

lec 02

Anything that we will do in world is to make our application  
Scale Better

## Simple Client-Server Architecture

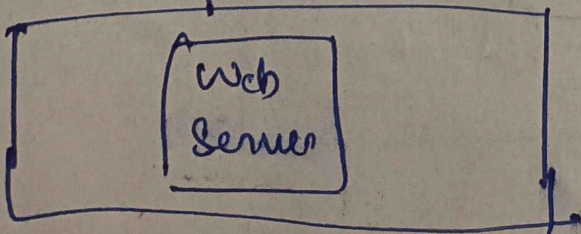
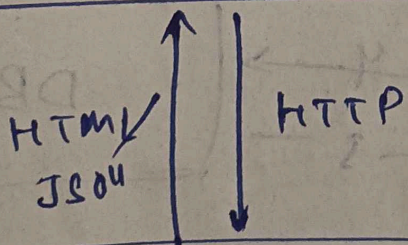
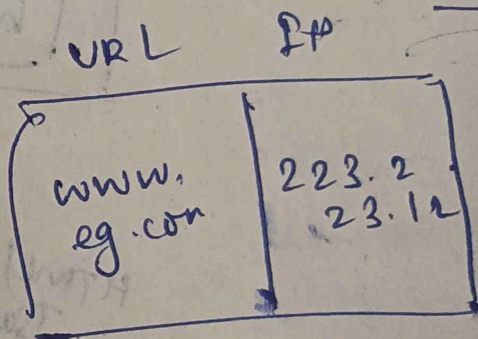
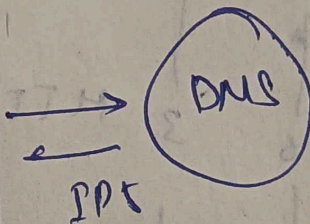
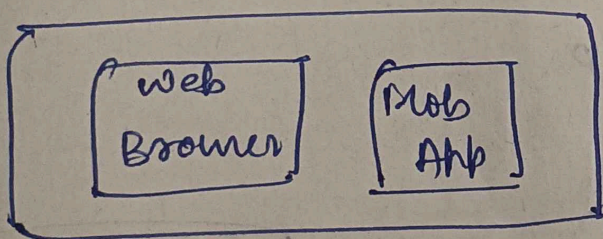
CLIENT

SERVER

Tabhi Client, Server ko  
konse request bhejta h  
to us directly hi  
bhejta, us  
"DNS" ke through  
bhejta h.

IP address  
lke dega  
Server ka

Architecture eg.





HTTP mai kya kya hota h?

C. POST, GET, PUT etc.

URL: www.eg.com

Header:

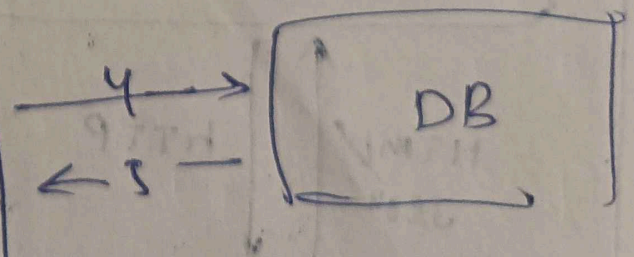
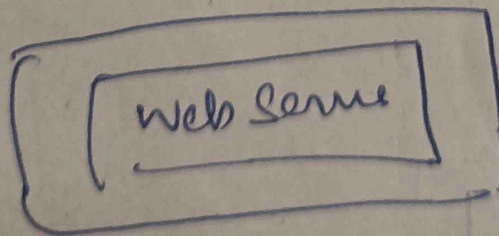
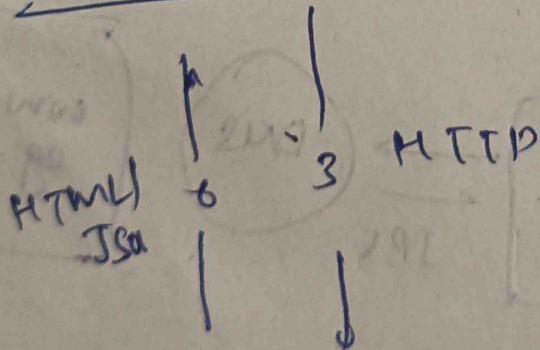
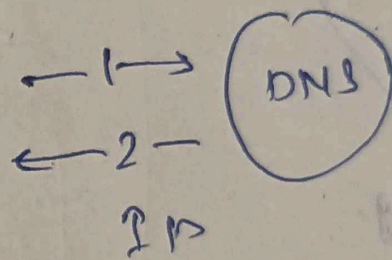
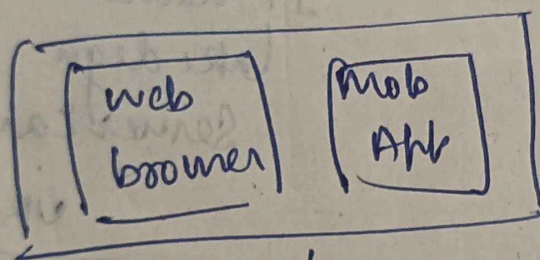
Body: {

