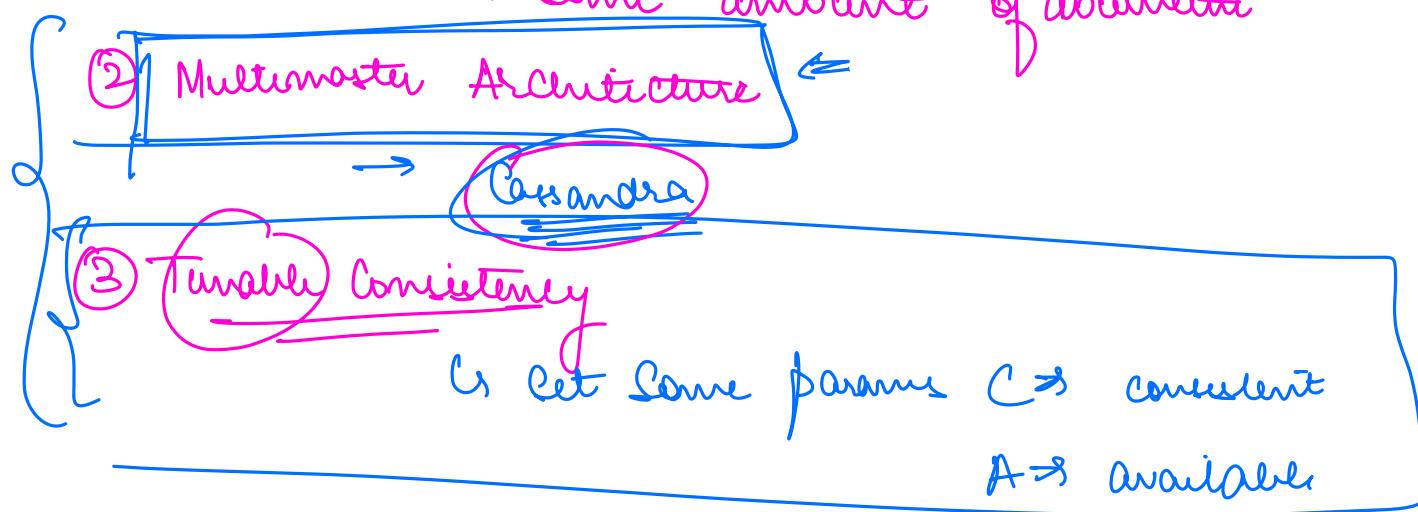
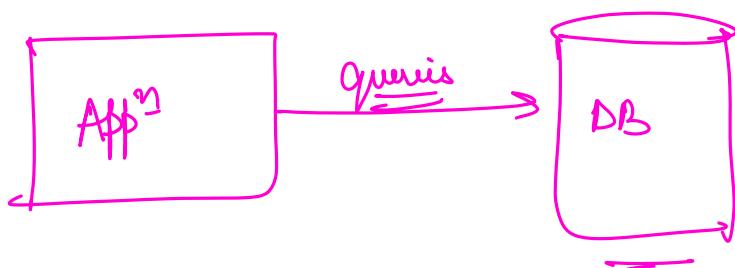


