

Agenda

- Designing Orchestrator
- Multi Master

9pm - 11:30pm

10-15 mins

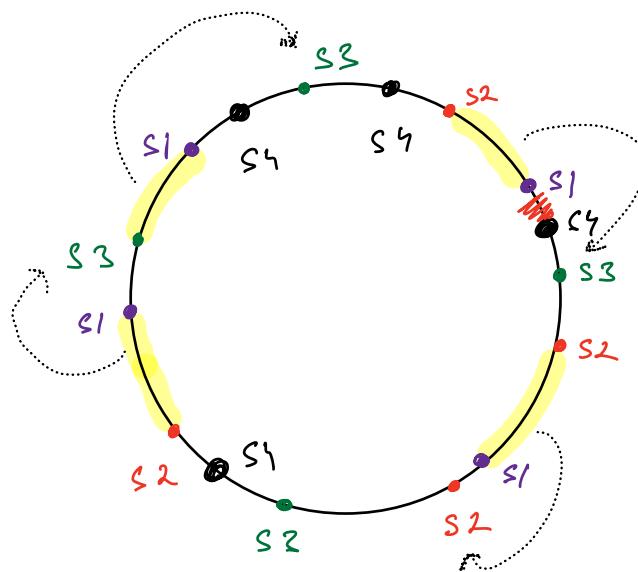
doubt session → 1pm

SQSL / No SQSL

↳ sharding happens automatically
How?

Consistent Hashing

→ what happens when a server goes down → consistent hashing automatically redistributes load
But how is data copied?



Sharded based on
Some key → user-id

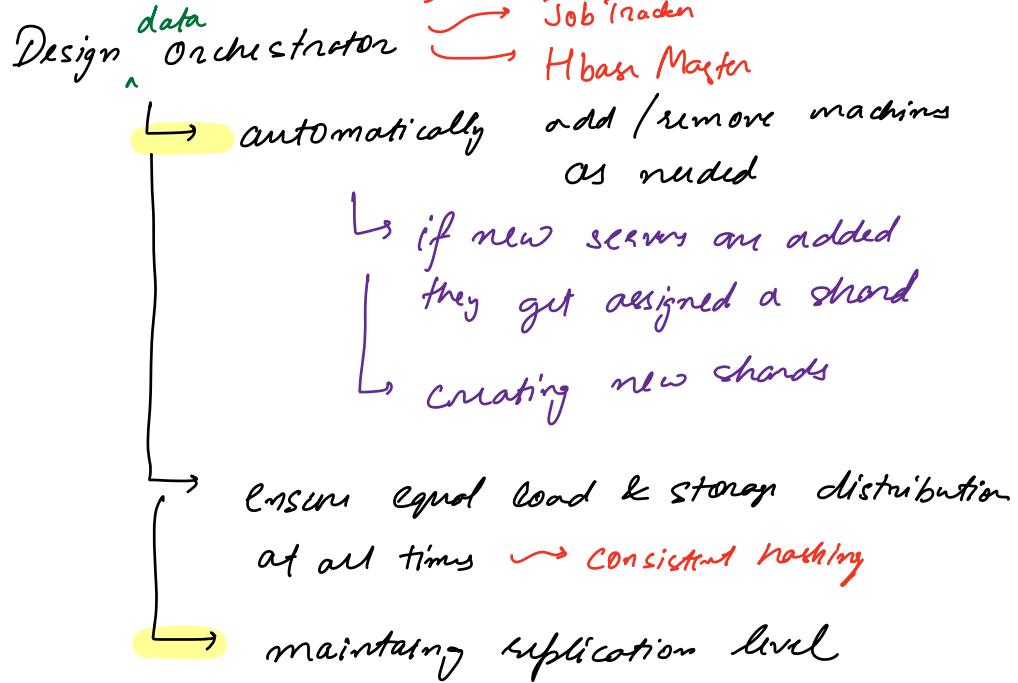
3 shards → S1, S2, S3

k hashes → 4

if down

Replication → how do we maintain replication?
↳ master-slave

Problem Statement



Given config

- ↳ sharding key → user-id
 - ↳ replication level → ?
 - ↳ 10 servers
- Orchestrator
manages everything else

Suppose \rightarrow 10 servers are already working

given 2 new servers
idle

Monitor memory/disk utilization for each shard

↳ cross a threshold

↳ move some data to new server.

Reacting to load on a shard

Drawbacks :

If we let sharding happen at
any time (reacting to load)



