

## Lec 19 / Database Scaling patterns

Aim: step by step manner, when to choose which scaling option.  
which scaling option is feasible practically at the moment.

\* Case study.

- Cab booking App

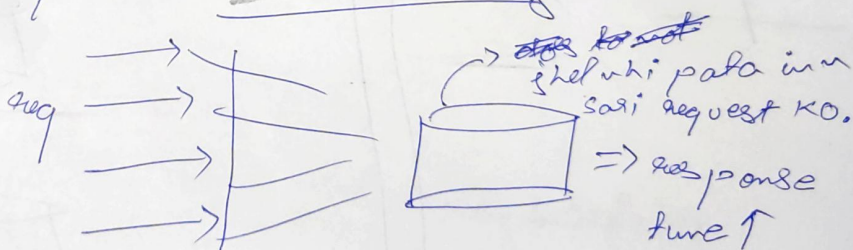
- Tiny start up
- ~ 10 customers onboard.
- a single & small machine DB stores all customers, trips, locations booking data and customer trips history.
- ~ 1 trips booking in 5 min

\* suppose app became popular

& 10 booking / min later

then feedback aaya ki app slow chl rha h

pta kiya th problem API Latency nikla



\* Your app becoming famous, but... the problem begins.

- Req. scales upto 30 booking per minute.
- Your tiny DB system has started performing poorly.
- API latency has increased a lot.
- Transactions facing Deadlock, starvation and frequent failure.
- Sluggish App experience.
- Customer diss satisfaction.

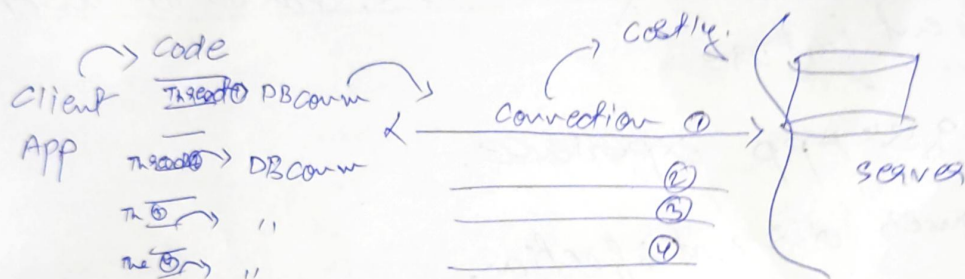
\* Is there any solution?

- We have to apply some kind of performance optimisation measures. we might have to scale our system going forward.



## \* Pattern 1 | A very optimisation & connection pool implementation

- Cache frequently used non-dynamic data like booking history, payment history, user profile etc.
- Introduce DB redundancy (or may be use NoSQL)  
=> less response time
- use connection pool libraries to cache DB connection



Is bhi koi connection create karta h ye costly hota h in terms of time complexity. So, it better ki, ~~code wale~~ <sup>code wale</sup> is bhi DB conn. ki call ko th baar baar ~~new~~ <sup>new</sup> connection establish na karo. So it's better ki ek connection ka pool bana ke rakho connections 1, 2, 3, 4

or iski bhi connection h ye inki threads ko baar baar reuse krke

=> Is bhi ~~performance~~ <sup>performance</sup> ko issue aa rha h connection create krne wale wo gyab ho jayega.

Is bhi library ko use kr rha h, un sab library wale connection pool ki libraries ~~are~~ <sup>are</sup> h unko use kr ke cache DB connection.

- multiple application threads can use same DB connection.

& Good optimisations of now

- Scaled the business to one more city, and now getting ~100 bookings/minute

And then phir se wohi problems wapas aagye

10 bookings/min increase hua tho.

Now,

\* Pattern (2) vertical scaling & scale-up.

- upgrading our initial tiny machine.
- RAM by 2x and SSD by 3x etc.
- scale up is pocket friendly till a point only.
- more you scale up, cost increase exponentially.

\* Good optimisation as of now.

