

* Scale from 0 to 1M users *

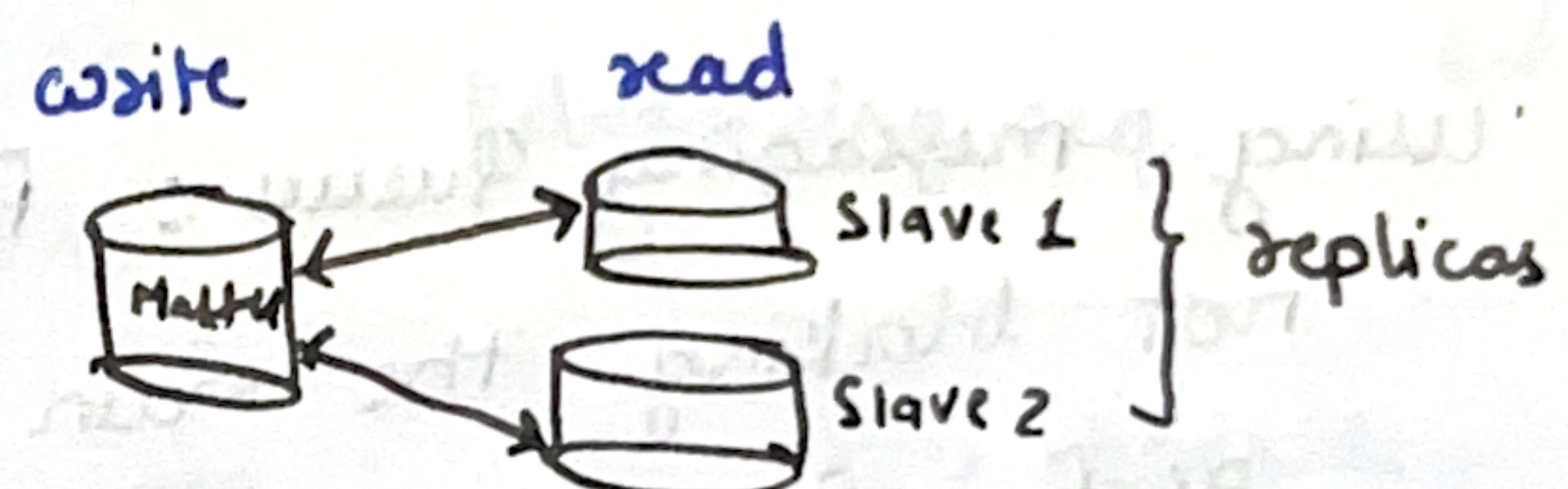
application servers

① ⇒ separate out mid tier & DB to let them scale independently

② ⇒ Bring in LB and multiple App Servers
↳ also brings security as LB & app servers talk on private IP.

