\n Revoking lease; no longer the leader")
# Here we're killing the process. Revoke the lease and exit
lease.revoke()
sys.exit(1)

def wait_for_next_election(client):
    election_event = Event()

    def watch_callback(resp):
        for event in resp.events:
            # For each event in the watch event, if the event is a deletion it means that the key expired / got deleted,
            # which means the leadership is up for grabs.
            if isinstance(event, etcd3.events.DeleteEvent):
                print("LEADERSHIP CHANGE REQUIRED")
                election_event.set()

    watch_id = client.add_watch_callback(LEADER_KEY, watch_callback)

    # While we haven't seen the that the leadership needs change, just sleep
    try:
        while not election_event.is_set():
            time.sleep(1)
    except KeyboardInterrupt:
        client.cancel_watch(watch_id)
        sys.exit(1)

    # Cancel the watch; we see that the election should happen again.
    client.cancel_watch(watch_id)

    # Try to insert a key into etcd with a value and a lease. If the lease expires that key will get automatically
    # deleted behind mthe scenes. If that key was already present this will raise an exception.

def try_insert(client, key, value, lease):
    insert_succeeded, _ = client.transaction(
        failure=[],
        success=[client.transaction.put(key, value, lease)],
        compare=[client.transaction.version(key) == 0],
    )
    return insert_succeeded

def do_work():
    time.sleep(1)

if __name__ == "main":
    server.name = sys.argv[1]
    main(server_name)

```

```
leader_election — Python leader_election.py server2 — 93x24
~/Documents/Content/Design_Fundamentals/Examples/leader_election — Python leader_election.py server2
Clement's-MBP:leader_election clementmihai$ python3 leader_election.py server2
New leader election happening.
I am a follower.
^CClement's-MBP:leader_election clementmihai$ python3 leader_election.py server2
New leader election happening.
I am a follower.
LEADERSHIP CHANGE REQUIRED
New leader election happening.
I am a follower.
LEADERSHIP CHANGE REQUIRED
New leader election happening.
I am a follower.
```

```
leader_election — Python leader_election.py server2 — 93x24
~/Documents/Content/Design_Fundamentals/Examples/leader_election — Python leader_election.py server2
Clements-MBP:leader_election clementmihai$ python3 leader_election.py server2
New leader election happening.
I am a follower.
^CClements-MBP:leader_election clementmihai$ python3 leader_election.py server2
New leader election happening.
I am a follower.
LEADERSHIP CHANGE REQUIRED
New leader election happening.
I am a follower.
```

```
[Clement's-MBP:leader_election clementmihai]$ python3 leader_election.py server3
~/Documents/Content/Design_Fundamentals/Examples/leader_election — Python leader_election.py server3
[Clement's-MBP:leader_election clementmihai]$ python3 leader_election.py server3
New leader election happening.
I am a follower.
LEADERSHIP CHANGE REQUIRED
New leader election happening.
I am a follower.
```

```
leader_election — Python leader_election.py server2 — 93x24
~/Documents/Content/Design_Fundamentals/Examples/leader_election — Python leader_election.py server2
Clement's-MBP:leader_election clementmihai$ python3 leader_election.py server2
New leader election happening.
I am a follower.
^CClement's-MBP:leader_election clementmihai$ python3 leader_election.py server2
New leader election happening.
I am a follower.
LEADERSHIP CHANGE REQUIRED
New leader election happening.
I am a follower.
```

```
[Clement's-MBP:leader_election clementmihai]$ python3 leader_election.py server3
New leader election happening.
I am a follower.
LEADERSHIP CHANGE REQUIRED
New leader election happening.
I am a follower.
□
```

```
[Clement's-MBP:leader_election clementmihai]$ python3 leader_election.py server4
New leader election happening.
I am a follower.
LEADERSHIP CHANGE REQUIRED
New leader election happening.
I am the leader.
Refreshing lease; still the leader.
Refreshing lease; still the leader.
Refreshing lease; still the leader.
```

each of these amts pts. to the relevant record
(row) in the main dB table
Thus, we get over largest amt by simply
goin' to the end of the dB table

Simple way of lookin at dB

- auxiliary data struc of main dB that speeds up gettin data from main dB (read opns. faster)
- auxiliary data struc so ~~is~~ extra space reqd. for storage
- whenever you write data to main dB, ~~data~~ you also gotta write to dB (write opns. slower)
- In practice, we create an idx on 1 or multip cols. in our main dB

NOTE : Eventual Consistency

- A consistency model which is unlike Strong Consistency
- In this model, reads might return a stale ~~now~~ of the sys.
- An eventually consistent ~~db~~ datastore will give guarantees that the state of the data will eventually reflect writes within a time period (eg:- 10 mins, 30 days)
- Eg:- Google Cloud Datastore