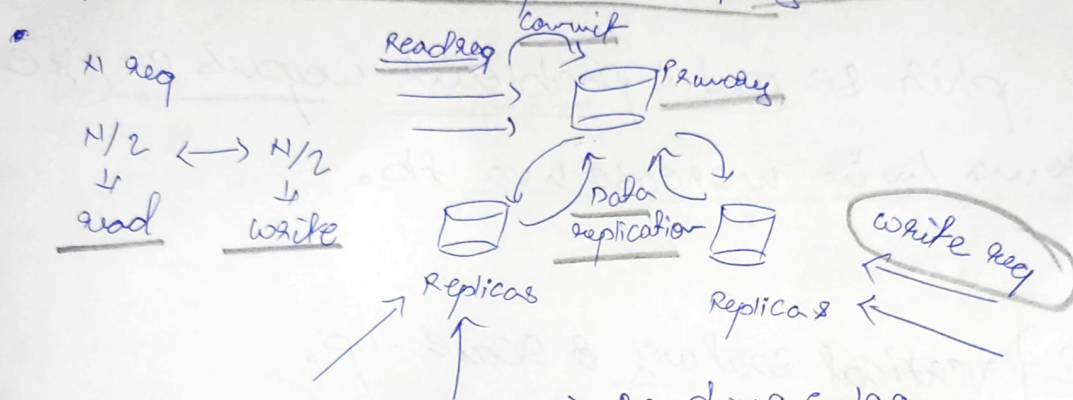
- & then chlo a growth hua tho phir
- Business is growing, you decided to scale it to 3

more cities and now getting 300 bookings per minute.  
 & then phir se problems aagye lekin ab scale up  
thi nhi kr sktte kyunki us waqt th but paisa  
jata h ~~th~~ th ab ek better option dhkna  
hoga.



## \* Pattern 3 / Command Query Responsibility Segregation / CQRS

- The scaled up big machine is not able to handle all read/write requests.
- Separate read/write operations physical machine wise



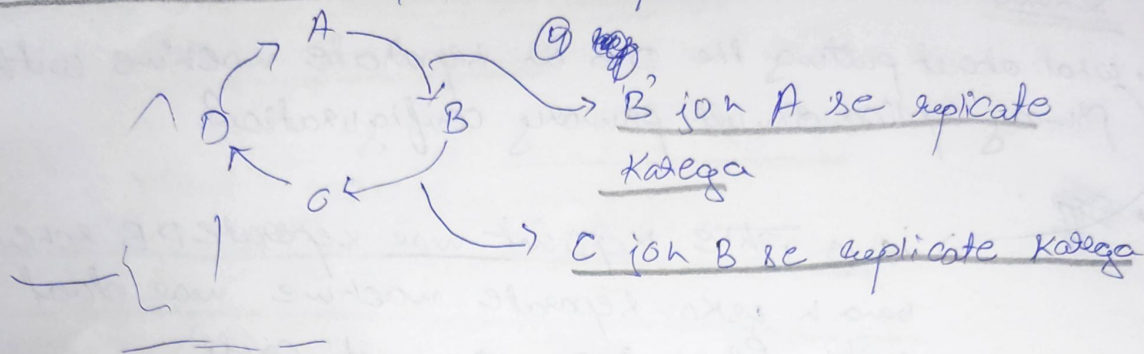
- All read queries to replicas.
- All write queries to primary
- Business is growing, you decided to scale it to 2 more cities.
- Primary is not able to handle all write requests.
- Log b/w primary and replica is impacting user experience

as app is scale katch in Pattern 3 bhi  
hoga tb.  
next Pattern



## \* Pattern ④ / Multi primary replication

- why not distribute write req. to replica also?
- All machines can work as primary & replicas.



Ab agar koi bhi write req aai th ye kisi bhi random node ko de denge. (any node will handle it).

i.e. agar 1 req. aai th 1 req distribute hojenge, th for that time pe replication chhi nhi

or read karne ke liye read ko broadcast karenge

i.e. chao node wase se phire reply kiya usko chli jaegi

- Multi primary configuration is a logical circular ring.
- write to any node.

~~• write data to any node.~~

- Read data from any node that replies to the broadcast fast

• How scale to 5 websites & your system is in pain again.  
~ 50 req/s

then

next Pattern



## \* Pattern 5) Partitioning of Data Functionality

• What about separating the location tables in separate DB schema?

