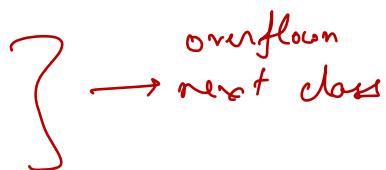


Agenda

1. What & Why of HLD ✓
2. Case Study - [del.icio.us](#) ✓
3. Scaling challenges ✓
4. Stateful vs Stateless servers
5. Consistent Hashing
6. HLD curriculum overview

support@scaler.com

+91 7351769231
(whatsapp)



Traditionally write & execute code on 1 machine (your machine)

↳ handle 100s of requests / sec

How do you scale to million / billion users

↳ how to architect the solution

for senior roles

HLD is a must!!

Google's Staff Soft Eng position.

↳ file → a bunch of strings (search queries)

↳ sort this file in asc.

There is 50 Peta Bytes of data.

10^3 - KB
 10^6 - MB
 10^9 - GB
 10^{12} - TB
 10^{15} - PB

50 Million GB
of data.

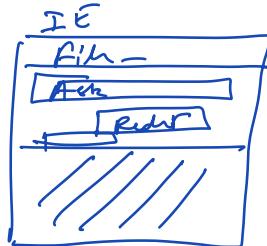


Del.icio.us Case Study

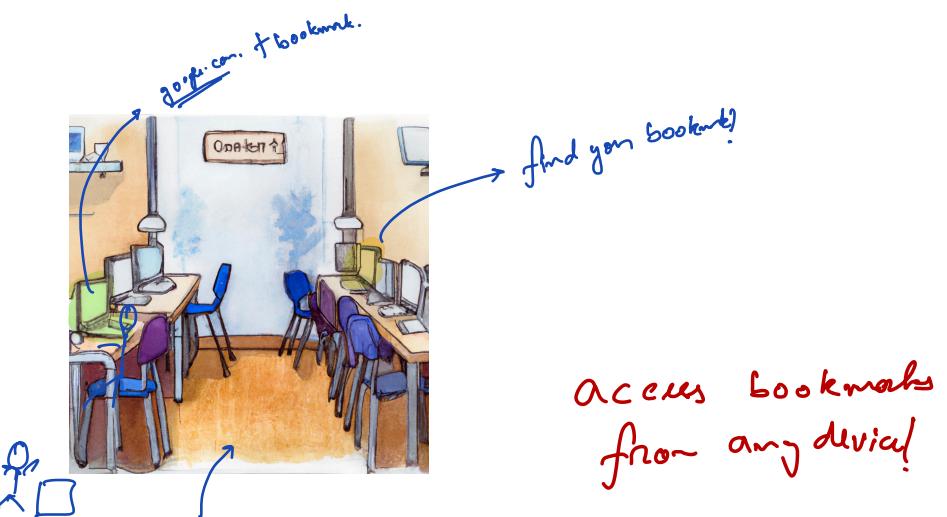
Bookmarking Service.

? Launched in

Joshua, 2003

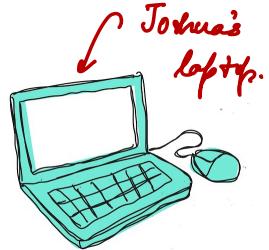


Motivation



Features → Minimal Viable Product (MVP)

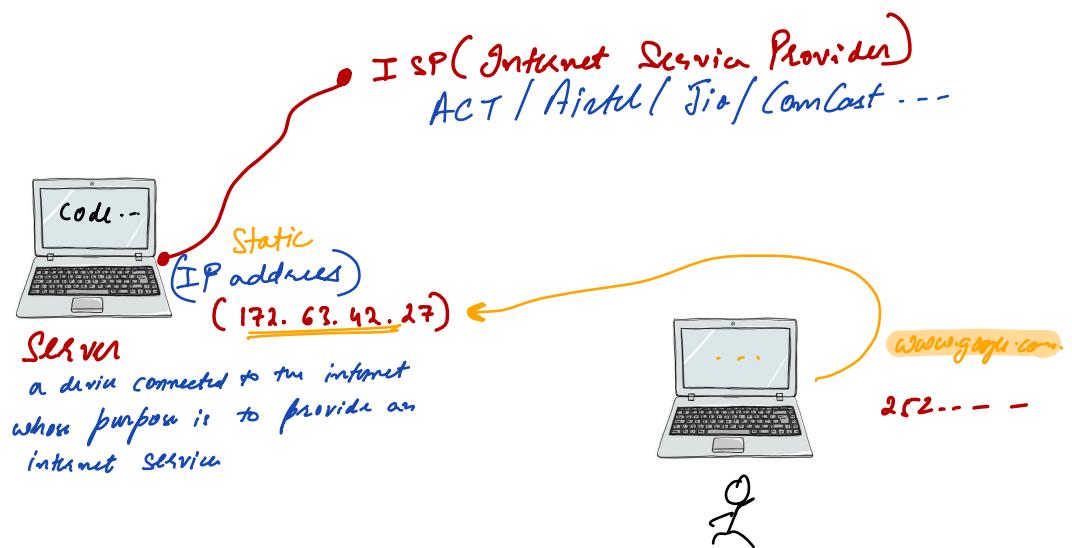
1. User registration & Login (Auth)
2. add-bookmark (url, user-id)
3. get-all-my-bookmarks (user-id)



