

Recap

Why of HLD

Case Study - Bookmarking website - Delicious

MVP - login/add/view

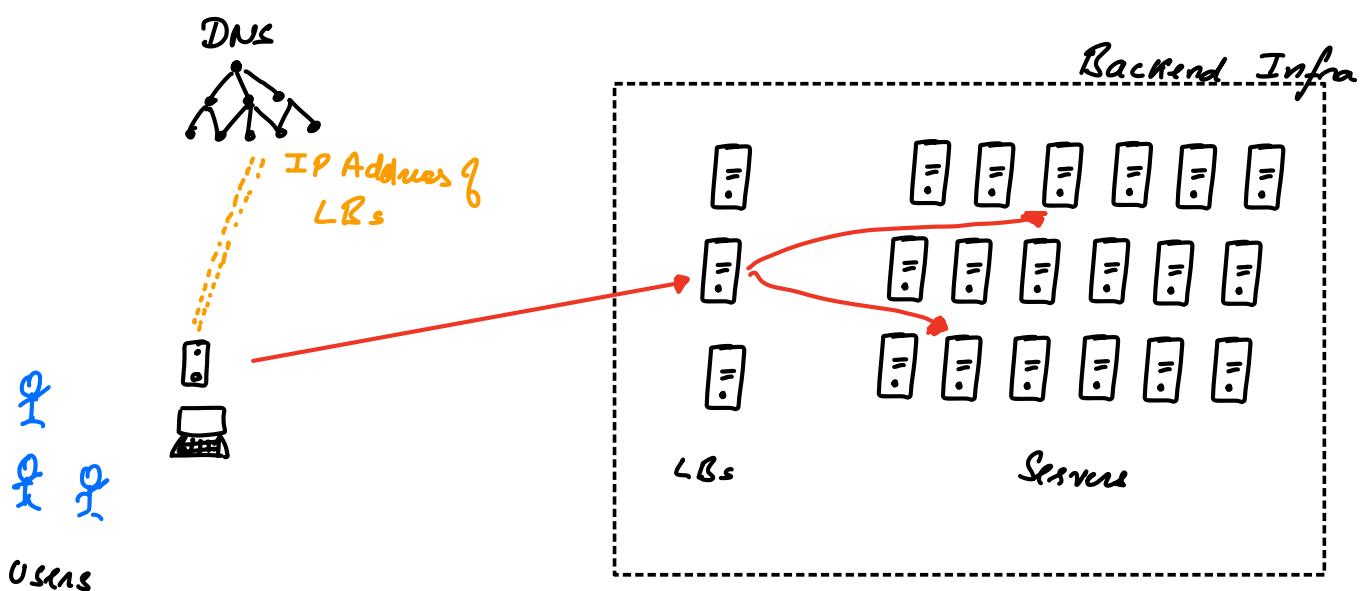
DNS
└ bottleneck
SPOF

Scaling
└ Vertical
Horizontal

Load Balancer
└ unified view of backend
ensures equal load distribution

└ Healthcheck
Heartbeat

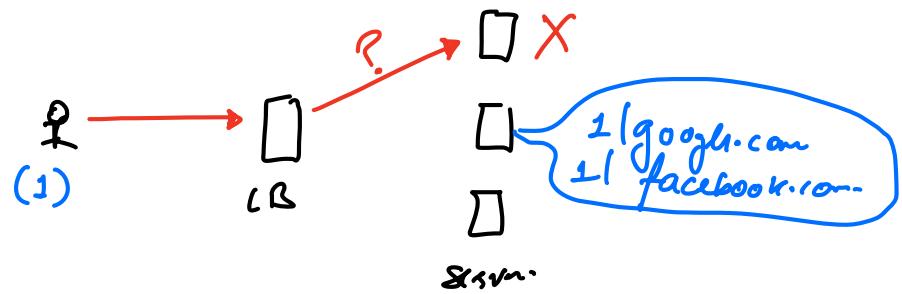
multiple LBs - DNS assists



Q1 Which request gets forwarded to which server?

Routing Algo

Routing will depend on how the
data is stored



Q2 How will you distribute the data? (what data is where?)