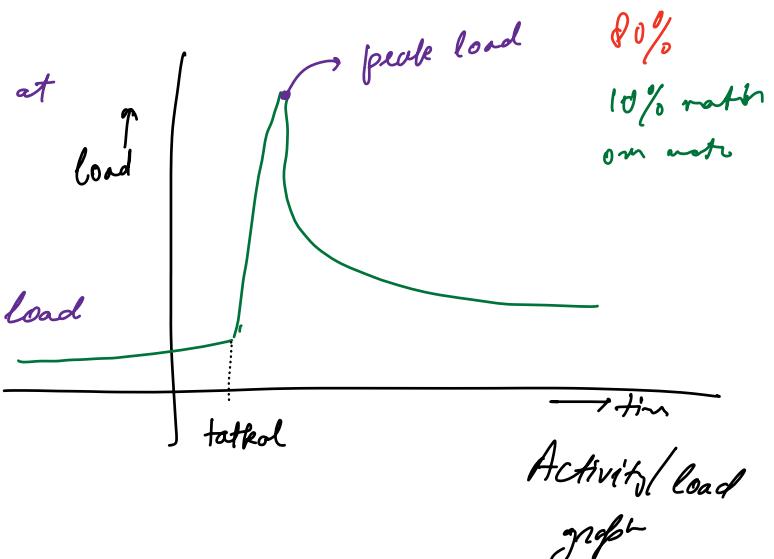
Create shard during peak load



move data from 1
machine to another



more traffic \rightarrow worse situation



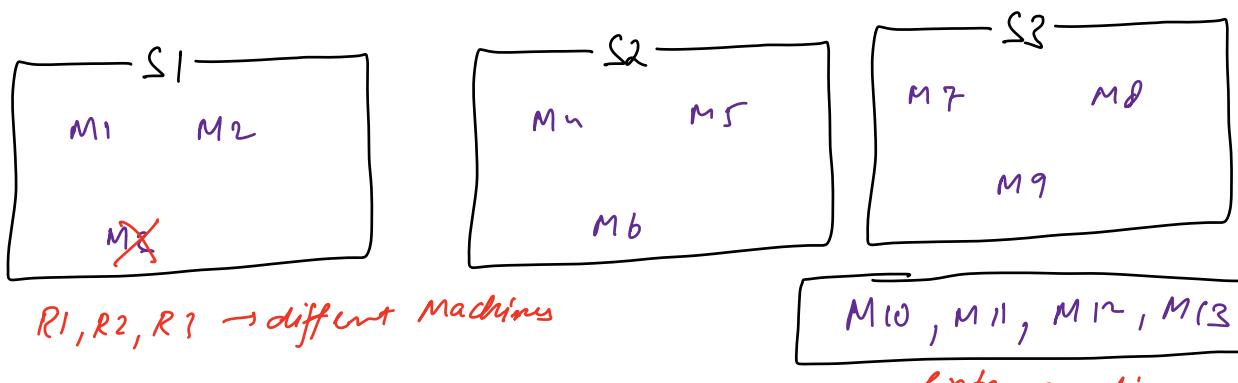
- Let us not react to load at time of load
 - when we have extra servers available → assign a role to them
-

Replication levels

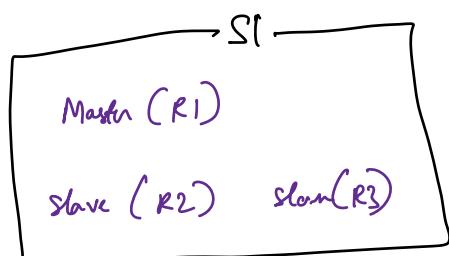
→ 3

3 shards → S₁, S₂, S₃

Rajeev DB



Master-Slave for replicas



highly available (inconsistent)
 ↳ write → Master

eventually consistent
 ↳ read → any slave

highly consistent
 ↳ write → Master + 1 slave
 ↳ read → any slave

highly consistent
 ↳ writes → all slaves + master.

Asymmetric Ranking System
↓
highly available
& eventually consistent

Questions

How to add or remove machines?

- ↳ if extra machines are available then
 - what should system do?
- ↳ what happens when a machine dies

Possibilities when you have extra machines

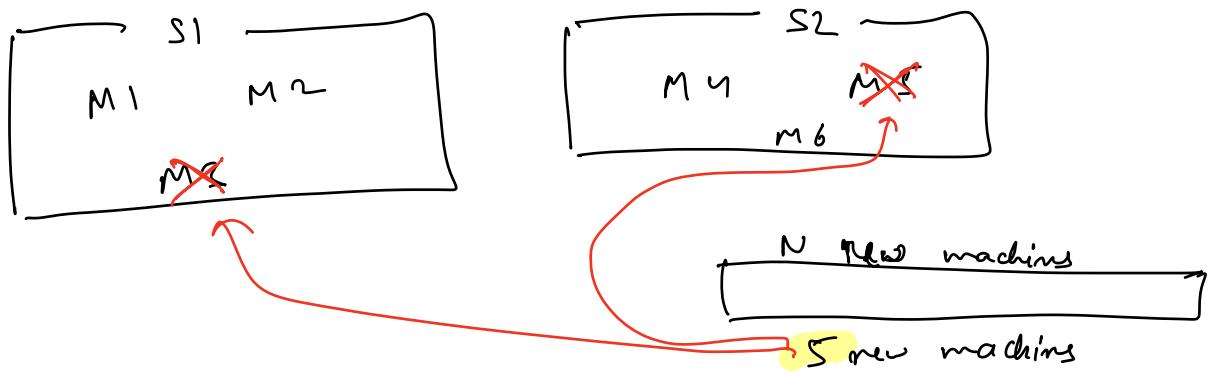
replication level = 3
of shards = 3

- Create new shard ($N \geq$ replication level)
- add machine to existing shard
- leave them idle