to be continued..

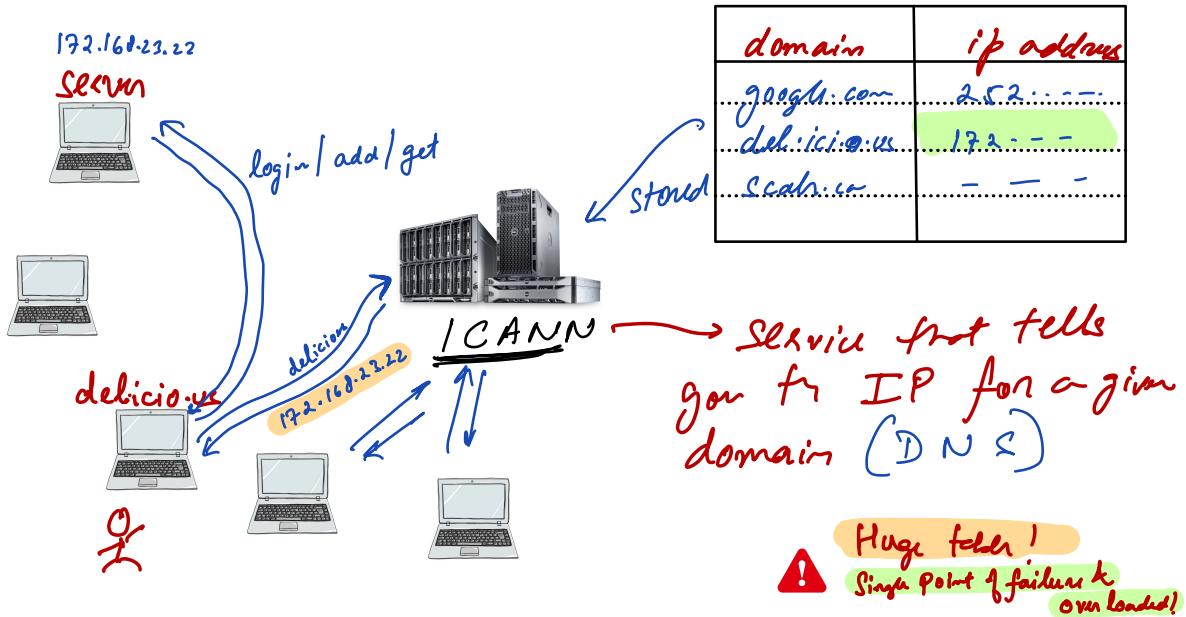


Fantastic websites & where to find them



? Quiz - how does the browser know?