① Can we store all data on a single server? No!
(bookmarks)

∴ we need to split the data
Partitioning

Data Partitioning → splitting

① Vertical Partitioning - database Normalization

② Horizontal Partitioning

- ↳ within the same server
- ↳ partition across multiple servers Sharding

DSA / DBMS / LLD

Users

	<u>id</u>	<u>name</u>	<u>phone</u>	<u>avatar-url</u>	<u>email</u>	<u>passhash+salt</u>
1	A	-	-	-	-	-
2	B	-	-	-	-	-
3	C	-	-	-	-	-
4	D	-	-	-	-	-

→ **Vertical (Normalization)**

<u>accounts</u>			<u>profile</u>		
<u>id</u>	<u>email</u>	<u>passhash+salt</u>	<u>acc-id</u>	<u>name</u>	<u>phone</u>
1	-	-	1	A	-
2	-	-	2	B	-
3	-	-	3	C	-
4	-	-	4	D	-

→ **Horizontal (Distribution)**

Users 1

<u>(id)</u>	<u>name</u>	<u>phone</u>	<u>avatar-url</u>	<u>email</u>	<u>passhash+salt</u>
1	A	-	-	-	-
4	D	-	-	-	-

Users 2

<u>(id)</u>	<u>name</u>	<u>phone</u>	<u>avatar-url</u>	<u>email</u>	<u>passhash+salt</u>
2	B	-	-	-	-
3	C	-	-	-	-

can be done on a single server too (to improve write through put)



Screen 1



Screen 2

when done across servers = Sharding.

How to Distribute the data?

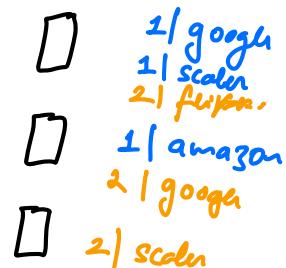
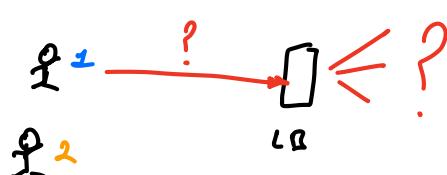
question - write it down
aha moment - take note of
that

① Can we distribute data randomly? No!

∴ there must be some path & logic



1 user's data should be
on only 1 server



fan-out mode

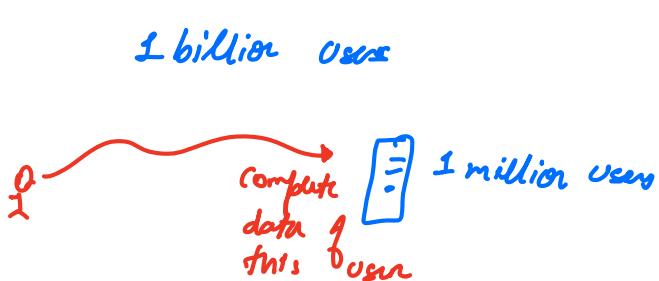


Sharding Key - user-id

Sharding always happens according to a key

each user's data will be on only 1 server.

Note: 1 server can still store multiple user.



1000 servers

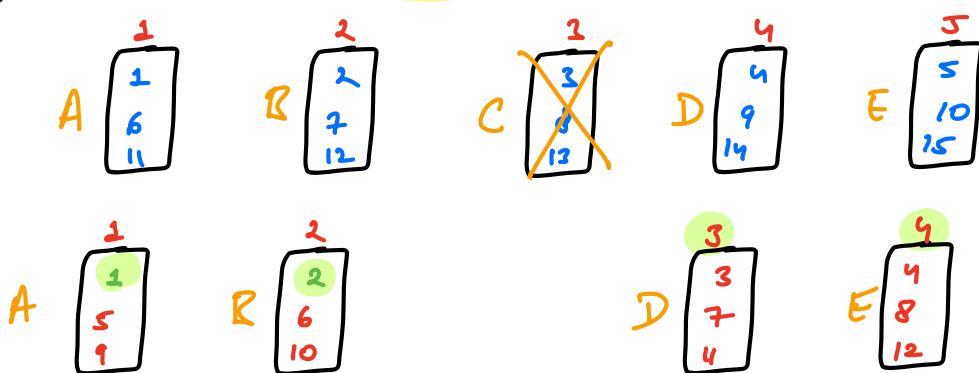


Task: Come up with a good mapping from User-id → to Server-id

this determines how data is distributed
& this determines our Routing algo.