Ashnith's Algorithm

- ① N Extra machines
 $N // R$ new shards
 $N \% R \rightarrow$ idle

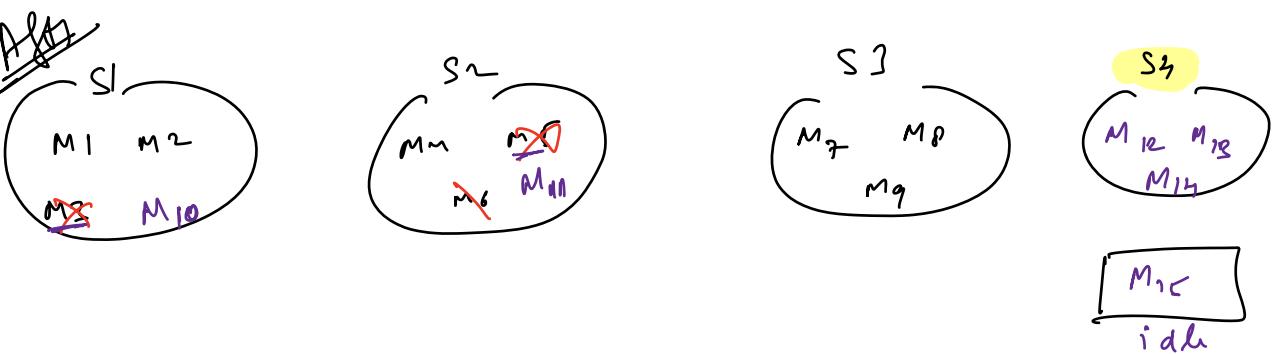
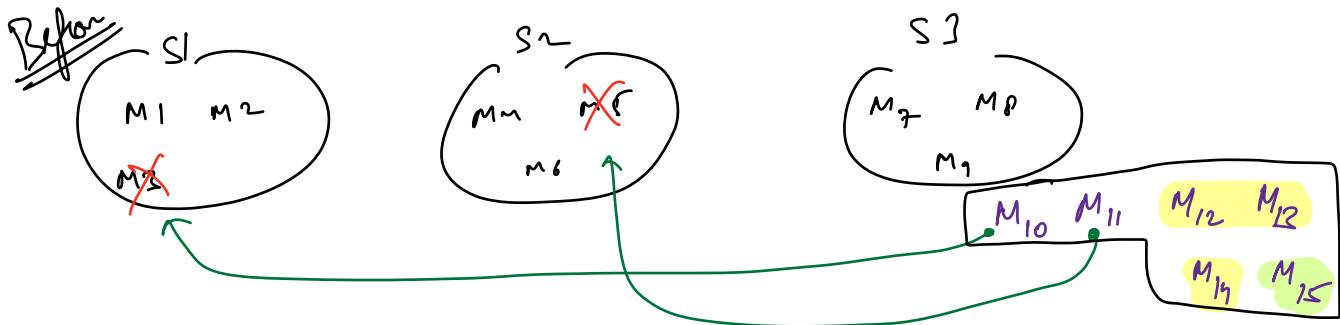
$n = 1001$
333 new shards
2 machines remain



Improved Algorithm Algo

- ↳ priority → maintain data consistency (don't want to lose data)
 - replace downed machines
- ↳ if there are machine left over after replacing downed machines
 - ↳ create $N/3$ new shards
 - ↳ $N \% 3 \rightarrow$ idle

Situation



Insight

↳ don't immediately create new shards if you have enough machines

↳ "... some machines might break in fact"

at all times we should have \times machines available just in case some machines in a shard break.

How to calculate how many machine should be have as backup?

1000 shards \rightarrow 1000 machines available as backup

Machine \rightarrow goes down once every
2 months \rightarrow 2 days to repair

Big Table

NoSQL

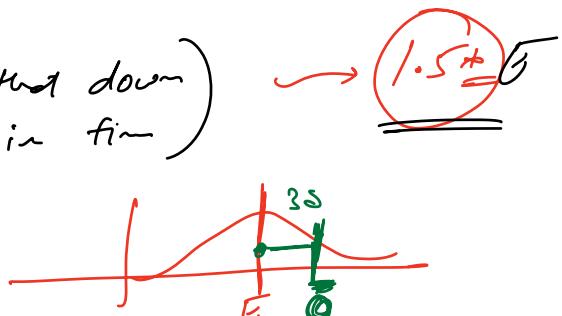
in a period of 60 days \rightarrow available for 58 days

$$\text{Prob that a machine is down} = \frac{2}{60} = \frac{1}{30}$$

S shards each with Replication level of R

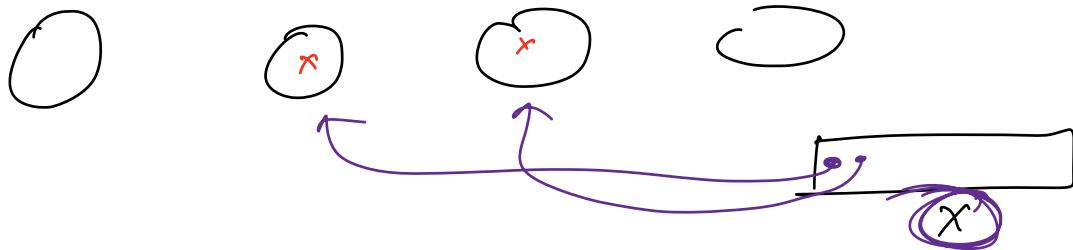
\hookrightarrow total $S + R$ machines

Calculate \rightarrow E (machines that down)
 \downarrow
avg.



\rightarrow how many backup mach. should I have
so that with 99% prob I have
sufficient backups

whatever remains \rightarrow create a new shard if
 $(N - x)$ possible



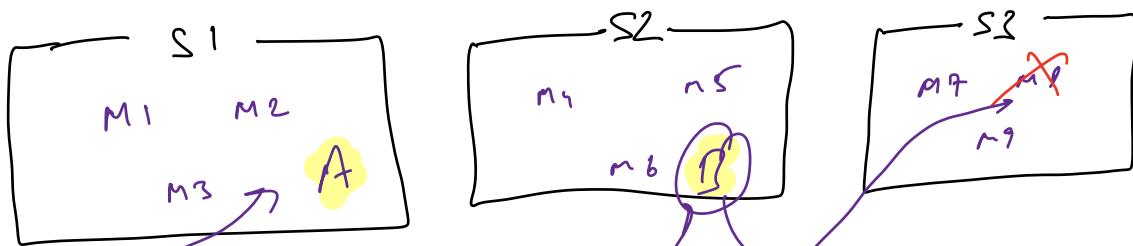
idle machines?

$$(N - x) / 3 \rightarrow \text{shard}$$

new baday

$$(N - x) \% 3 \rightarrow \text{redundancy}$$

x machines for backup \rightarrow idle?