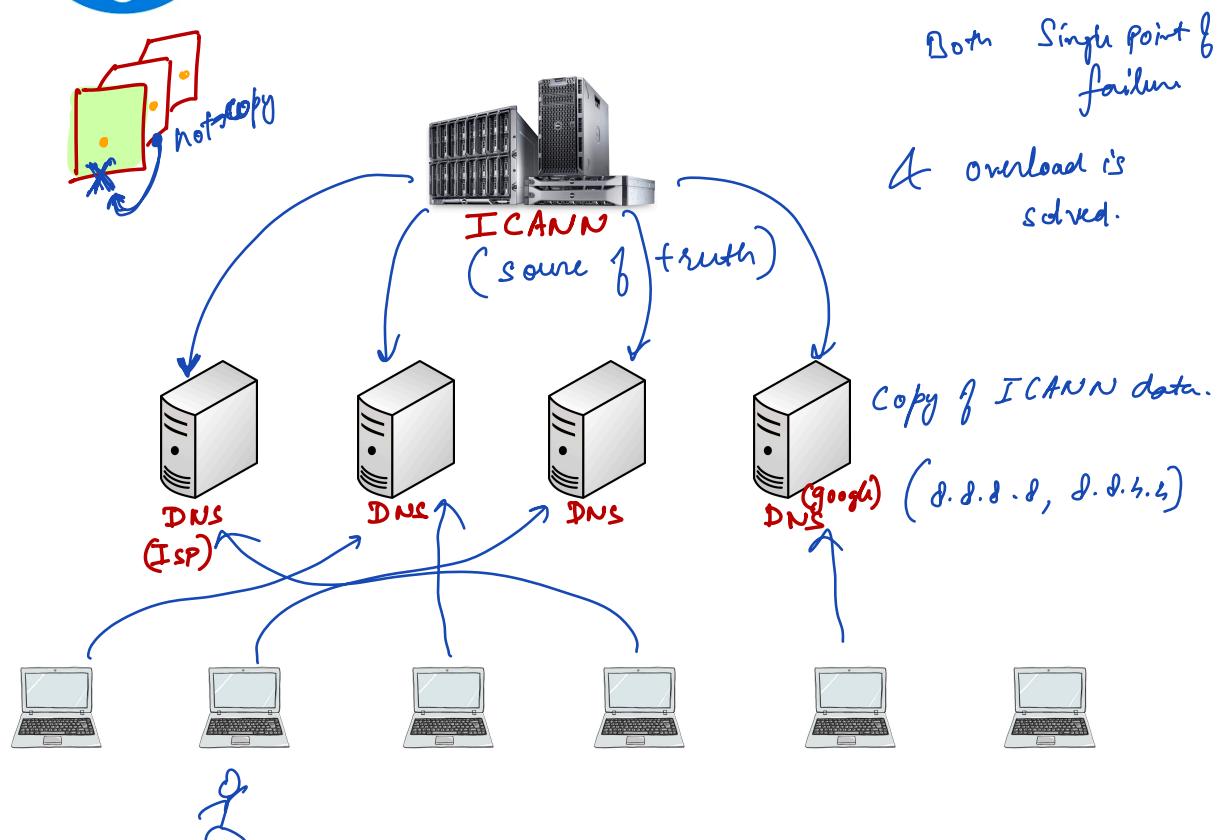
↳ somewhere there is a mapping





Domain Name Service

<https://www.cloudflare.com/en-gb/learning/dns/what-is-dns/>

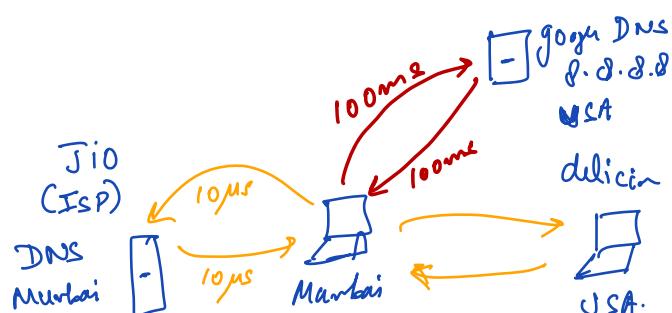


Who maintains the DNS?

domain registrar

godaddy
namecheap
google domains

Demos





Del.icio.us

continued..

loop?

System Configuration



RAM 128 MB RAM

Disk 50GB

CPU Intel Pentium Dual Core
3.2 GHz.

Postgres



del.icio.us DB



App/OS/take



Issues

① goes to sleep? System is shutdown
website goes down!!

② game / downloading website runs slow

③ Power outage website goes down



(disk)

Quiz - space

$$\begin{aligned} 2 \text{ bytes} &\rightarrow \# 2^{16} = 65536 \text{ values} \\ 4 \text{ bytes} &\rightarrow \# 2^{32} \approx 4 \text{ Billion} \\ 8 \text{ bytes} &\rightarrow \# 2^{64} \approx 16 \text{ BB} \end{aligned}$$

user-bookmarks	
4 bytes	1000 chars
user-id	url
2	google.com
21	dbdb.io
2	Scaln.com
101	netflix.com

index

1 KB per bookmark

1 Million bookmarks / day

amount of data / day

$$= 1 \text{ KB} + 1 \text{ M}$$

$$= 1 \text{ GB of data}$$

1 row $\approx 1 \text{ KB}$
(1000 bytes)

? Quiz-time

1 GB data/day

40 GB HDD available

$$\text{time to run out} = \frac{40 \text{ GB}}{1 \text{ GB/day}} = 40 \text{ days}$$



Solutions

1.