- ① Is the data load distributed **evenly** (**Skewness**)
- ② Is it **easy** to add & remove servers (**minimize data movement**)
- ③ How **fast** is it to calculate the mapping
- ④ Do we need to store anything to calculate this mapping.

~~Attempt 1~~ Round Robin



even?

$N=5^4 \rightarrow$ via healthcheck

fast?

void handleRequest(Request req) {

store anything? No

User-id = req.User-id

Server-id = $(User-id \% N) + 1$

add / remove servers X bad

}

forward to server-id

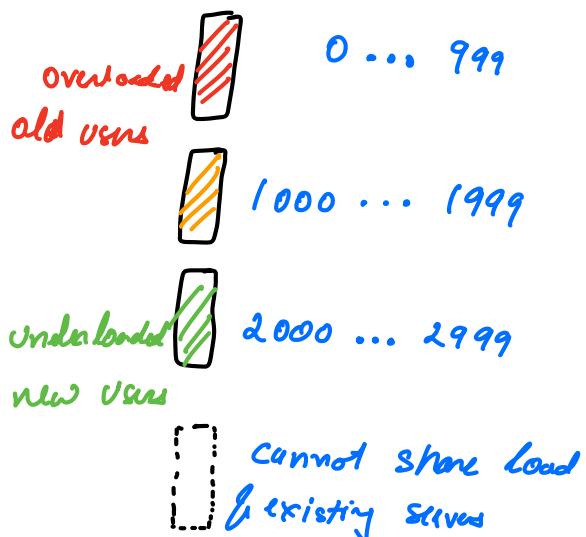
pretty much all user have been mapped to some other server.
a lot of data needs to be moved.

Q How is the data moved? who does that? How do we take backups?

latency class - Magic Fairy (Bon Pari)

↳ Database Distribution
latency class

~~Attempt 2~~ - Range Based

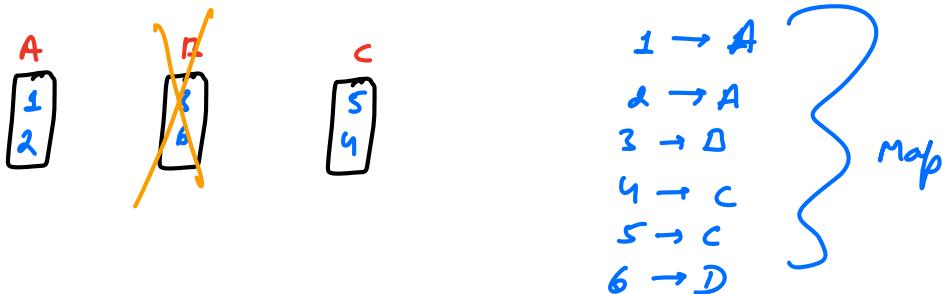


- adding new user mandates adding more servers bad!
- old user generate more load than new users

~~Attempt 3~~ - Hash Map Based

→ The thing about chance is that it is fair — Torver
randomly assign user to servers but
 Keep track of the assignment via a Hash Table

HashMap < user-id, server-id > map



① fast

② even

③ add/remove

loop through map

find all users that are part of
 clustered server & merge them.

~~④~~ store? Yes!

Rad

① Can we store in RAM? No — restarts?
 volatile

∴ we must store on disk — v.v. slow

Random read in RAM $\leq 10\text{ns}$
 (CL latency)

(DDR5
5600 MHz
CL 20)

Random read in HDD $\geq 10\text{ms}$

(7200 RPM
mag disk)

$\approx 2\text{ms}$

RAM is 1 million x faster than disks
 (for random read of 1 Byte)

(Samsung
980 Pro
gen 3)

② How large is this mapping?