"userID": 1,

"Name": "Aditya";

}

Introducing DB Layer





Which DB to me?

Relational

SQL Based

mysql,  
Postgresql,  
oracle sql

Non-Relational

Key-value Store

eg:- Amazon  
DynamoDB

Column Based Store

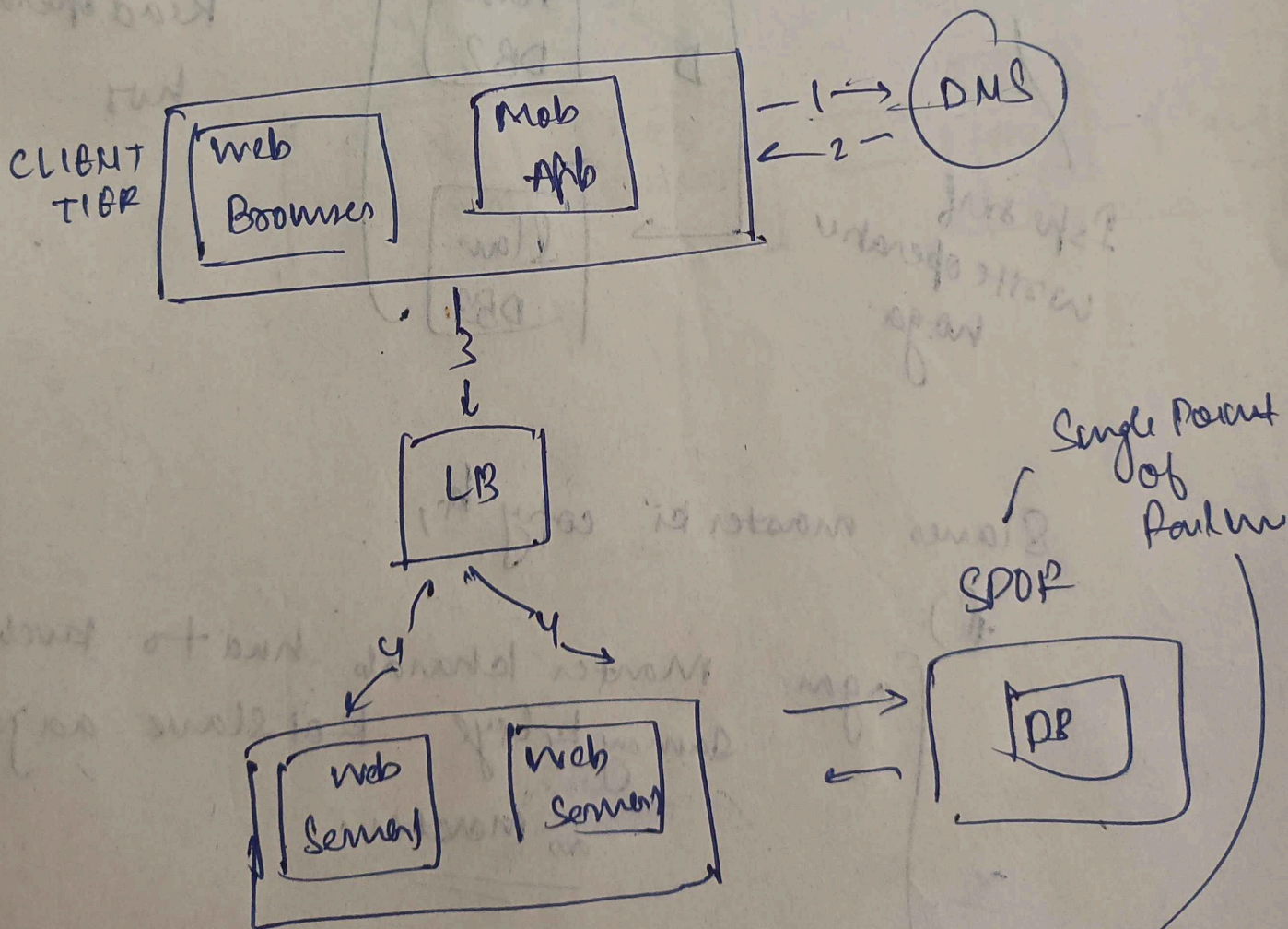
eg:- Cassandra DB

Document "