Upgrade Machine



128 Mb
40 GB
duo-con

30,000 Rs

Personal
Computing



512 Mb
100 GB
duo-con

2 lakh Rs

≈ 100 days (3 months)

dedicated
IBM
hardware



2 GB
500 GB
duo-con

10 Lakh Rs

≈ 500 days (< 2 years)



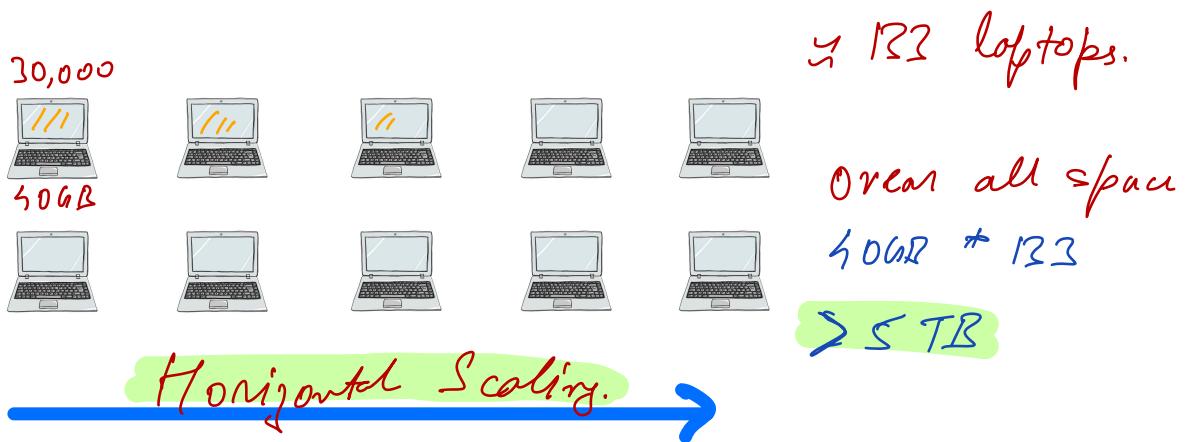
5 GB
1 TB
20 duo-con

2 crore Rs

a few years.

↑
SCALING
VR T1 CAV

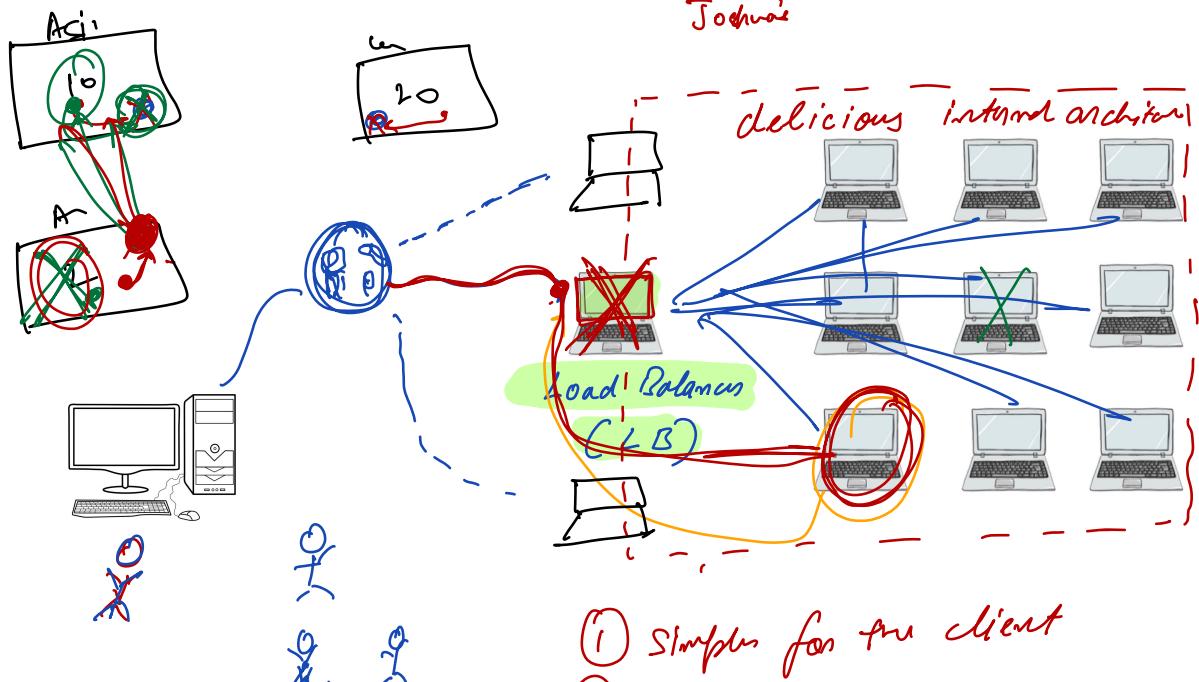
2. Buy a lot of cheap laptops.