if an app server fails the redistribute traffic among others

③ ⇒ Database Replication

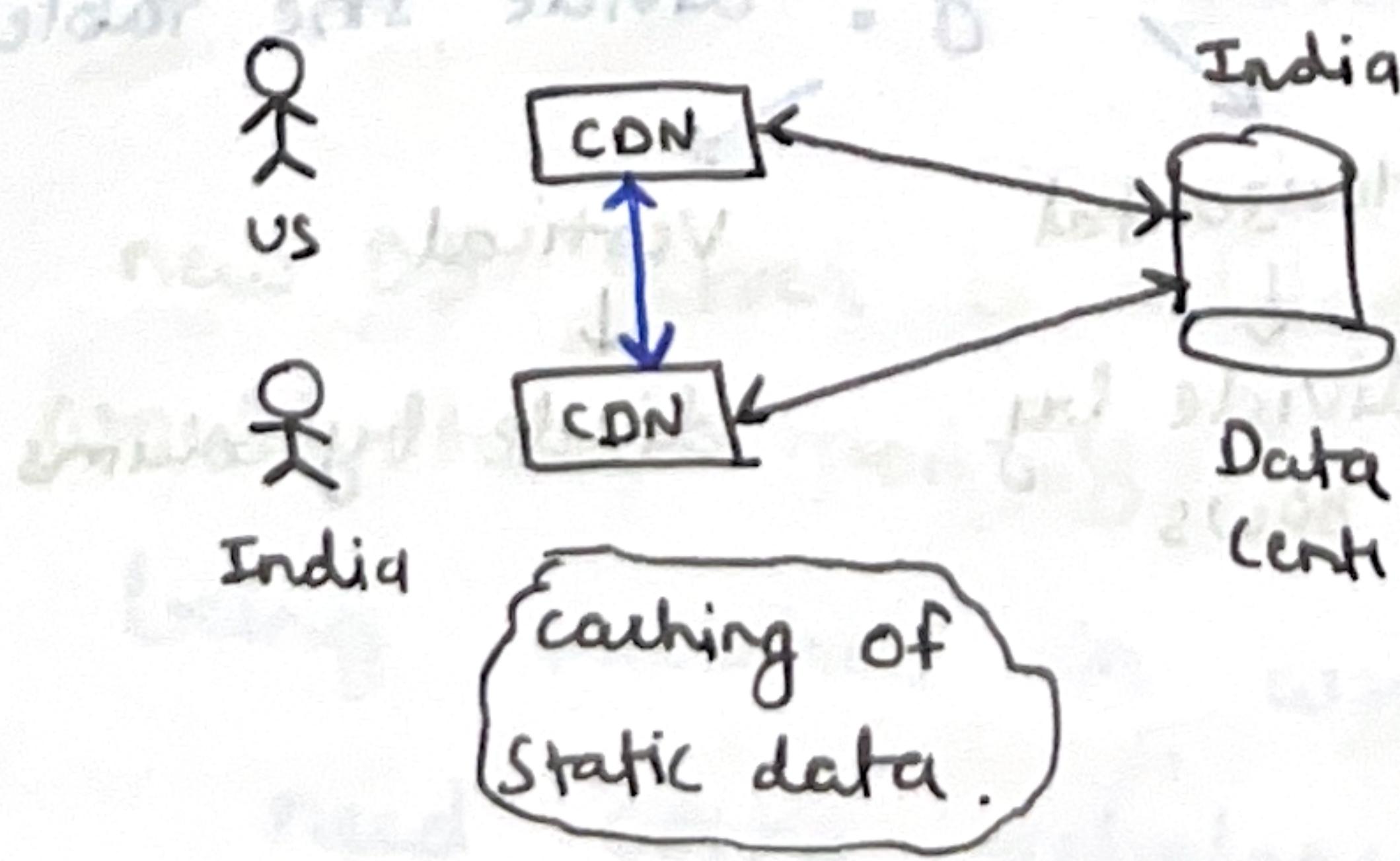


⇒ if master DB fails then one of the slave DB becomes master.

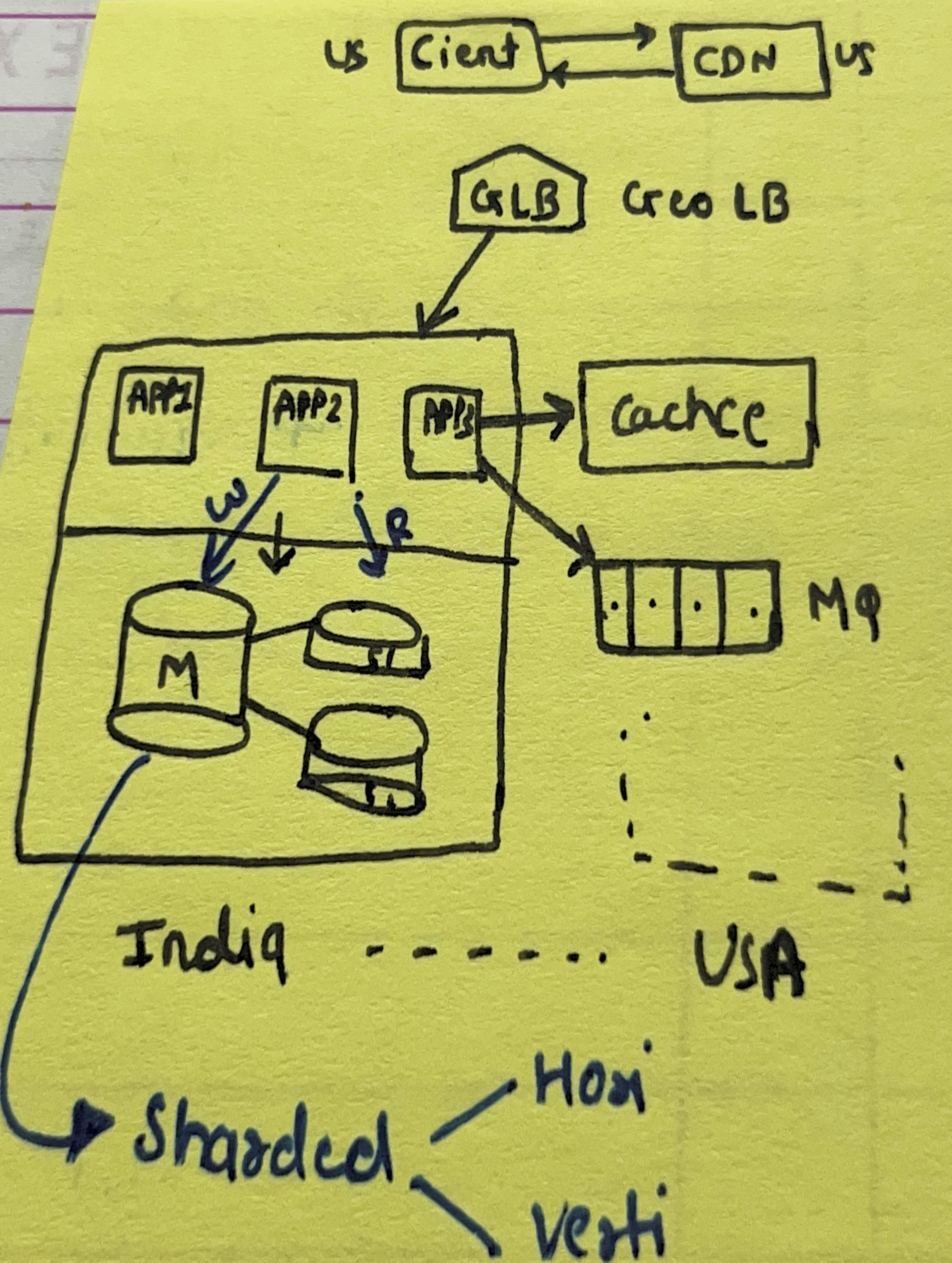
④ ⇒ APP ↔ DB are expensive because it's a network call,