Replication factor = 3

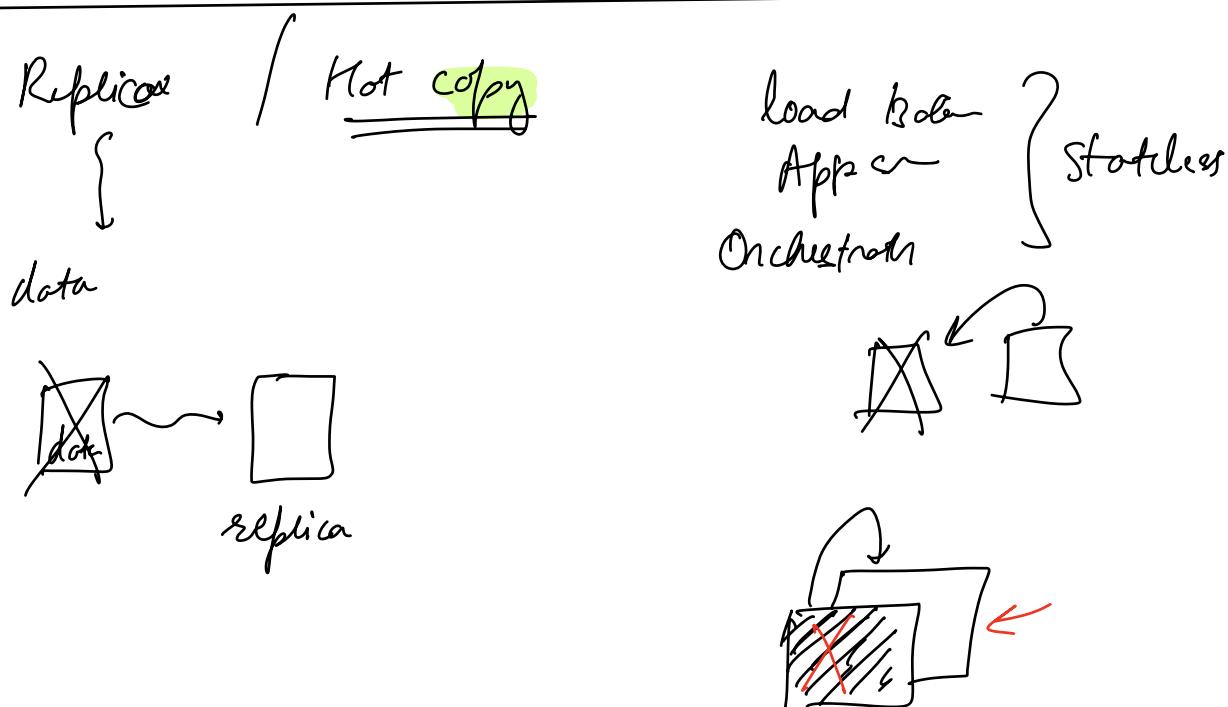
x (backup) = 2 machines

15 extra machines

13 rest \rightarrow 2 new shards \Rightarrow 12

3 mch's → idle?
 → work?
 A, B, C

Instead of having idle machines
 ↳ use them as extra replicas
 ↳ contract: if a machine goes down in
 some other shard
 this extra replica will be
 claimed.