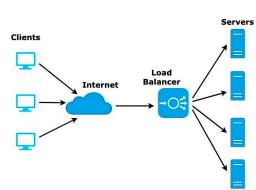
- | Scale-up
Vertical | scale-out
Horizontal |
|---|---|
| <ul style="list-style-type: none">- simple- hard limit to how much we can scale | <ul style="list-style-type: none">- difficult to manage- unlimited |
| <p>① technology is not there</p> <p>② Price will grow exponentially.</p> | <p>Price grows sub-linearly
∴ economy of scale.</p> |
| <ul style="list-style-type: none">- first choice- AWS (0-12+6-112x log₂ 109 \$ per hour)
1 million user | <ul style="list-style-type: none">- last choice↳ 448 GPU
12 TB of RAM
100 Gbit/s |

10:55 → 11:00

? Quiz - which ip?



- ① simpler for the client
- ② better security
- ③ possible to change & scale internal architecture freely.



Load Balancing

<https://docs.google.com/document/d/1DxQzLpu1XPemRWsewNWtKL6E4uwKHQhBp7GX6Sg7qI/edit>

① Load Balancer!
↳ reroutes requests to the appropriate server.
↳ distributes load evenly across servers.

- ② isn't the LB a single point of failure?? Not really
↳ easily switchable.
- ③ isn't the LB overloaded?
→ ① simple computation.
→ ② multiple load balancers per sec

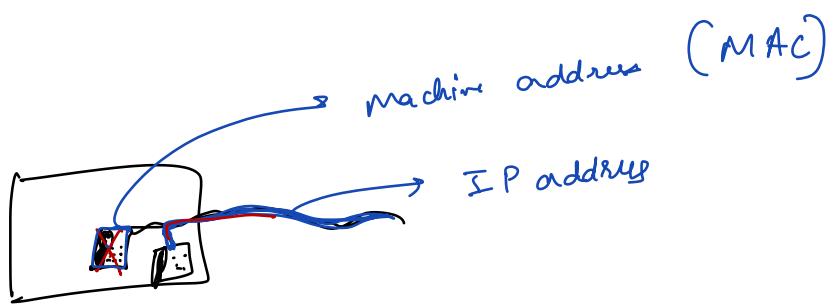
? How does the Load Balancer track which servers are up & running? & ip mappings in DNS.



Heartbeat (push-based) → servers send a periodic signal to LB → "I am alive!!"



Health check (pull-based) → LB sends requests to each server periodically to check if it is alive.



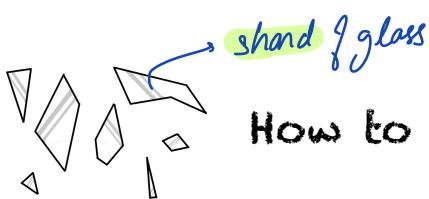


Further Challenges

? What if the Load Balancer goes down? ✓

? Which machine should we send the request to?
Routing Algorithm .

<https://docs.nginx.com/nginx/admin-guide/load-balancer/http-load-balancer/>

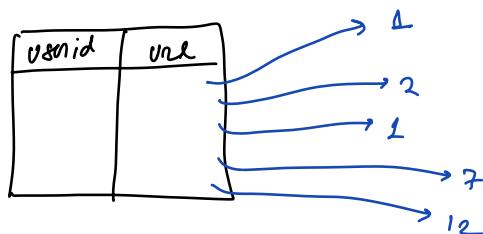


How to store the data?

? Can we store all data on 1 machine? X split the data across machines

Sharding

? Split randomly?