Google has 50 users
≈ 100,000 servers

8b 4b
user id → server id

12 bytes / user entry

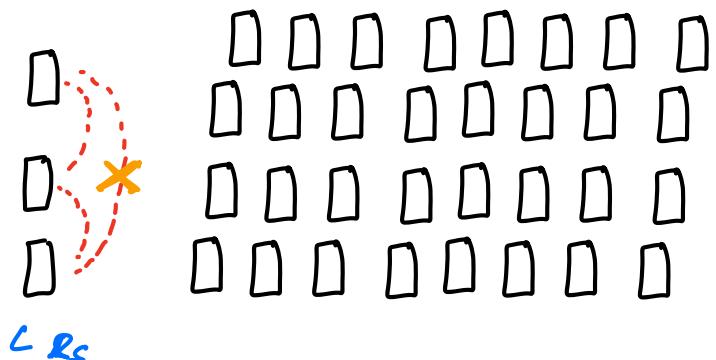
$$\frac{12 \text{ bytes}}{\text{user}} * 5 \times 10^9 \text{ users} = 60 \text{ Billion bytes}$$
$$= 60 \text{ GB}$$

③ What if there are multiple LBs?

is it okay to have diff mapping in each LB? X

mappings should be in sync

n/w overhead



Keeping lots of data in sync

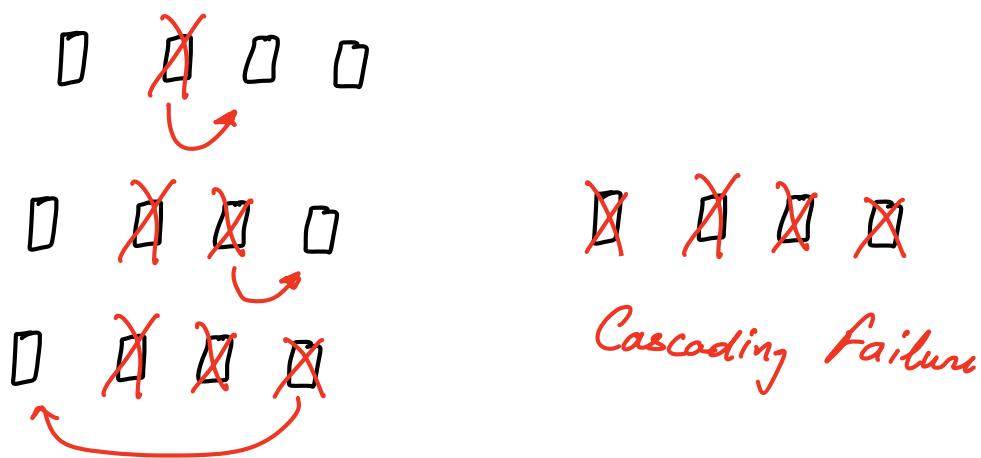
- fast
- consistently
- distributed

IMPOSSIBLE!
(CAP)

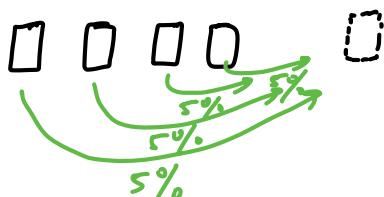
Solution - Consistent Hashing

relies on
Randomness

- ① even load distribution
- ② fast to compute
- ③ won't need to store & sync any mapping
- ④ adding / removing servers will be easy
- ⑤ if a server crashes its load will be handled by all remaining servers equally!



- ⑥ if you add a new server it will automatically relieve the load of all other servers equally.



Hash Functions

input → output
domain range

$$f(x) = x^2 + 3$$

- deterministic

No matter how many times you call it it always produces same result.

a function which has a finite range

& behaves pseudo-randomly is a hash difficult to predict what the output will be

Mapping (inp → out) is random but deterministic.

hash (x)

- represent x in digits
- └ split it in middle
- └ reverse each part
- └ cube it
- └ join
- └ mod M
- └ $\% (10^9 + 7)$