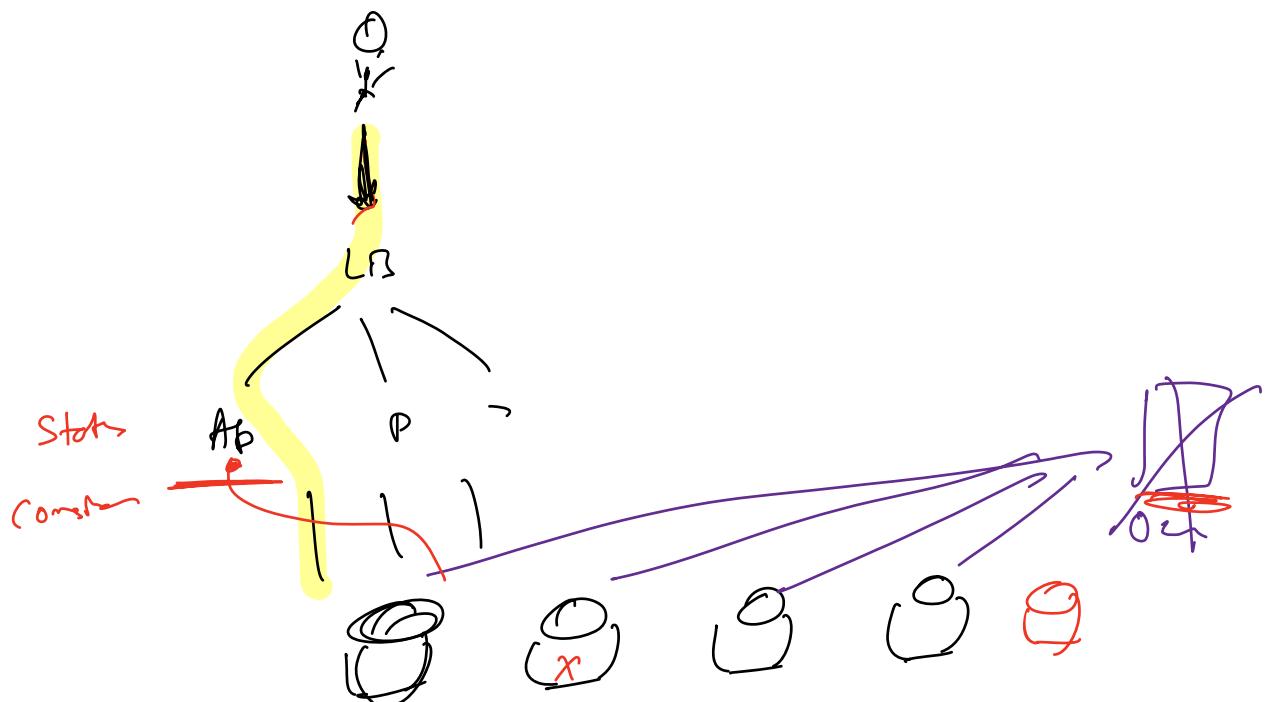


How is data moved

10:40 → 10:50

↳ new shard

↳ a new read replica is added?

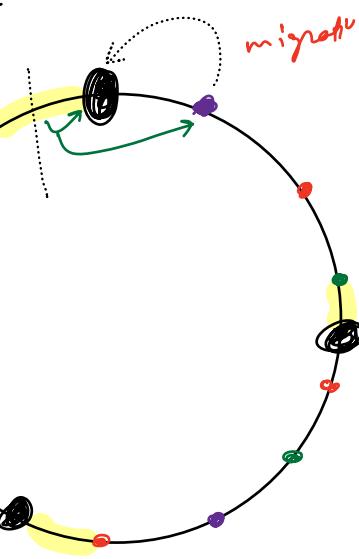


Seamless Shard Creation

- new shard \rightarrow 0 data
- data migration \rightarrow 2 step

Cold start

S₄ shard
has to be
added.



Staging Phase

- \rightarrow don't really add shard S₄ to the consistent hashing
- \hookrightarrow upper layers (db client) no idea that new shard will be added.

Simulation

- determine what keys (User-id) $\xrightarrow{\text{range of keys}}$ will be directed to new shard once it gets added.

- \rightarrow Start copying data \rightarrow Starts at 11pm
 $\left\{ \begin{array}{l} \\ 15\text{min} \end{array} \right.$

↓
11:15

→ shard is warmed up.

However any write requests b/w
11 pm & 11:15 pm still go to
old shard

So new shard is missing that data.

Real Phase

↳ shard actually goes online
(db client are now
aware of a shard)

→ we need to catch up with the
data from 11:00 pm to 11:15 pm

↳ Available → new shard start serving
requests

copying this
will take another

1 min

Consider → 1 min → reject any
requests

Timeline

$$T_1 \quad (\parallel pm) \quad \rightarrow$$

Orchestrator decides that new shard will be added

Simulation Phase

Start copying data.

T2 (11:15 pm) → copy done

New shard → live

T3 (11:16 pm) → Δ update on completion
ful old stand to new
data

SQL / NoSQL

→ Write ahead Log

WAL

$$x = 10$$

W

$$\underline{x} + = 10$$

w
R

$$\begin{array}{l} \text{oldest} \\ x = 10^{15} \\ y = 20^7 \\ w = 20 \\ z = 30 \end{array}$$

118

$$\begin{array}{|c|} \hline \text{With } x = 15 \\ \text{With } w = 20 \\ \text{With } y = ? \\ \hline \end{array}$$

$$\begin{cases} x = 15 \\ y = 20 \\ z = 20 \end{cases}$$

11:1

WAL → general DB operation to provide

- ↓
- Durability
- Atomicity

Roll back in case
Power failure

→ improve with speed (SQLite)

~~X backup machine~~ → estimation X can be challenging.