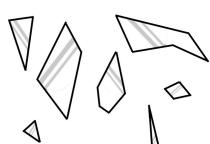
how will we retrieve it?



① store the data in a sharded manner

② also have to read it back.

does reading from all app servers for each request work? X bad.



sharding

based on user-id



Simple mod

round-robin

$U_0 \rightarrow S_0$	0	0	$U_5 \rightarrow S_0$	5	1
$U_1 \rightarrow S_1$	1	1	$U_6 \rightarrow S_1$	0	2
$U_2 \rightarrow S_2$	2	2	$U_7 \rightarrow S_2$	1	4
$U_3 \rightarrow S_3$	3	4	$U_8 \rightarrow S_3$	2	0
$U_4 \rightarrow S_4$	4	0	$U_9 \rightarrow S_4$	3	1

10 users
5 servers
 $User-id \% N$

① does L8 now know which req → which server?
✓
↳ Server crash.
(# of servers change)

add a new server \Rightarrow # servers (N) = 6
every server's data has to be changed.
Server crashes \Rightarrow S3 down want
↳ only users going to S3 get affected.



Mapping Table

Inside Load balancer

User-id	serverid
6	1
6	2
6	3
6	4

Range

a new server ✓
server crashes ✓

Cones?? this mapping table
can get very large.

10 servers.

$U_1 \rightarrow U_{100} \rightarrow S_1$
 $U_{101} \rightarrow U_{200} \rightarrow S_2$
 ...

$U_{901} \rightarrow U_{1000} \rightarrow S_{10}$

$U_{1001} \rightarrow U_{1100} \rightarrow S_{11}$

$U \rightarrow S$ ✓

server crashes ✓

only few users
going to that
server are
affected

new server gets
added.

load on existing server
keys increasing.

user-id % N = 5

A ⁰	B ¹	C ²	D ³	E ⁴
0	1	2	3	4
5	6	7	8	9

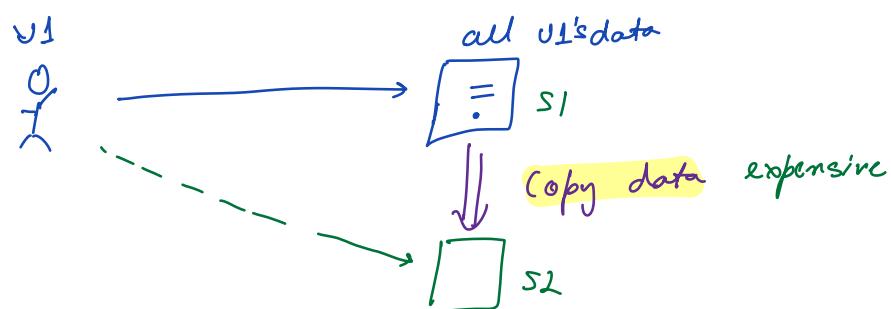
user-id % N = 9

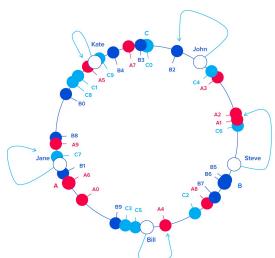
0	1	2	3
4	5	6	7
8	9		

Servers go down

↳ ① which rig goes to which server.

↳ ② restore data (to prevent data loss)
↳ latm (replication & redundant)