## Complete NoSQL

① Problems with Master Slave and Charding  
↳ some amount of downtime



①

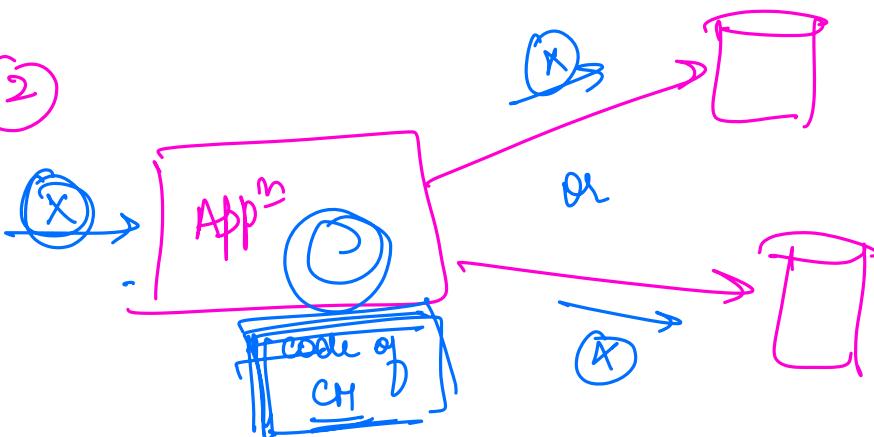


bottleneck

(i) size of data

(ii) # of queries

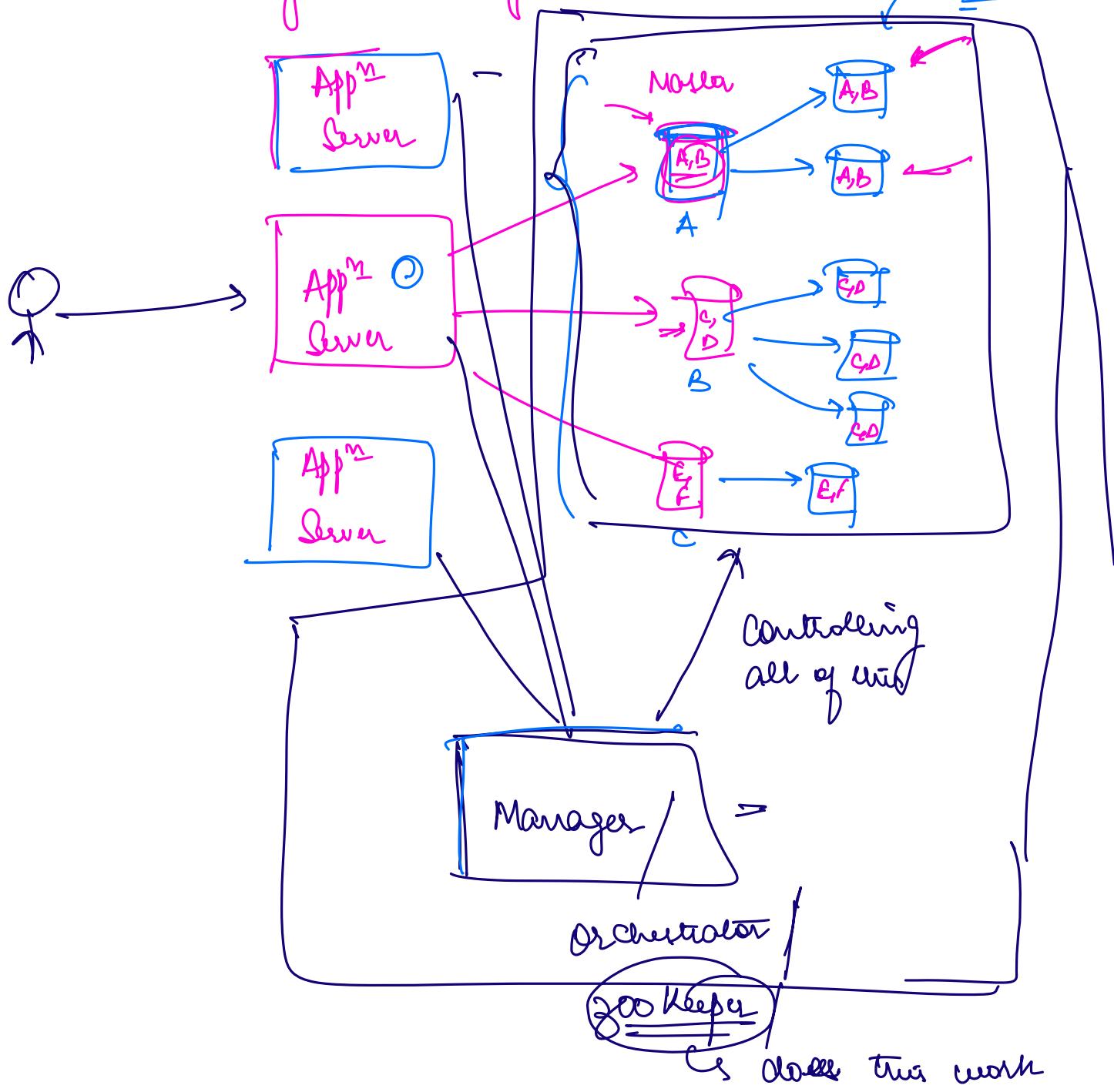
②

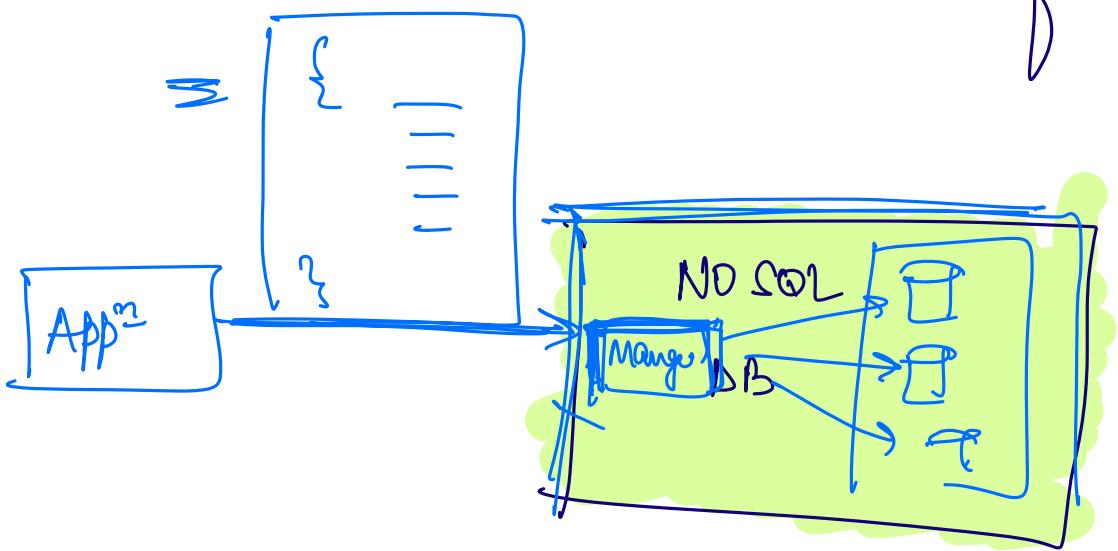


For DB by default don't support sharding / replication.

There has to be someone who does!

- (i) creating a new shard if needed.
- (ii) Create another replica if needed
- (iii) Given a key, tell which shard has it

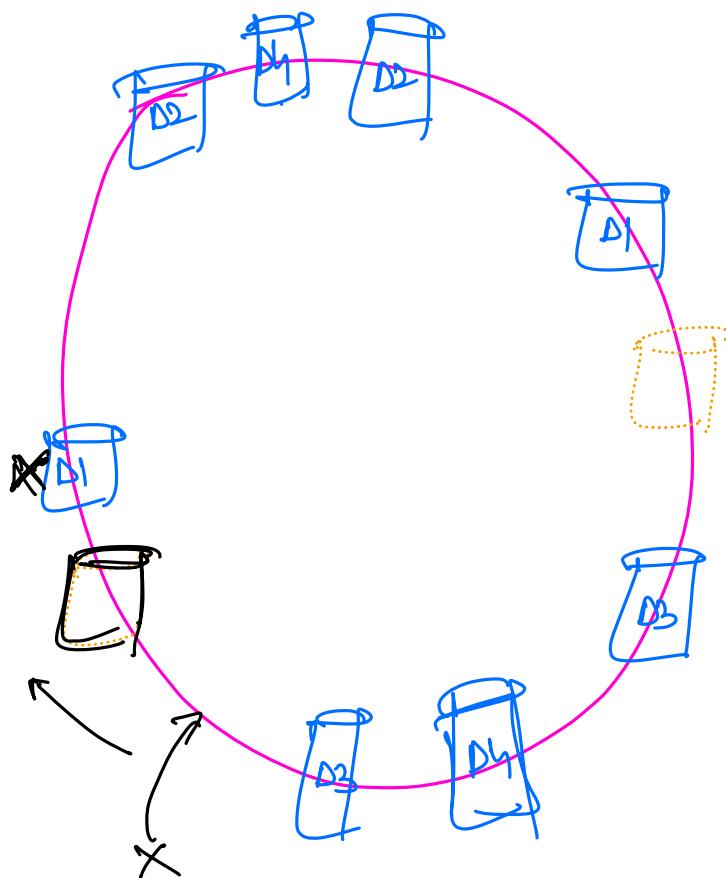




We consider a NoSQL DB as a blackbox where we only care about the  $f^m$  that DB offers.

How a DB Manager may work

- ① Create new Shards
- ② more replicas of each shard
- ③ CH work



A machine is not able to handle all of data:



[ ADD another shard ]

- ① They bring machine to chain right away.



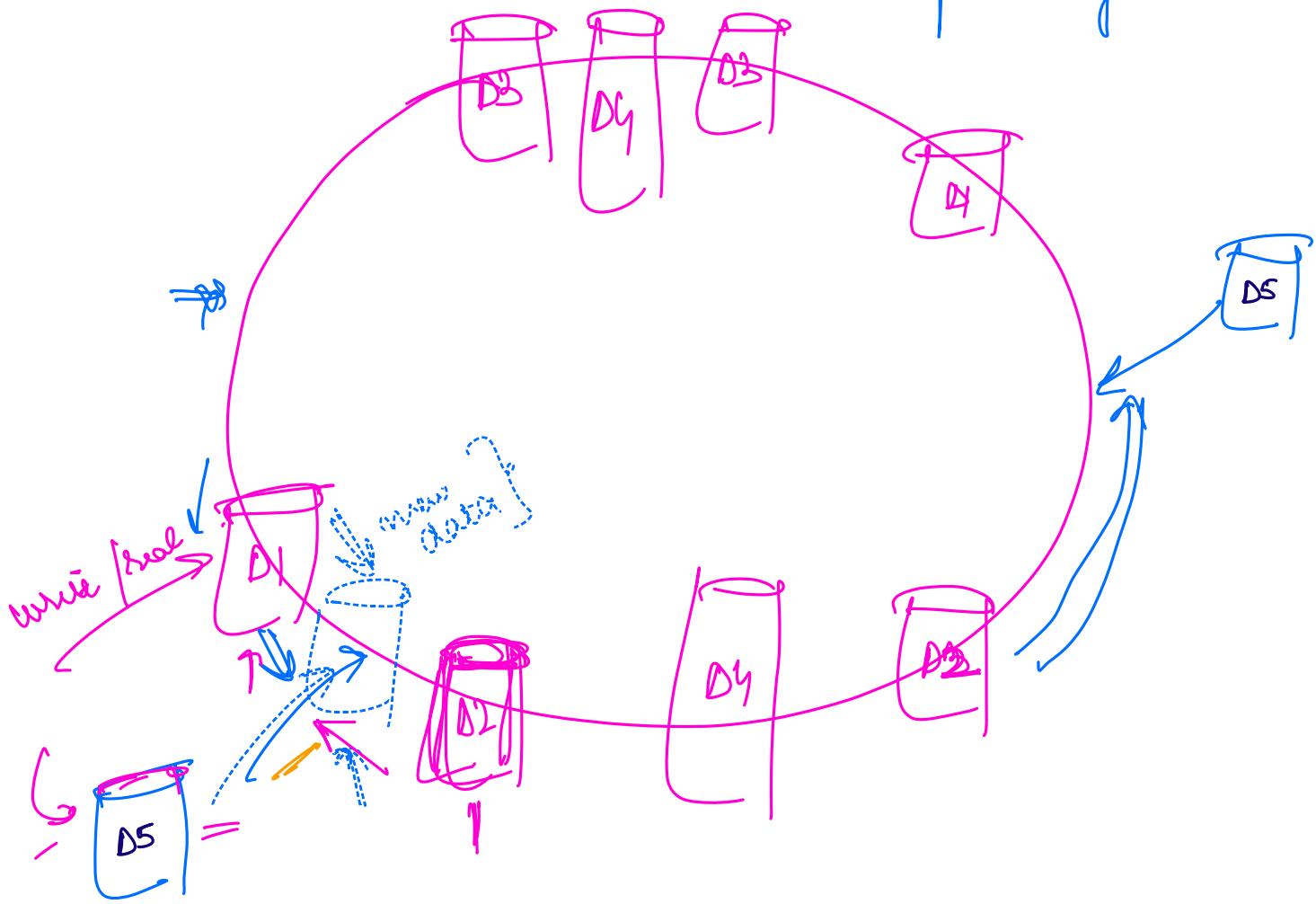
$\Rightarrow$  X doesn't exist

Available  $\Rightarrow$  ✓

Consistent  $\Rightarrow$  ✗

② first I fill the data into new shard  
and then bring to the chain.