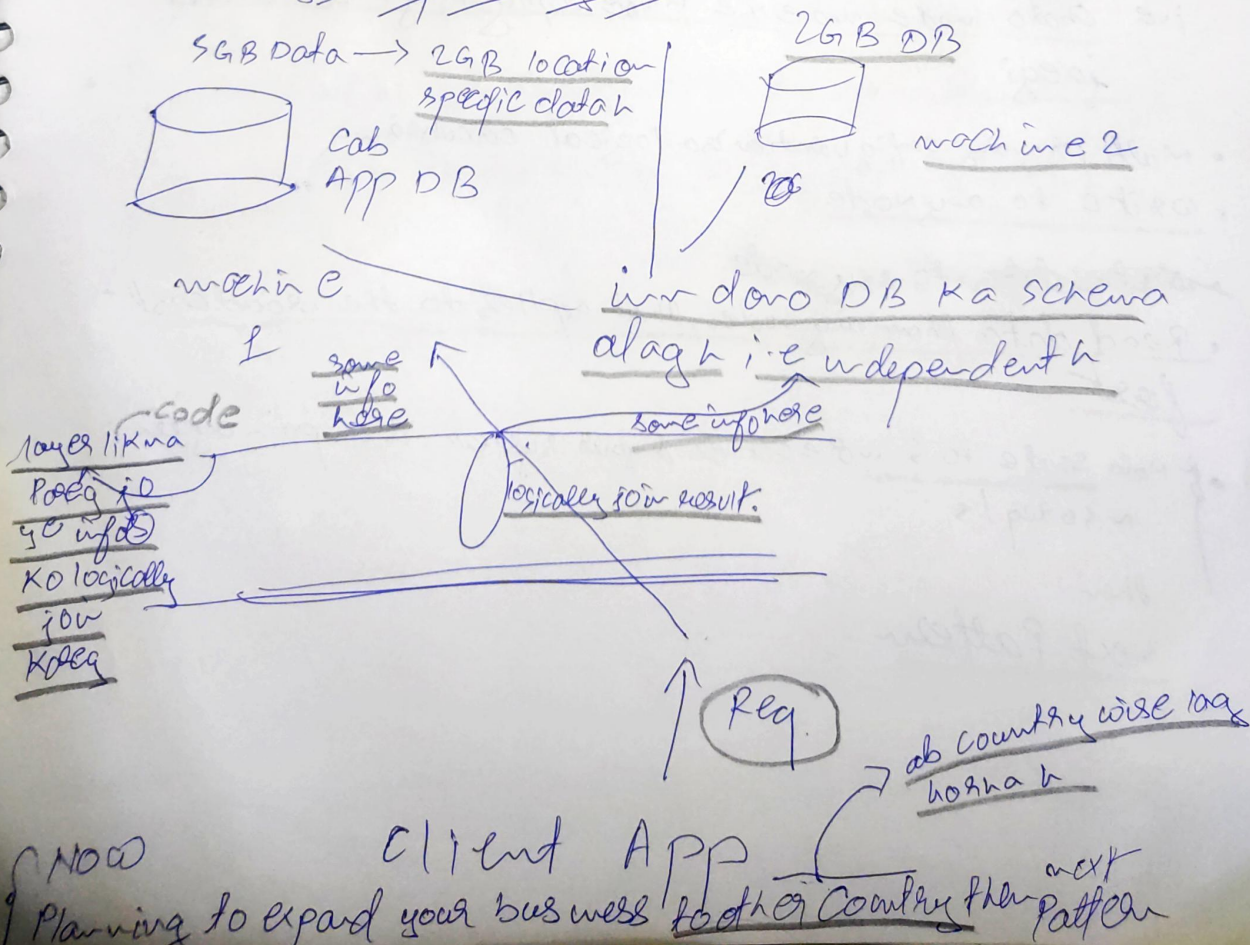
• What about putting the DB in separate machine with primary replica or multi primary configuration?

• ~~Diff~~ → ie jo phle report me separate DB schema bana kr usko separate machine me dal do. with 1m rep or multi confi

• Diff DB can host data categorised by different functionality.

• Backend & application layer has to take responsibility to join the results.