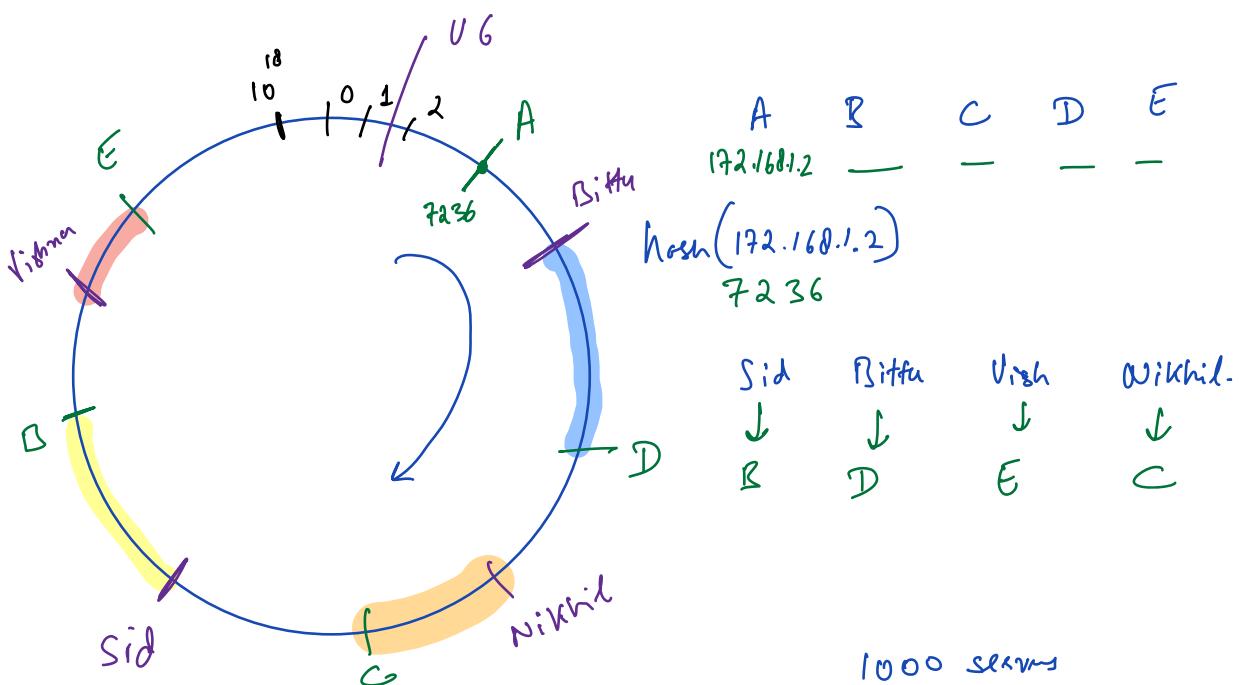
Consistent Hashing

hash both user-id & the server-ids.

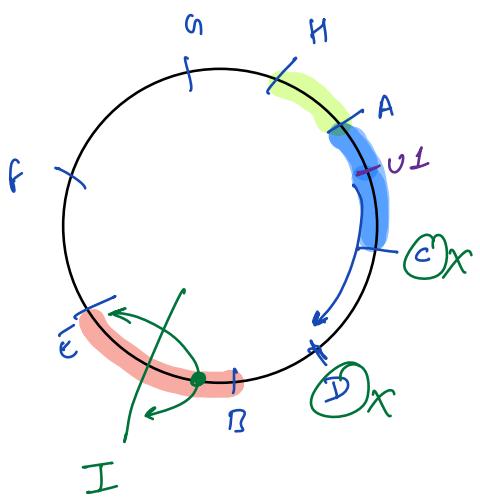
hash-function $\begin{pmatrix} \text{user-id} \\ \text{server-id} \end{pmatrix} \rightarrow [0 \dots 10^{18}]$
 sha256
 md5
 long int

$$\text{hash}(x) = (x^3 + x^2 - 1) \% 10^{18}$$

Pass all server-ids through the hash function.



$$10^{18} \quad \text{Probability of two same hash} \\ \frac{1000}{10^{18}} = 10^{-15}$$



① is load equally distributed?
NO!

② server C crashes.
all of C's load
gets pushed to server
D
D will also crash

Cascading failure!!

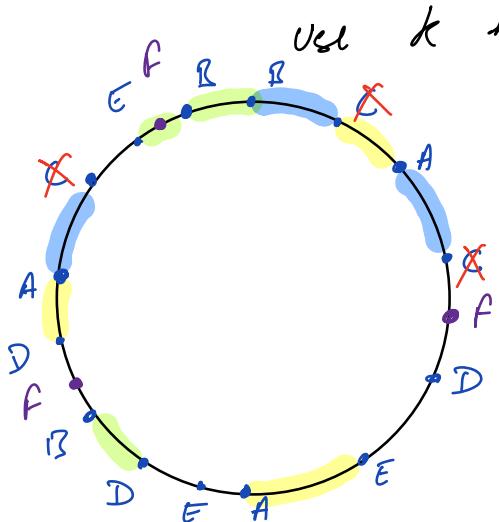
req → user id → server id

mappings should be consistent & easy to compute.

① new server gets added.
automatically some load is given to
that server.

② server goes down

Instead of using 1 hash function for servers
 $k=3$



5 servers $\frac{A}{h1} \quad \frac{B}{h1} \subseteq \frac{D}{h2} \subseteq \frac{E}{h2}$
 $\frac{A}{h2} \quad \frac{B}{h2}$
 $\frac{C}{h3}$

① is load equally divided?
② new server gets added.
load for multiple
server is reduced.