1224567
↓
1234 | 567
4321 | 765

00677568161 | 947697125

00677568161 947697125

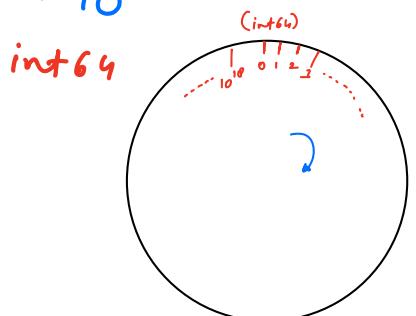
[]
b/w 0 ... (M - 1)

family of hash functions

Consistent Hashing

k hash functions for servers
 1 hash function for request
 (User-id)

} the range for all
 these is the same
 $0 \dots 10^{10}$



Algo

① for each server

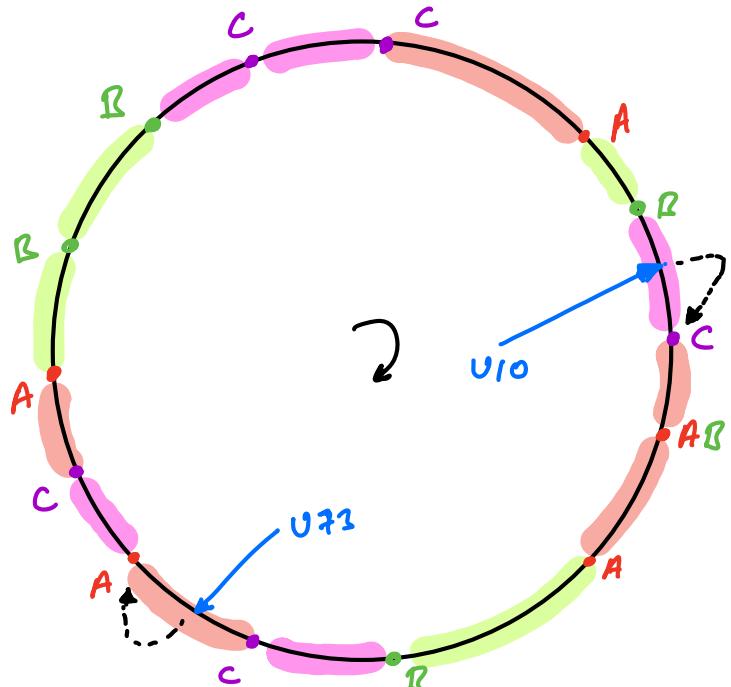
pass it through all k hashes
 each hash value will be a spot for this
 server on the ring

virtual spots

② when a req comes

extract User-id sharding key to pass it through the
 hash k to note the spot on ring

③ forward req to the nearest server moving clockwise



$N = 3$ servers A B C

$K = 5$ hash functions

$h_1(A) = 10$	$h_2(A) = 60$
$h_2(A) = 30$	$h_4(A) = 70$
$h_5(A) = 35$	
$h_1(B) = 15$	$h_2(B) = 80$
$h_2(B) = 20$	$h_4(B) = 85$
$h_5(B) = 50$	
$h_1(C) = 25$	$h_2(C) = 90$
$h_2(C) = 55$	$h_4(C) = 5$
$h_5(C) = 65$	

even

∴ of randomness

as K increases it becomes more even

$$K = 16 \quad \text{or} \quad K = 32$$

fast

① Stop anything? NO

$$\begin{array}{ll} h_1(A) = 10 & h_2(A) = 60 \\ h_2(A) = 30 & h_3(A) = 70 \\ h_3(A) = 35 & \\ h_1(B) = 15 & h_2(B) = 80 \\ h_2(B) = 20 & h_3(B) = 85 \\ h_3(B) = 50 & \\ h_1(C) = 25 & h_2(C) = 90 \\ h_2(C) = 55 & h_3(C) = 5 \\ h_3(C) = 65 & \end{array}$$