Staging Time: when I am getting my new  
shard warmed up (ready)



Staging Time  $\Rightarrow$  Transferring data from older  
shards to new shard

2/3 of peak capacity

Write

a) on new one:

Read

$\Rightarrow$  Old shard keeps on  
handling the

what if a new comes.

The old one won't have

value.

(b) For old one

-  
queue

→ Once the new shard is

live, I will have to  
copy the newly created  
data.

↳ till not done, my  
reads can be  
inconsistent

(c) Write to both

↳ writes will become  
slow

---

When a new shard is going to come  
up, there will be a lot of work to  
be done, to get that shard ready.

→ at this time, my system may be slow.

- ① System is overloaded
- ② at that time, you are adding an additional overhead of sharding

↓

System may get impacted heavily on perf.

### DELETE A SHARD

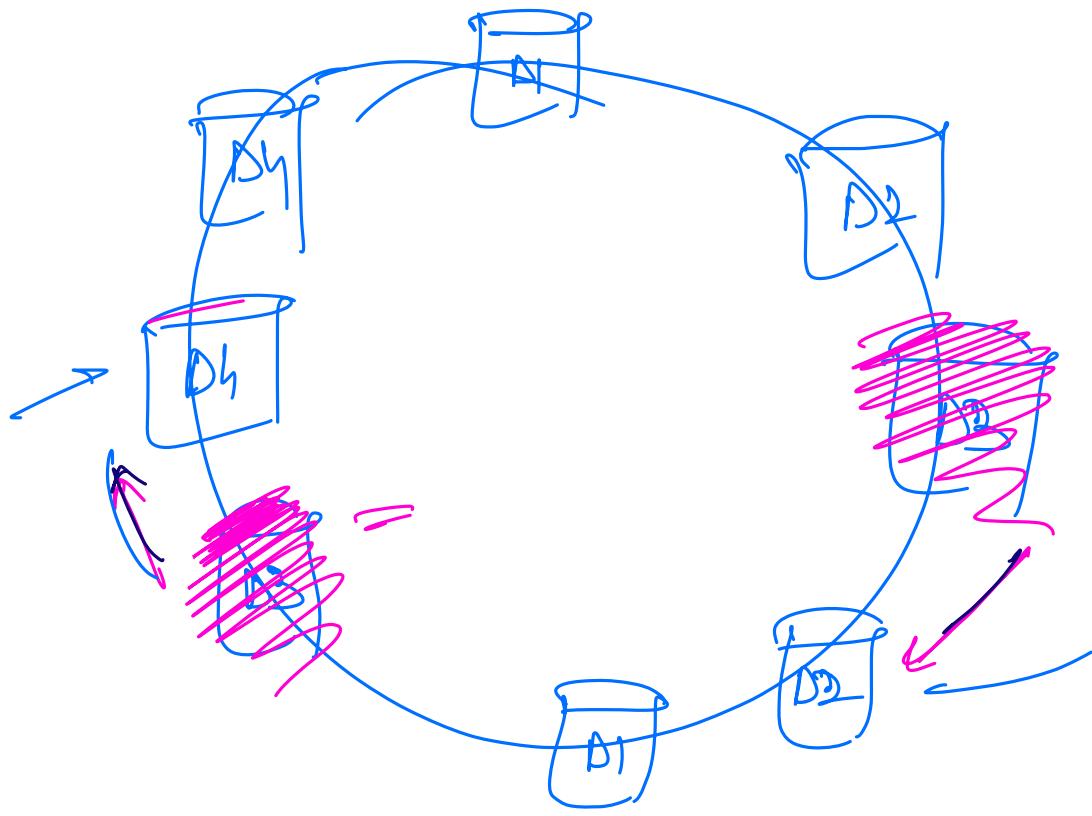
→ ① Load has reduced

→ ②

Master machine of the shard may itself die.