to overcome this we introduce **cache** → TTL (Time to live) ⇒ this will make data from cache to be false.

⑤ ⇒ CDN: content delivery network. (cache)



CDN closer to user reduces latency for static content, if cache miss in CDN then go to nearest other CDN if miss then go to original DB.



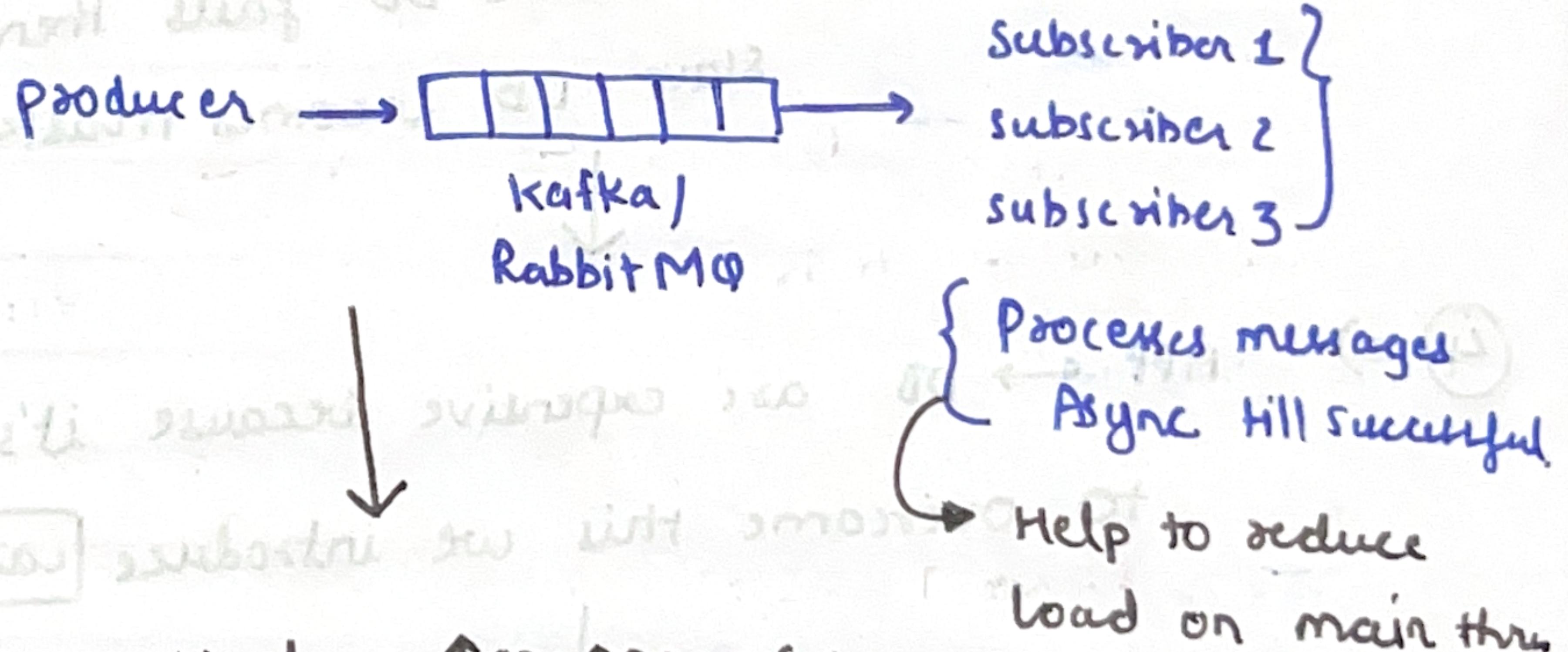
⑥ multiple datacenters based on users location

↳ same replica of application & database

⇒ here we have Geo Load Balancing, based upon latency to reach server.

* there is a replication bet' zones & regions to support Disaster recovery.

⑦ Using messaging Queues : for async exec of work & not blocking the main behaviour, this can improve perf of main service



⑧ DB scaling :

Vertical : ↑CPU↑RAM (limited)

Horizontal : add more nodes

↳ Sharding : divide the tables

Horizontal

divide by rows

Vertical

divide by columns

Drawbacks of Sharding:

- again and again de-sharding might be required.
- join can cause latency problems, so prefer denormalized table to avoid join.
- can be solved by consistent Hashing.