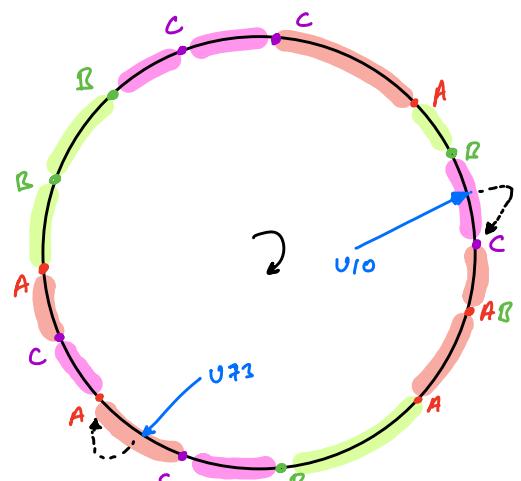
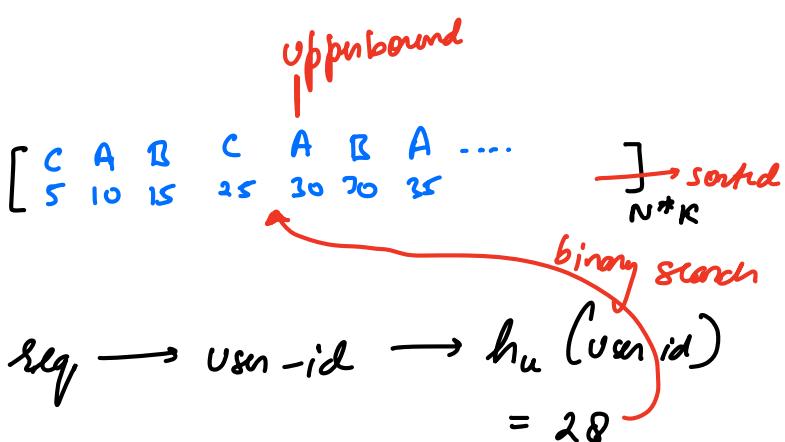
array = []

servers = [...] // healthcheck →

for each server

serverid → each hash → preconfig
→ values hash function
→ store value in array

Sort array by keys



$$\lg_2(N * K)$$

$$N = 100,000 \quad \left. \right\} \text{Google}$$
$$K = 32$$

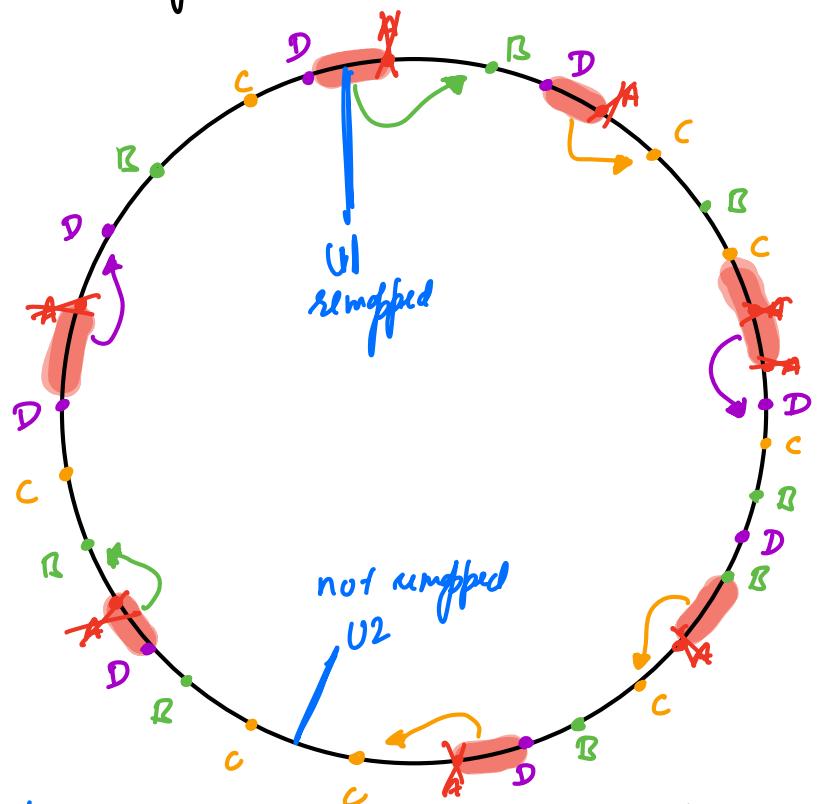
$$\text{Time } \lg_2 \left(\frac{10}{3,200,000} \right) \approx 21 \text{ ops} \approx \text{few } \mu\text{s}$$

$$\begin{aligned} 2^{10} &\approx 10^3 \\ 2^{20} &\approx 10^6 \\ 2 \cdot 2 \pi &\approx 10^6 \end{aligned}$$

$$\text{Storage} \quad 3,200,000 \text{ entries} \quad \text{if each entry is } 12B \quad = \frac{36MB}{RAM}$$

(serverid + hash)

What if server crashes?



server A dies!