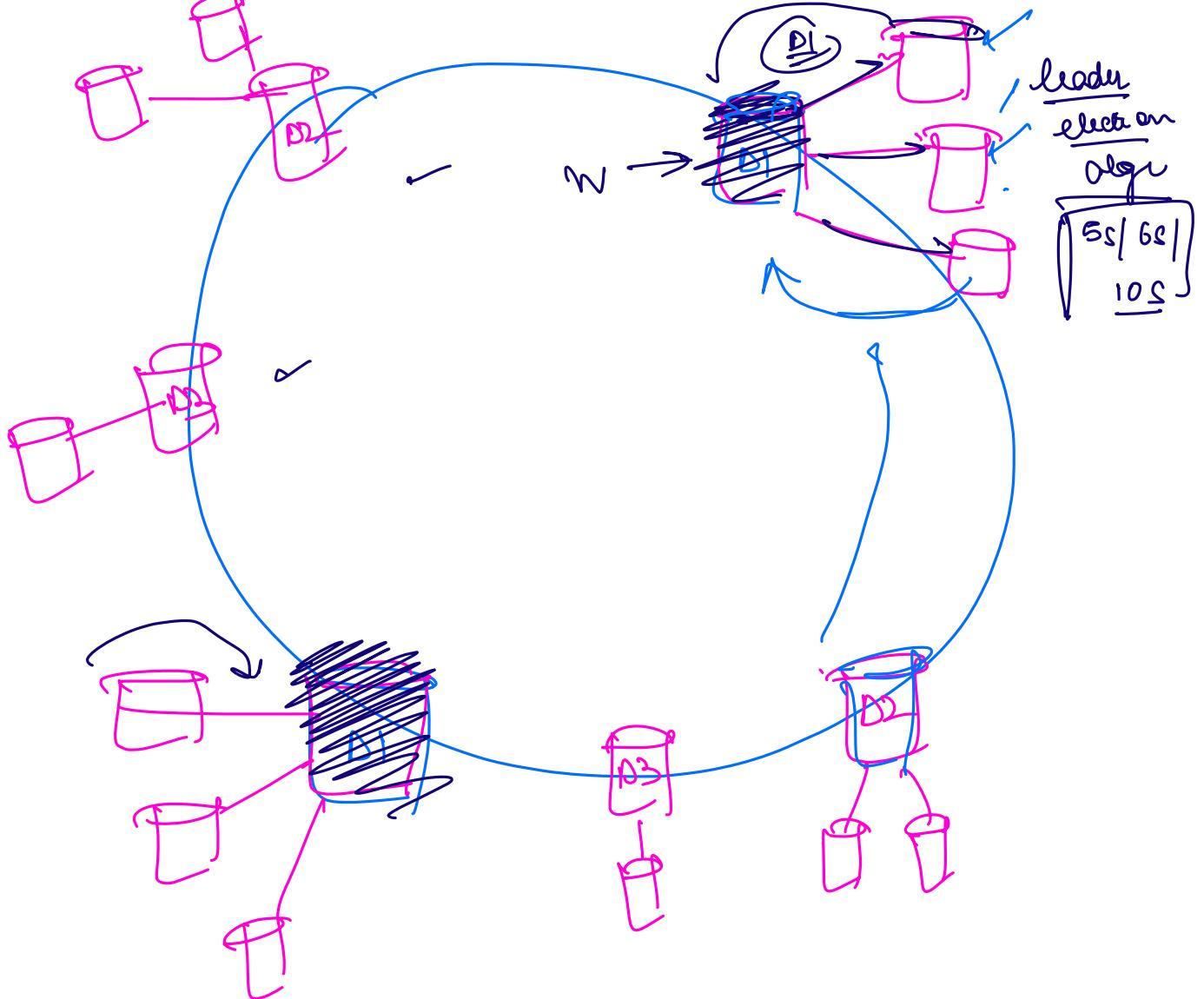
LOAD HAS REDUCED

~~not~~

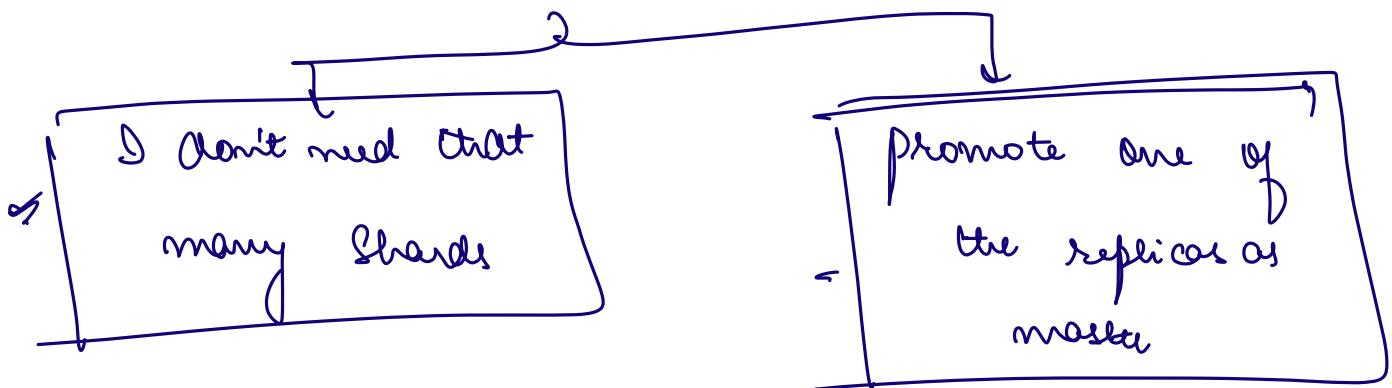


① We will be transferring the data from  
 ② D3 to D2 and D4.





⇒ Machine has died



Writer

- ① Write to all replicas

↳ writer will  
become slow.

↑s

- ② keep my write on  
hold

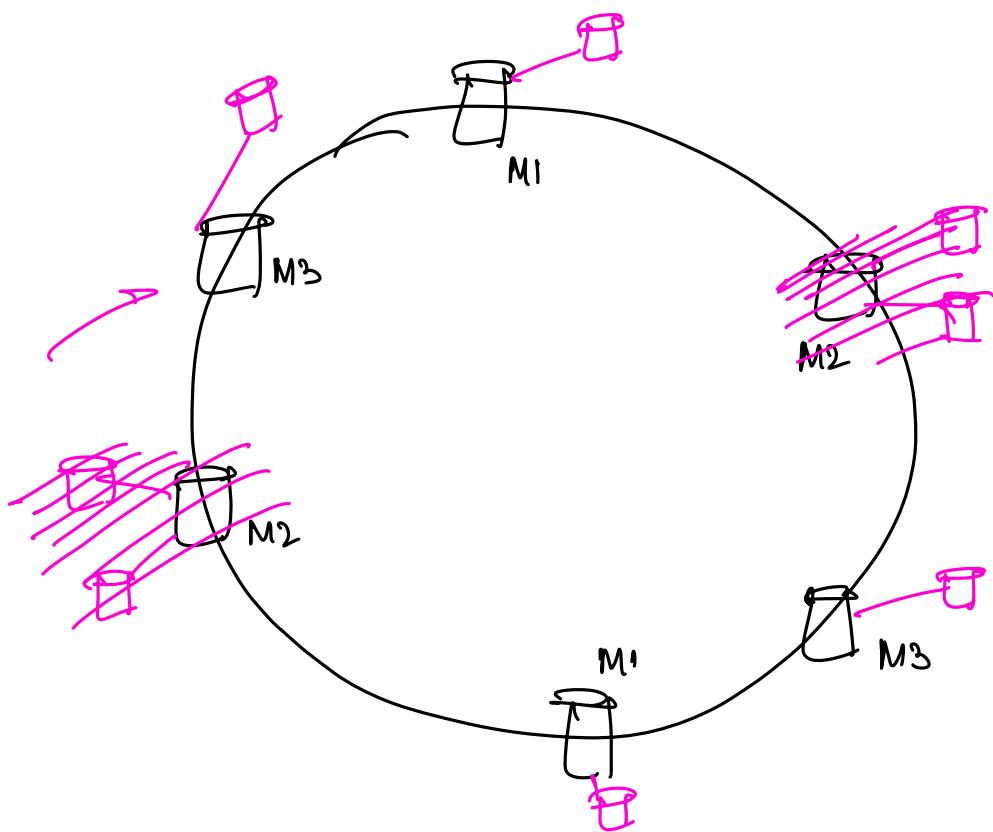
→ unavailable.

Read

→ can still happen on  
replicas

---

Is there a cleaner col<sup>n</sup>

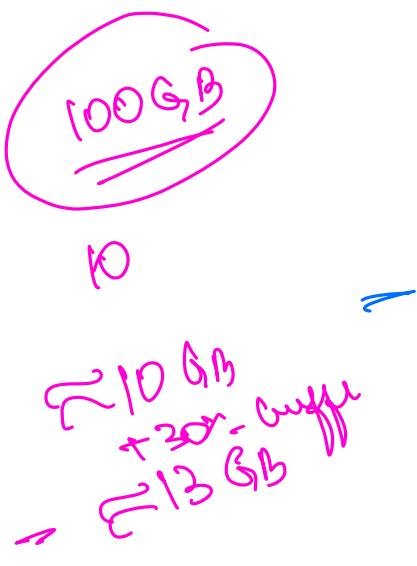


Till Now

| for every key  $\Rightarrow$  there was a dedicated shard

then why did I need replica machines?

{  $\Rightarrow$  To handle SPOF }  
 $\Rightarrow$  To handle more queries



$\Rightarrow$  If load inc, inc replicas.  
 $\Rightarrow$  if data size, inc # of shards.