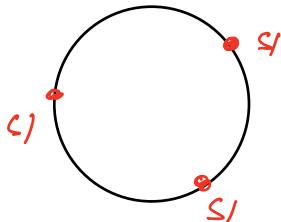
Dynamo DR / Cassandra

Multi-Master

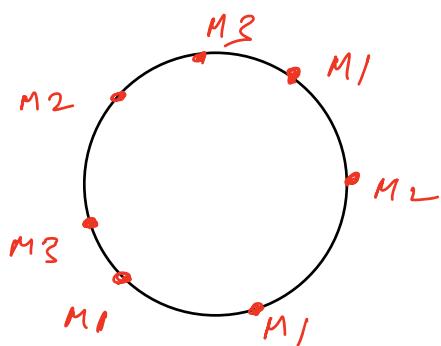
→ no slaves

→ every machine is a master

Central

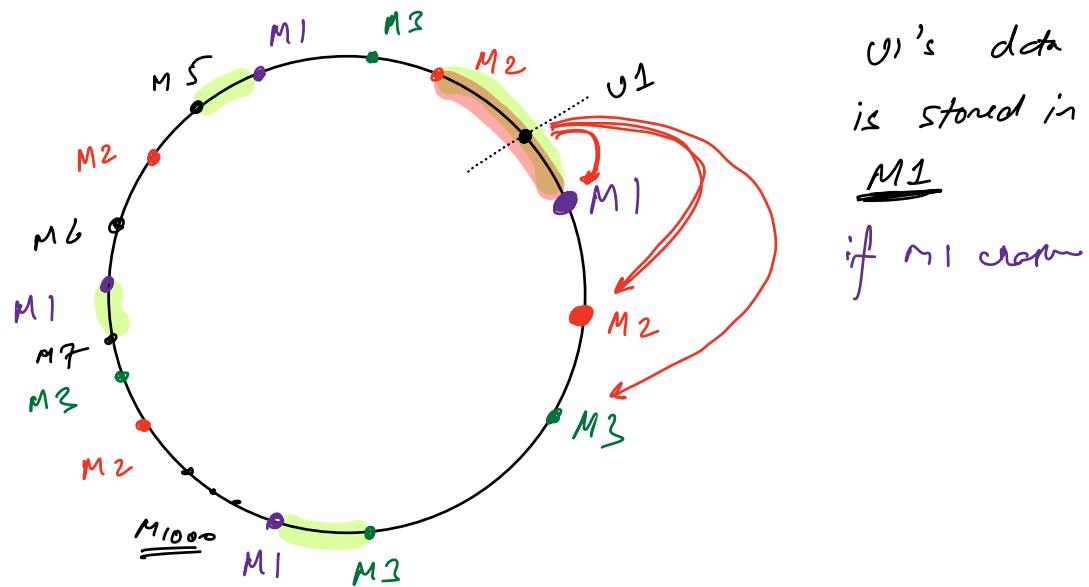


multi-master



machines are not grouped into shards → each shard is just 1 machine

→ Imagine replica factor is ?.

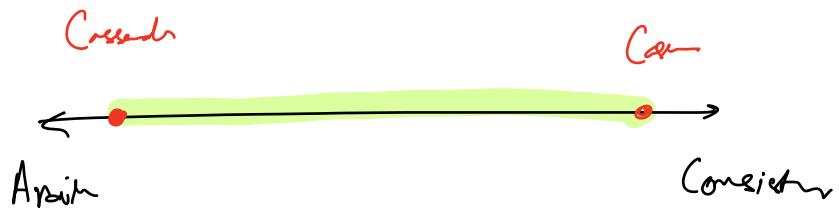


→ Optimal to store replicas in servers which are next in line to serve UI's request

→ Note: M2 is not a replica of M1 X

if stores replicas for any users whose req will go to M2 in case M1 dies.

Tunable Consistency



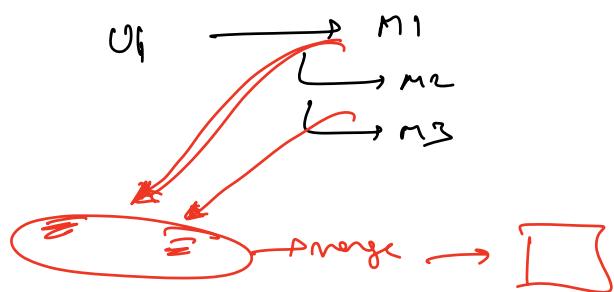
Configure two variables $R \& \omega$ ↗ Min # units
↓
Min # reads

$\chi \rightarrow$ replication level

$$\underline{\chi = 3}$$

$$\underline{R = 2}$$

$$\underline{\omega = 2}$$



$\chi = 3$

$R = 1$ how many copies must be read

$\omega = 1$ how many copies in tx with

M1	M2	M3
	$a = 1$	
		$b = 2$
	read(a)	efictor
$a = 1$ $b = 2$	$a = 1$ $b = 2$	$a = 1$ $b = 2$

UI 'data' \rightarrow write $a = 1$

\rightarrow write $b = 2$

\rightarrow read (a) \rightsquigarrow null

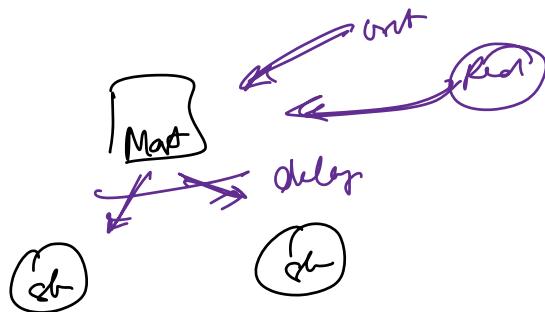
(consistent / Available?)

$R \uparrow \Rightarrow$ higher consistency
 $\omega \uparrow$

If R is high then before successful read
for read to access multiple copies
if copy is down \rightarrow read unsuccessful
More consistent but less available

$R + \omega \rightarrow$ the higher this is, the more consistent your system will be.

MAP → if you make a read → data is the latest data

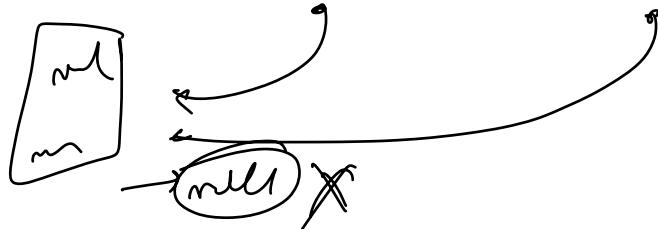
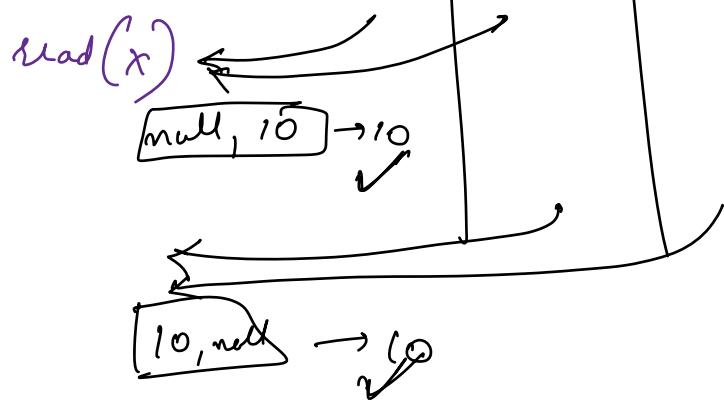