③ Slams crashes
the load gets dist
across multiple servers.

- ① HLD 101 & consistent Hashing
 - ② Caching
 - ③ Caching contd.
 - ④ CAP theorem & Master Slave
 - ⑤ SQL vs NoSQL
 - ⑥ NoSQL
 - ⑦ Case Study — Type ahead
 - ⑧ Multi-master
 - ⑨ Case study — Messaging.
Notification.
 - ⑩ Zookeeper + Kafka
 - ⑪ Case study — Elastic Search
 - ⑫ Nearest neighbor — Quad Trees
file storage (ss)
 - ⑬ Case study — Uber
 - ⑭ Popular Interview Question
 - ⑮ Hotstar
 - ⑯ Microservices
- ① Delicious
- ② theoretical
- ③ \Rightarrow Andur
- ④ \Rightarrow Mudit

Implement Consistent Hashing.

① multiple hash functions.

$$h_1(x) = (x^2 + x) \% 10^{10}$$

$$h_2(x) = (x^3 - x^2 + 2x) \% 10^{10}$$

$$h_3(x) = (\sqrt{x}) \% 10^{10}$$

,

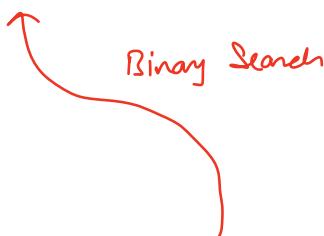
$$h_{bad}(x) = 10$$

②

R A D D < A D P < A D B D C < (n * k)

n servers

k hashes (16, 32)



③ req → hash(user-id) → id

1000 servers

16 hashes

16,000 values

log₂(16000) ≈ 16 steps

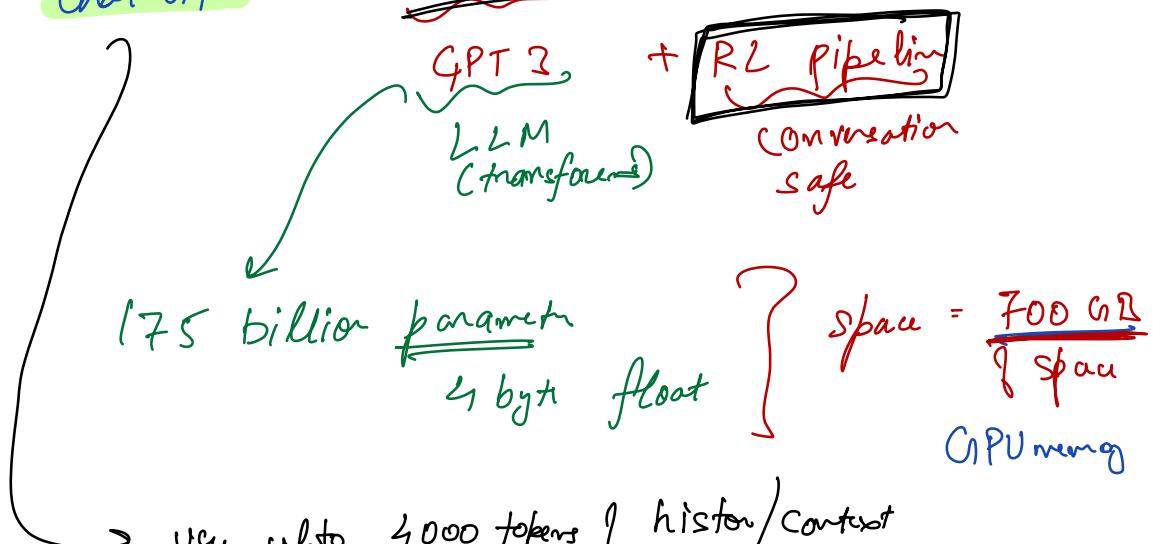
① Consts — IRCTC concurred.

②

Cheat Sheets → when to choose what Database
Common Cloud Provider

8+ years to explain for SDE3 role.

Chat GPT → GPT 3.5 architecture



usu upto 4000 tokens of history/context
↳ sub-word

every req can take $\frac{1s}{8 \text{ GB CPU}}$ compute time
↳ 0.04 \$ per req.
100 million req / day.

OpenAI was burning 5M \$ per day.

git / docker

+91 7351769231

Pragy @ scalar.com

support@scalar.com