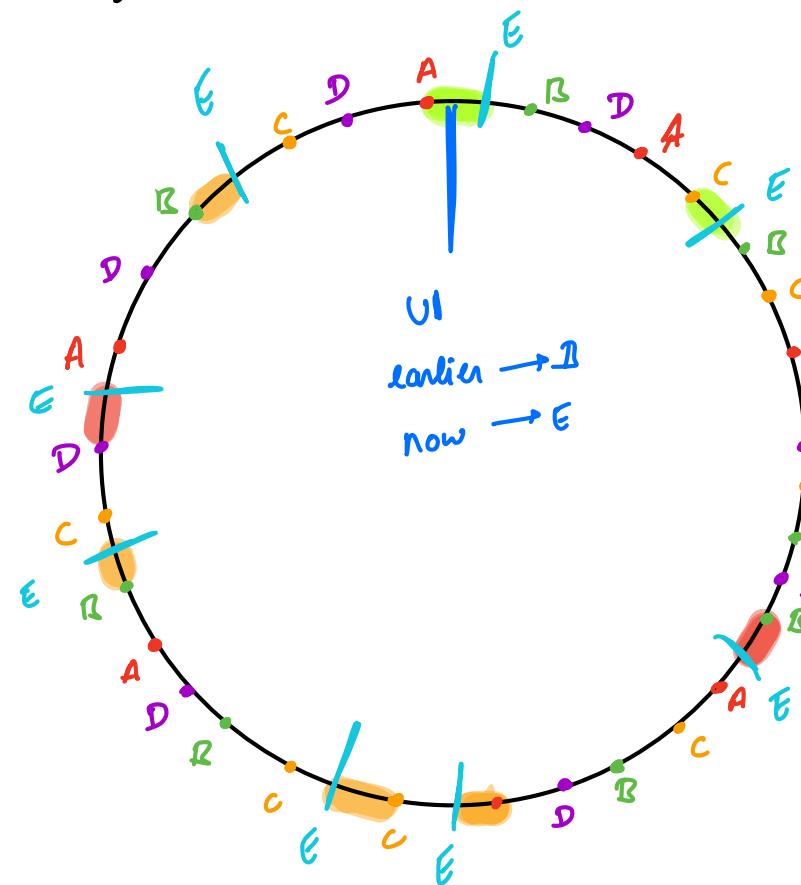
- ① LB will detect it
(all of them)
- ② Recompute from array
- ③ Continue working as if nothing happened!!

(Note : Data Movement
Magic Fairy)

the load of server A was taken over by servers
R C D

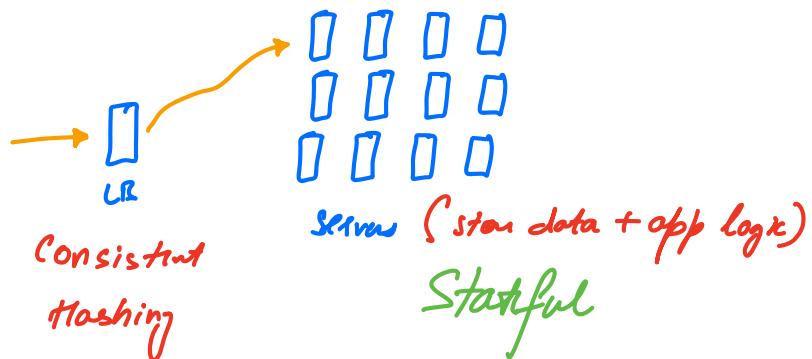
If a new server is added

add server E



It will take over some load from the all remaining servers

Consistent hashing → used when distributing data.