$$x = 3$$

$$R = 2$$

$$W = 1$$

m_1	m_2	m_3
	$x = 10$	
		$y = 20$

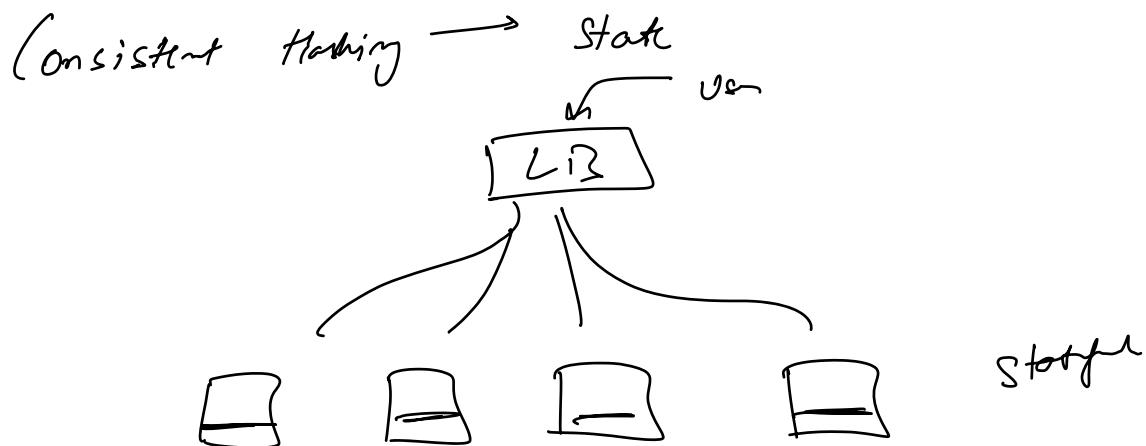


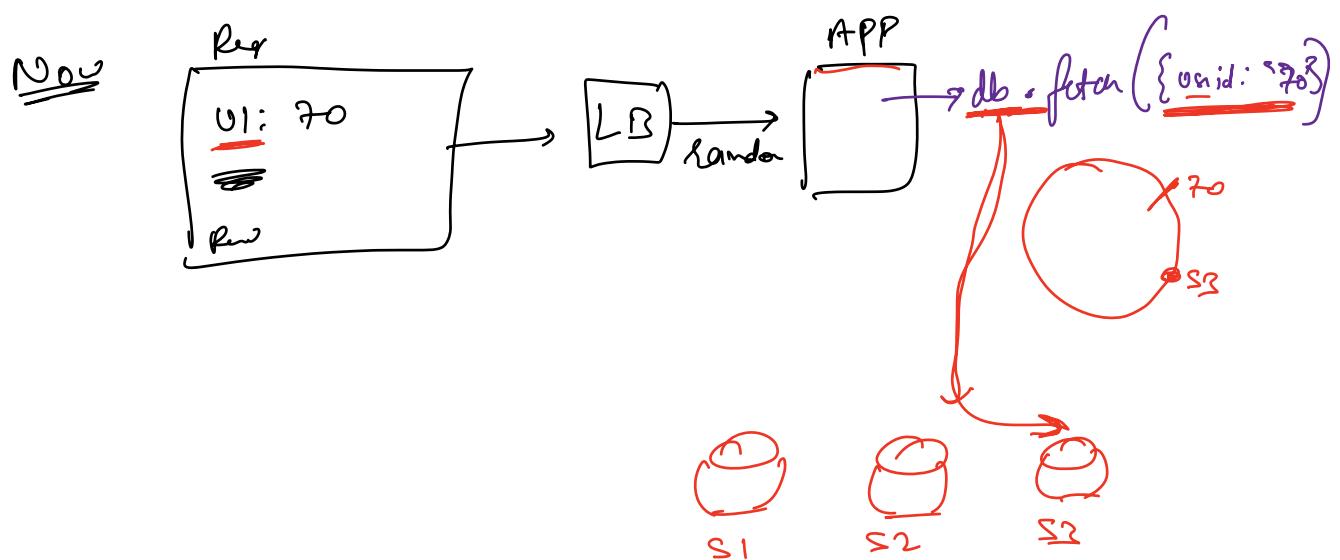
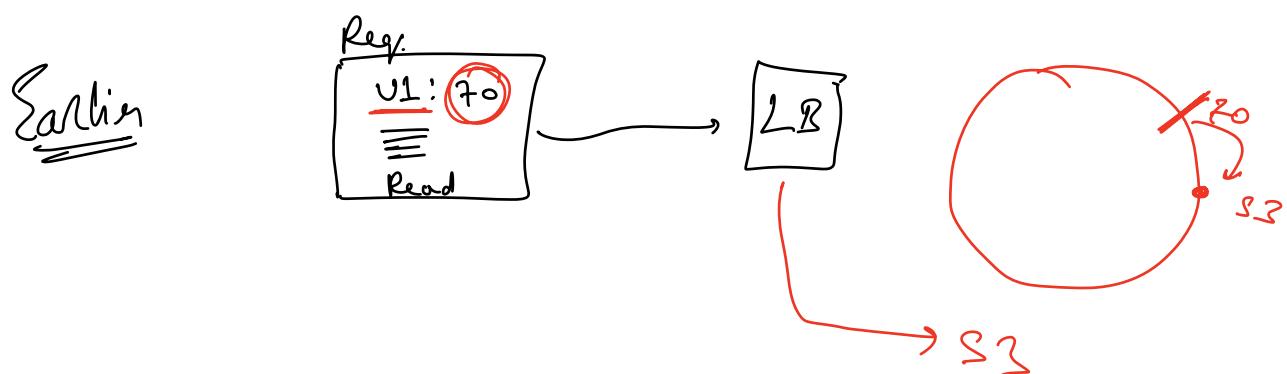
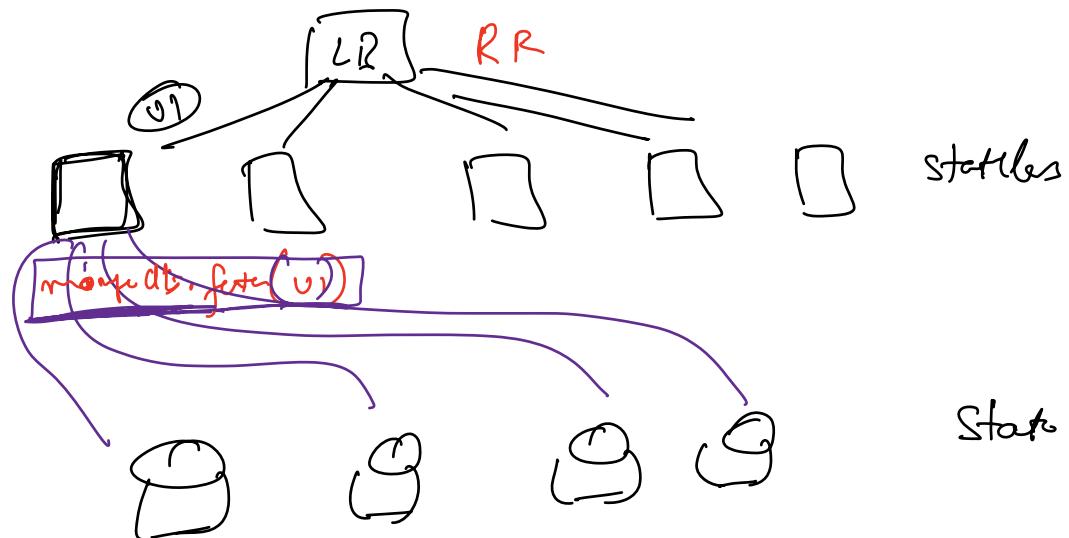
$R + \omega > \underline{x}$ \longrightarrow highly consist. syst.