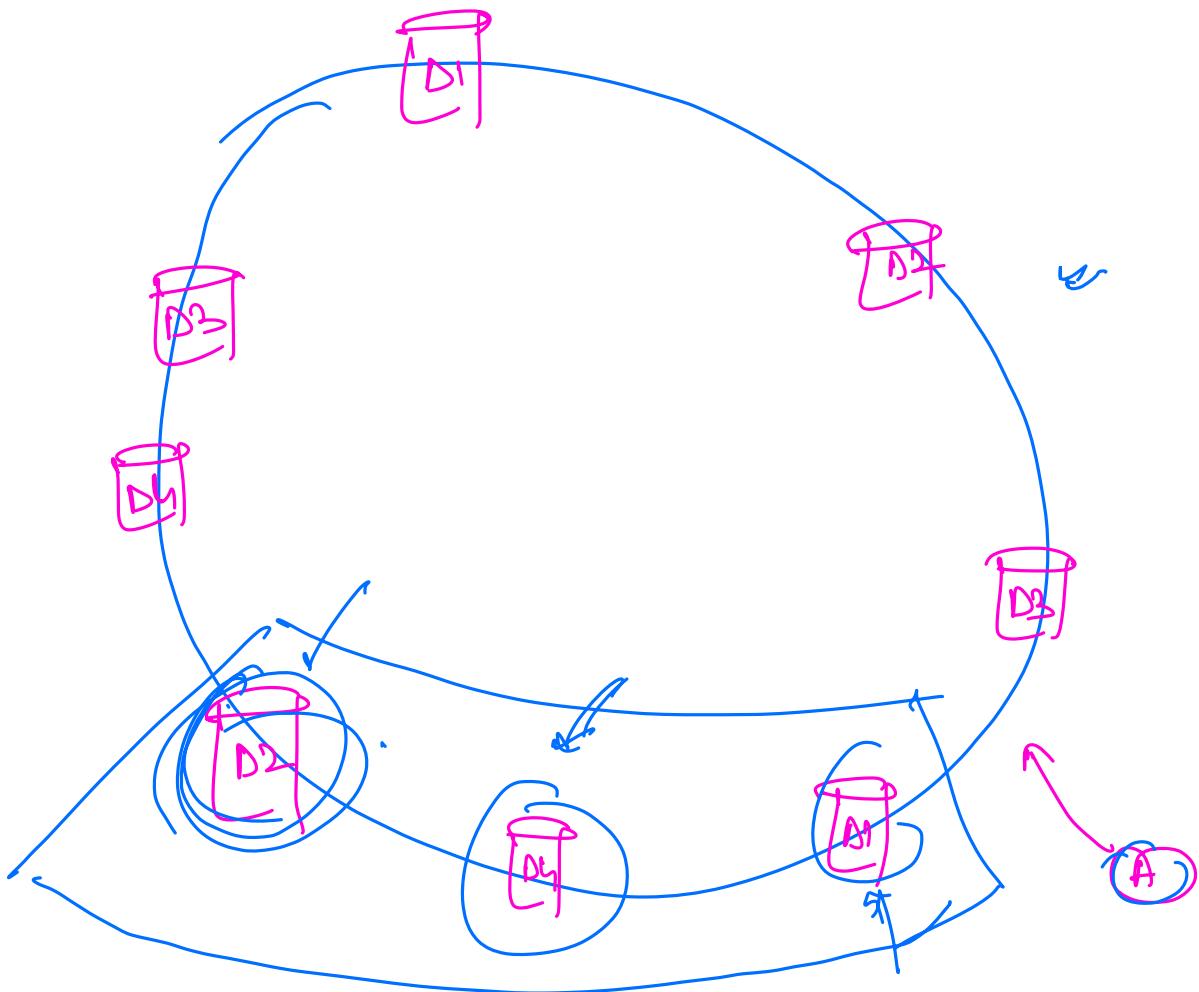
## Multi Master Architecture

$\rightarrow$  remove the need of replicas

$\rightarrow$  # of machines have reduce

$\rightarrow$  don't associate a key to only one shard-, instead associate that to  $x$  shards.

(assume  $x=3$ )

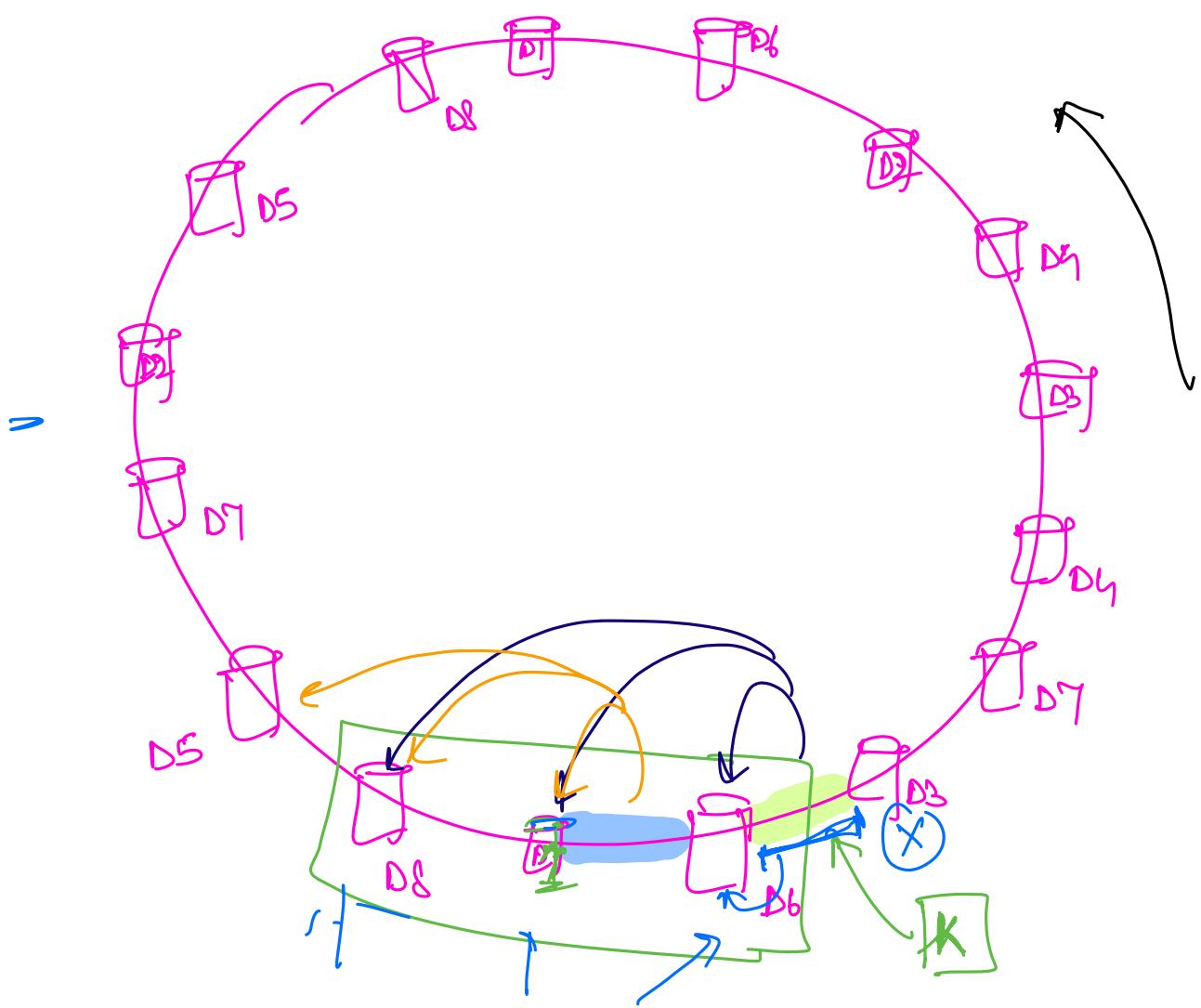


① Based on CH , Key  $\textcircled{A}$

↳ G will write its value to next K machines

→ Now instead of 1 owner, there are X owners

→ No need of replicas. Because anyway, data is being written at multiple places.



⇒ I have just made my slaves also as shards.

$$X = 3$$

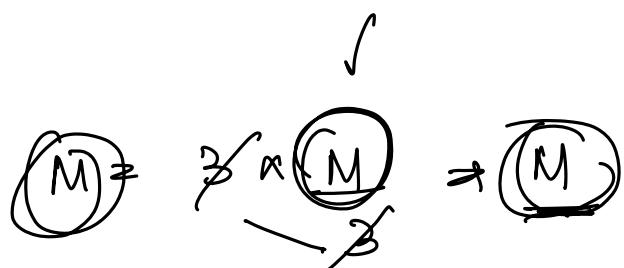
$K \Rightarrow \boxed{3 \text{ owners: } 6, 1, 8}$

Why we needed replicas earlier:

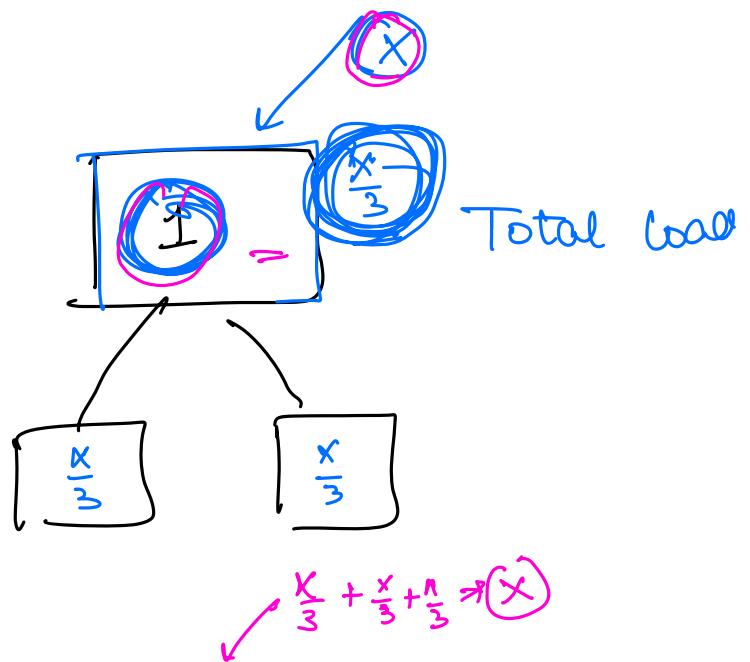
~~① Backup~~

~~② To share load of read queries.~~

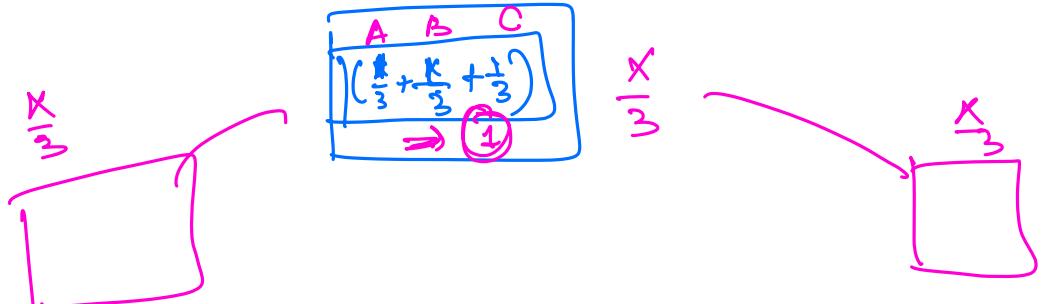
→ distribute amongst all masters that own the key.



earlier  $\Rightarrow$



now

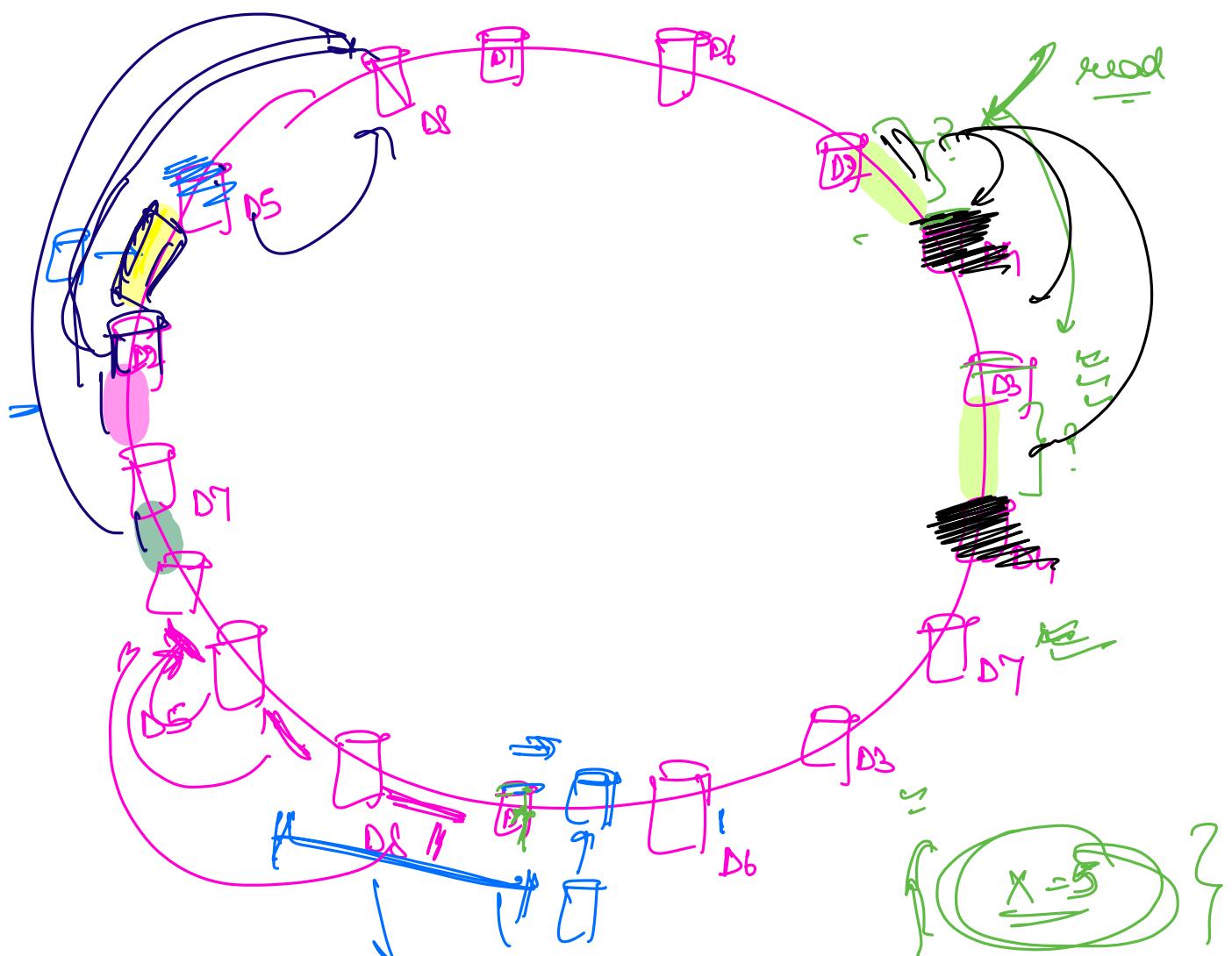


P

Problem with M-S  $\Rightarrow$  issues at time of adding/  
removing a chord.

for each machine, there will be  $\times$  owners

## DELETION OF A CHARD



for the keys:

Read

- no need of leader elec
- just go to D3

Write

just write to next

X ma chini

adding a new shard

Multi Master Architecture

⇒ Algo : Cassandra



Each Key has  
X masters

Master - Slave Architecture

each key has 1 master and  
each master may have slaves

In M-S

① Consistent

Write to master as well as all replicas

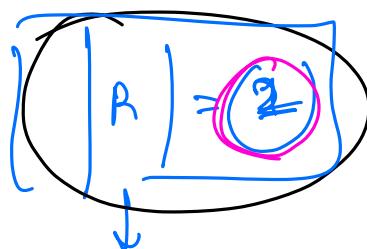
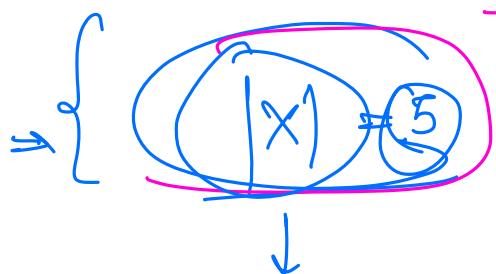
② Available

Write to master and more on.

In M-M

↳ tunable consistency

↳ how much consistent / how much available



$$|W| = 1$$

↳ how many masters for every key

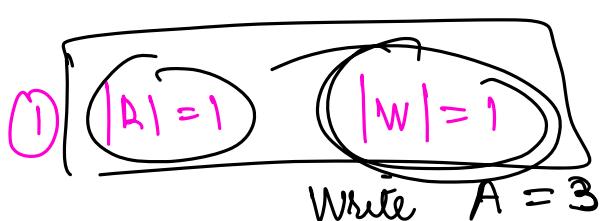
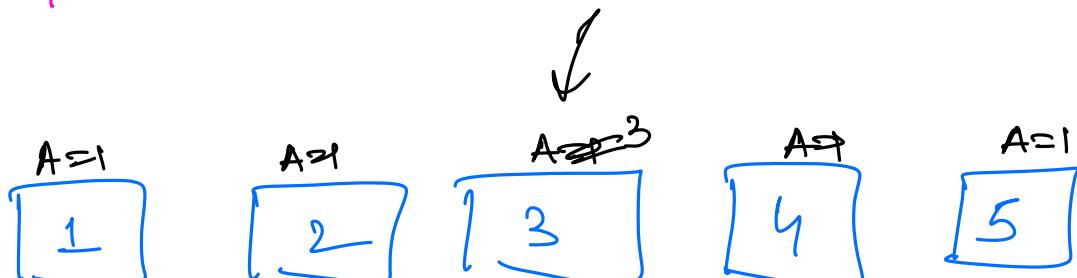
↳ how many machines to read from before returning to client