State = data

Stateful → containing data

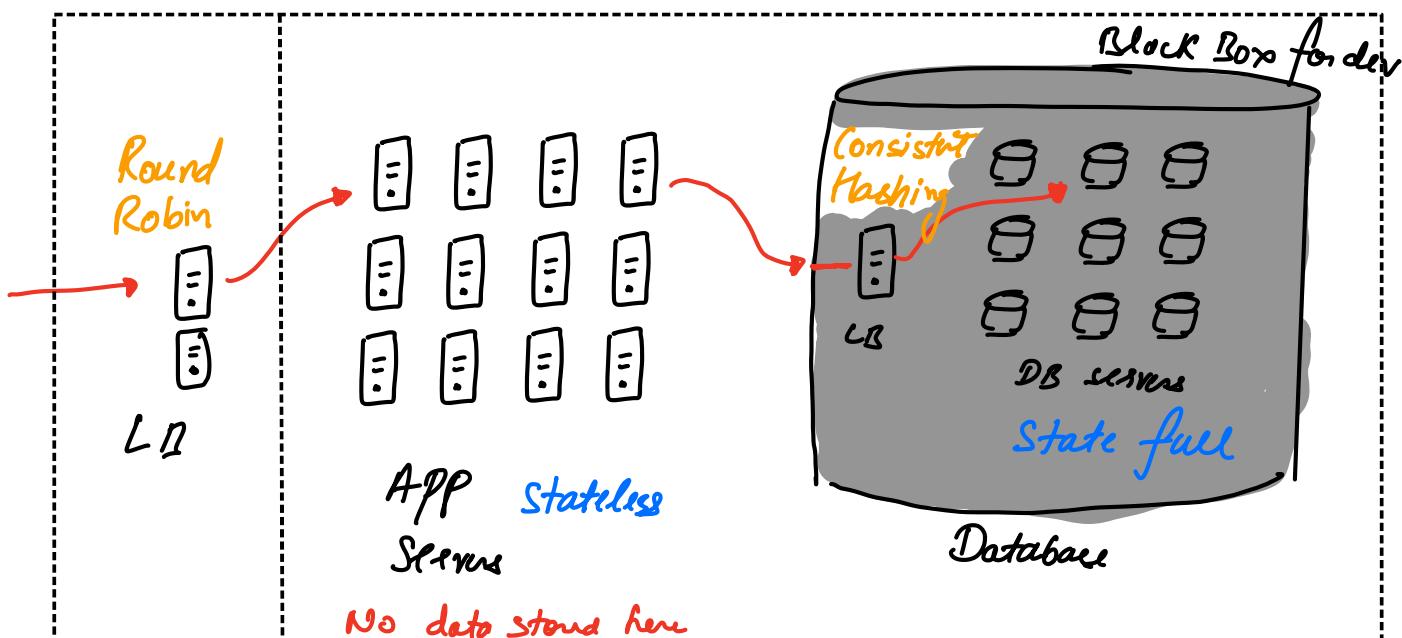
Modern System Design

Push data away from app servers into DB layer.

why?

if the same server has both app logic & data

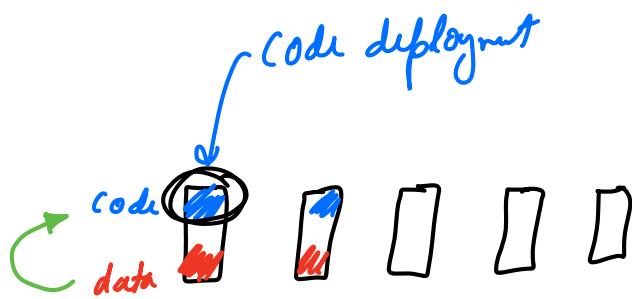
- ① tightly coupled — cannot scale independently
- ② code deployment — temporary data outage!



only app logic

Untrusted
n/a

Trusted VPC



CI/CD

→ @Microsoft

≈ 9000 code deploy / day

Next class — Caching

- ① What is why — types
- ② Case studies

Please ensure that you're done assignments — 5 MCQs