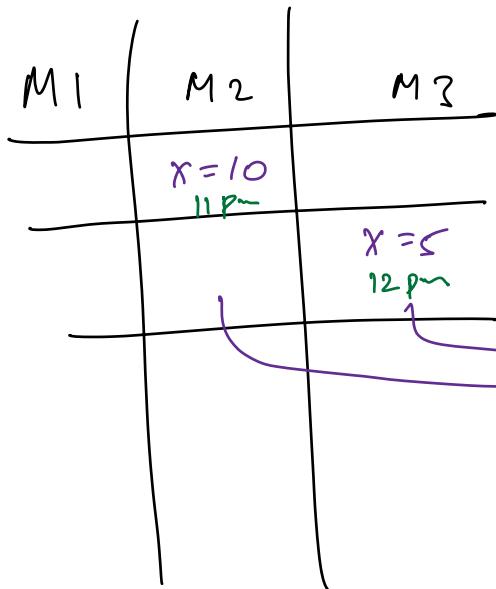
$$x = 10 \quad R = \underline{5} \quad \underline{\omega = 6}$$



8th node \rightarrow Ward







$$x = 2$$

$$\omega = 1$$

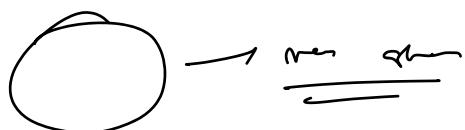
$$R = 2$$

read

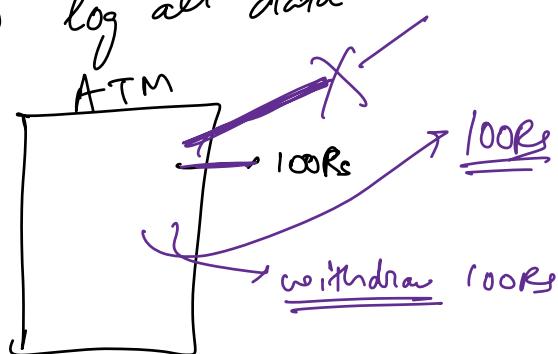
$$m_2 \rightarrow x = 10 \quad 11pm$$

$$m_3 \rightarrow \cancel{x = 5} \quad 12pm$$

N new reads $\rightarrow X$ bytes

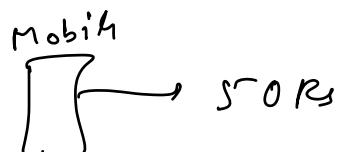


Banks log all data



$$100Rs - 50$$

Overdraft
due 10Rs



Server push

app → websocket connection with server.

whatsapp → 2 million ws

per server.