↳ if any 1 of master have same value return

↳ how many machines to write to before returning success to client

↳ later asynch with sync  
other

$$K = |K| = 5$$



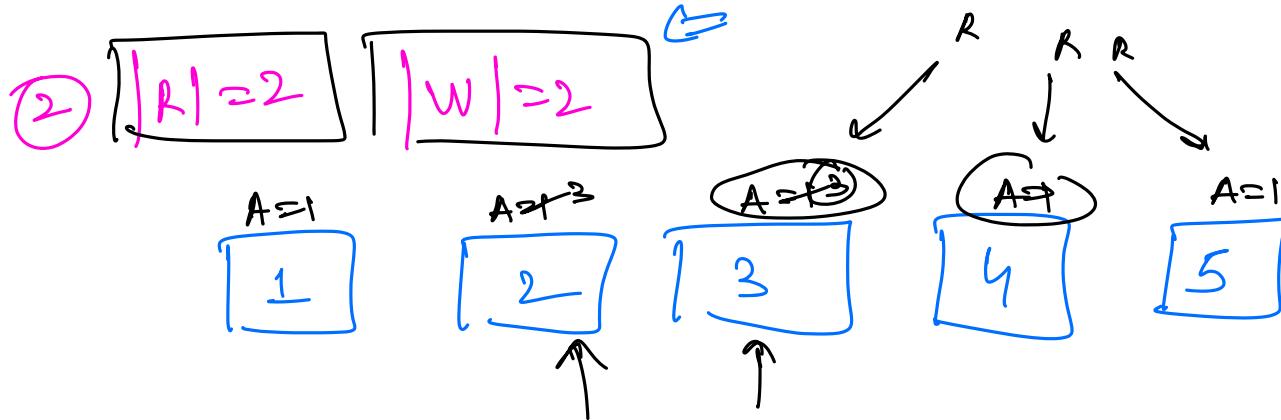
lower value of  $W \Rightarrow$  fast writes

Read A  $\longrightarrow$  1

Available  $\Rightarrow$  ✓

Consistent  $\Rightarrow$  ✓

less value of  $R \Rightarrow$  fast reads.



Write  $A=3 \Rightarrow$  higher latency than  $|W|=1$

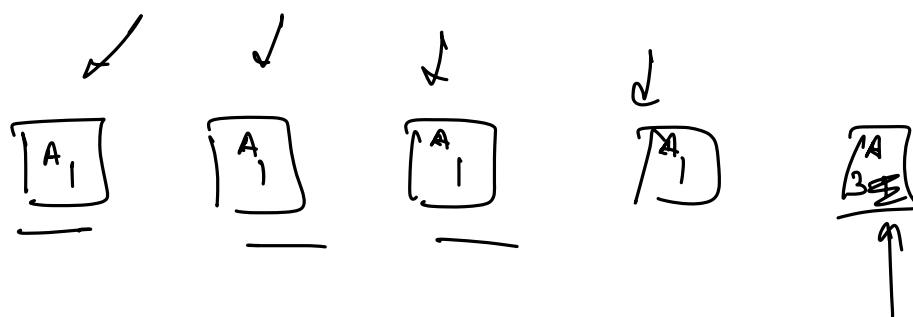
Read A  $\Rightarrow 3$

consistent but higher latency

1

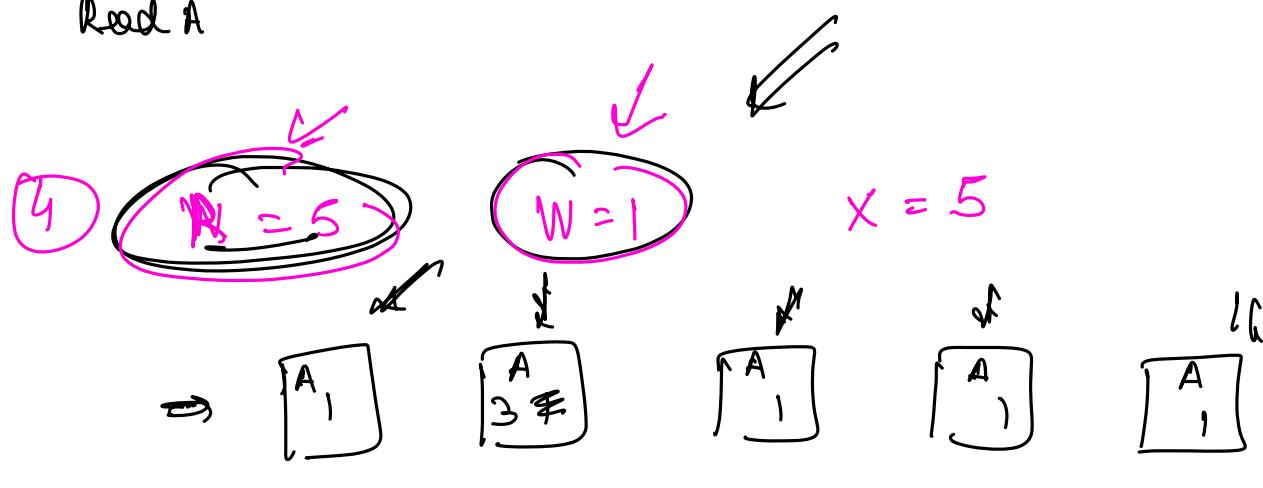
$\rightarrow$  available but higher latency

③  $|A| = 4$   $|W| = 1$   $|X| = 5$



Write A=3

Read A



Write  $A = 3$

⇒ Consistent  $\Rightarrow$

✓

Read A

Available  $\Rightarrow$

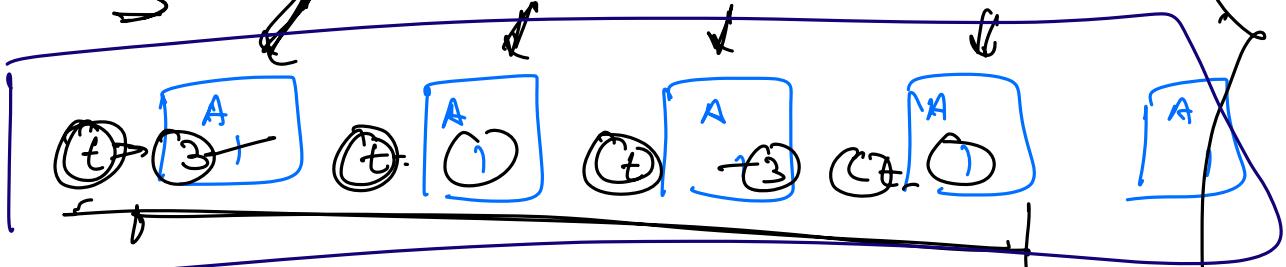
X

Consistent  $\Rightarrow |R| + |W| > X$

$|W| = 2$

$|R| = 4$

$|X| = 5$



Write  $A = 3$

Read(A)

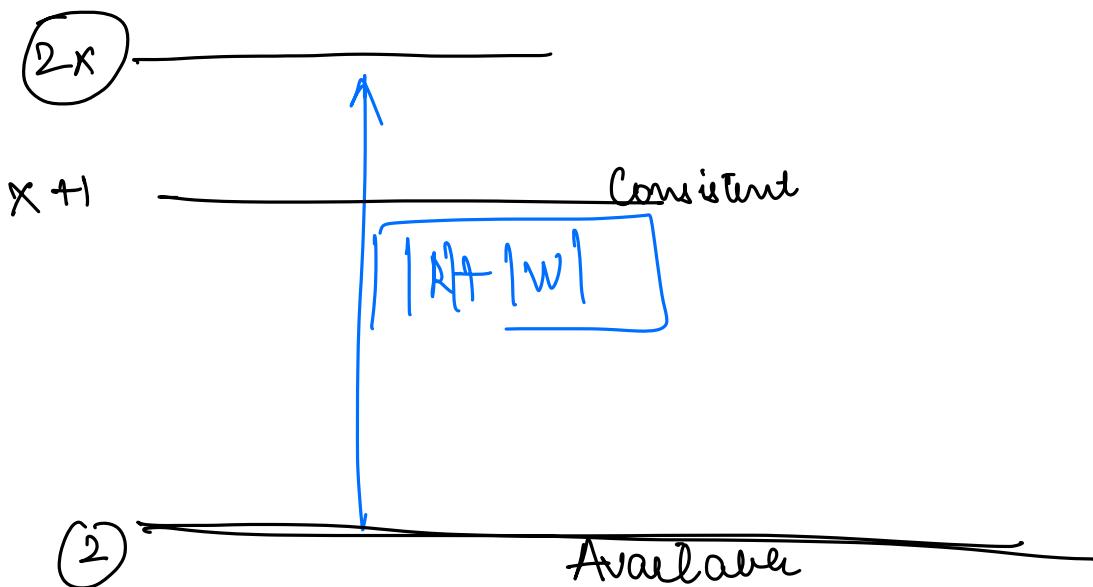
Consistent  $\Rightarrow |R| + |W| > |X|$

Fast writes  $\Rightarrow$  less W, big R  
Fast reads  $\Rightarrow$  less R, big W.