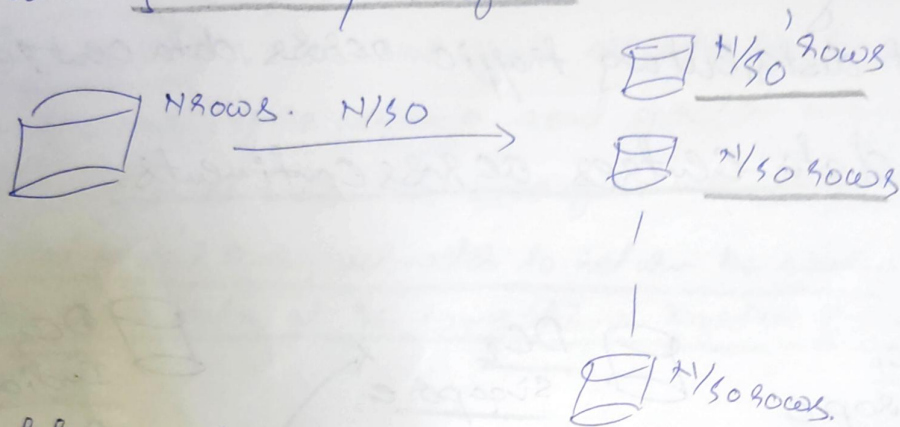
• Planning to expand your bus



## \* Pattern 6) Horizontal Scaling or Scale out

- sharding - multiple shards.

Allocate 50 machines - all having same DB schema - each machine just hold a part of data



- locality of data should be there.

In each machine we have local data and change takes machine we switch no karna pade.

- Each machine can have their own replica, maybe used in failure recovery.

- sharding is generally hard to apply { No pain  
No Gain }

Now  
scaling the business across continents.

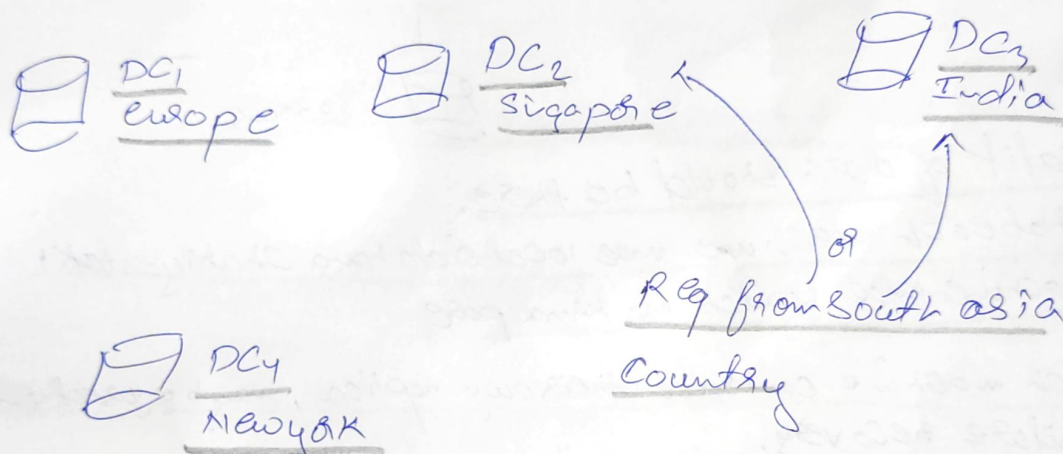
then again java server n wo thik n lek in behar log probs no smeth ab next pattern.



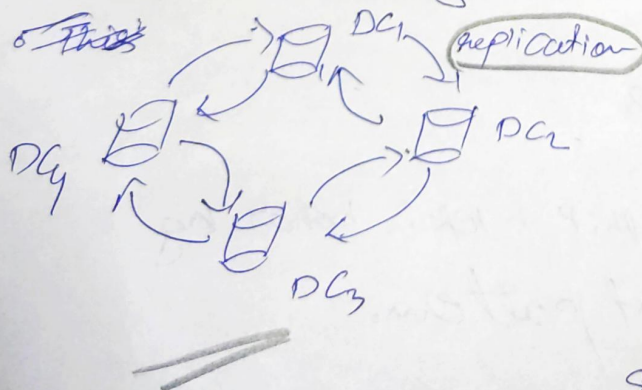
## \* Pattern 7) Data centre wise partition

- Requests travelling across continents are having high latency.
- What about distributing traffic across data centres?

=> Create data centres across continents.



- Enables cross data centre replication which helps disaster recovery.



Yani ki kuch kuch time wale  
min DC me replicas wale replication  
Rate change data ki

=> Agar Europe ka DC nhi  
chalta for some reason  
the public of Europe's req wo

Singapore or India bhej skta h  
yes, latency will come but  
availability will be high

But BB DC Europe sahi to page than again normal