Available

$$|W| = 1$$

$$|R| = 1$$



Cassandra

Fast Write

Consistent

Slow Reads

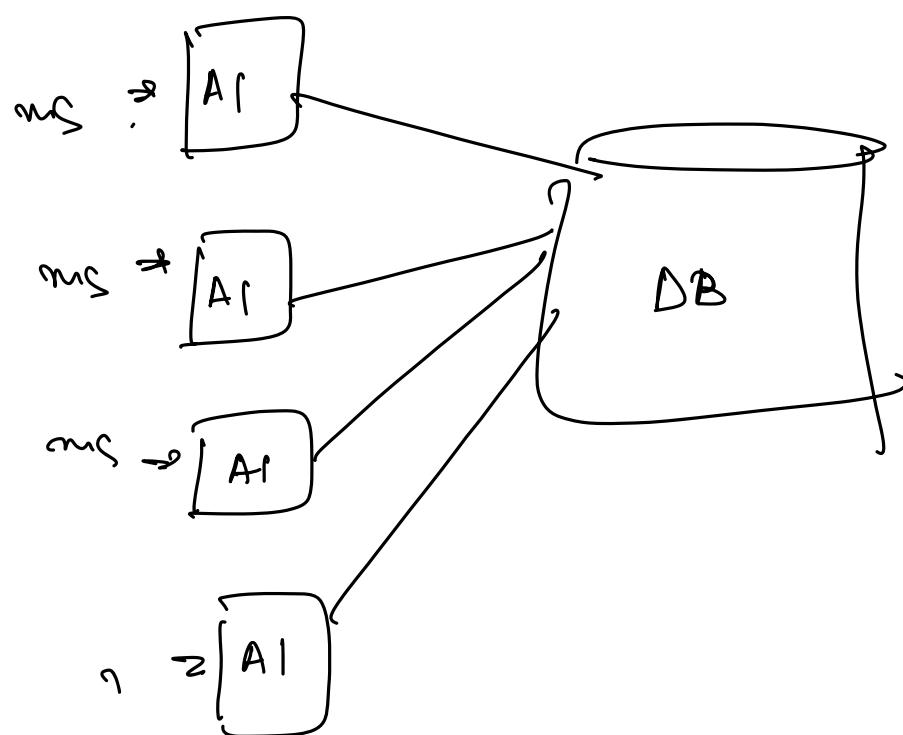
$$|W|=5 \quad |A| \geq 1$$

⇒ consistent

⇒ writer slow

⇒ reads v.v. fast.

"NO GLOBAL CLOCK"

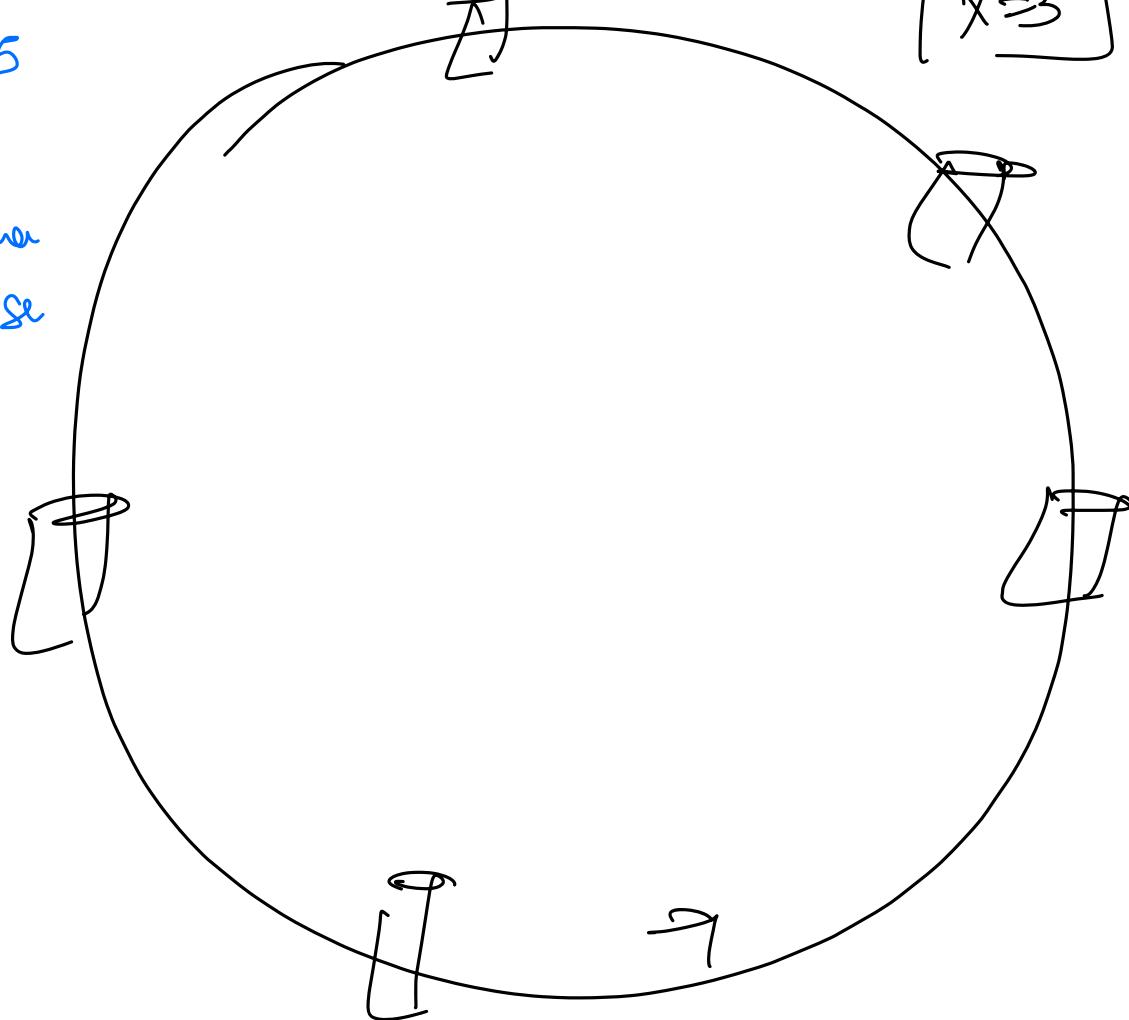


$x = 5$

6 gbar  
10 sec

[ ]

[  $x = 3$  ]



$$\checkmark \quad R \subseteq X \quad \checkmark \quad W \subseteq X$$

$$R + W \subseteq X$$

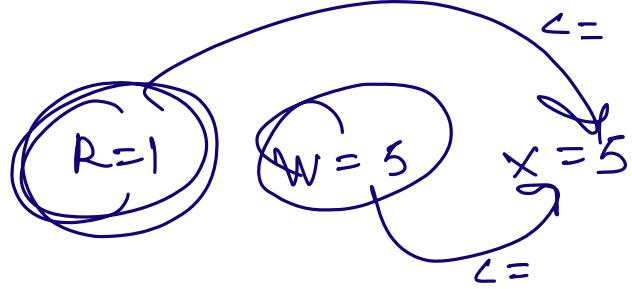
Avail ~  
Consume X

$$R + W > X$$

Consume  
Avail X

$$R \subseteq X$$

$$W \subseteq X$$



Coniston

W=1      R=5      X=5

Conis 1