eg:- mongo db

Graph Store

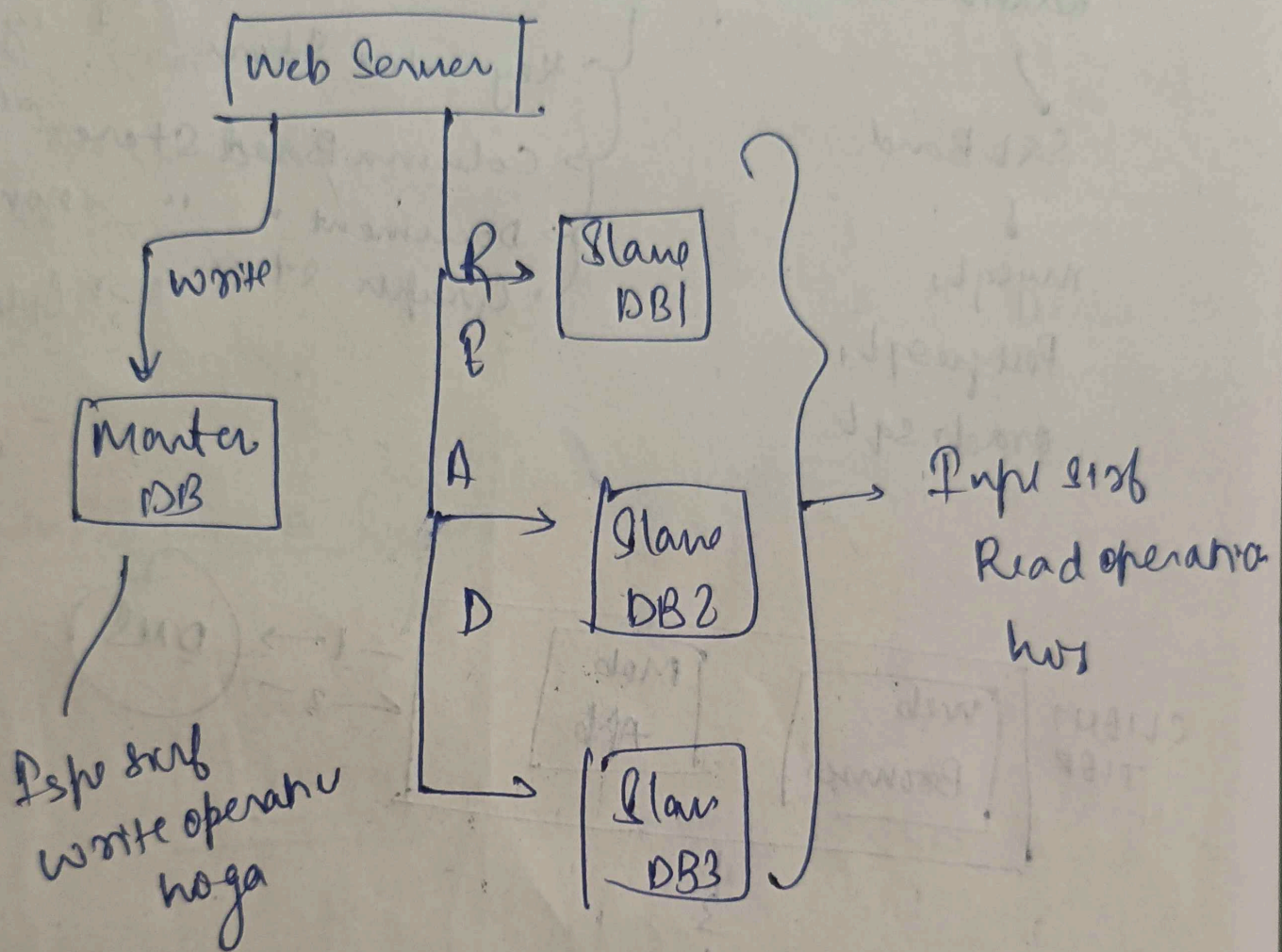
eg:- Graph QL



as ek hi DB,  
no gayato sab  
par on



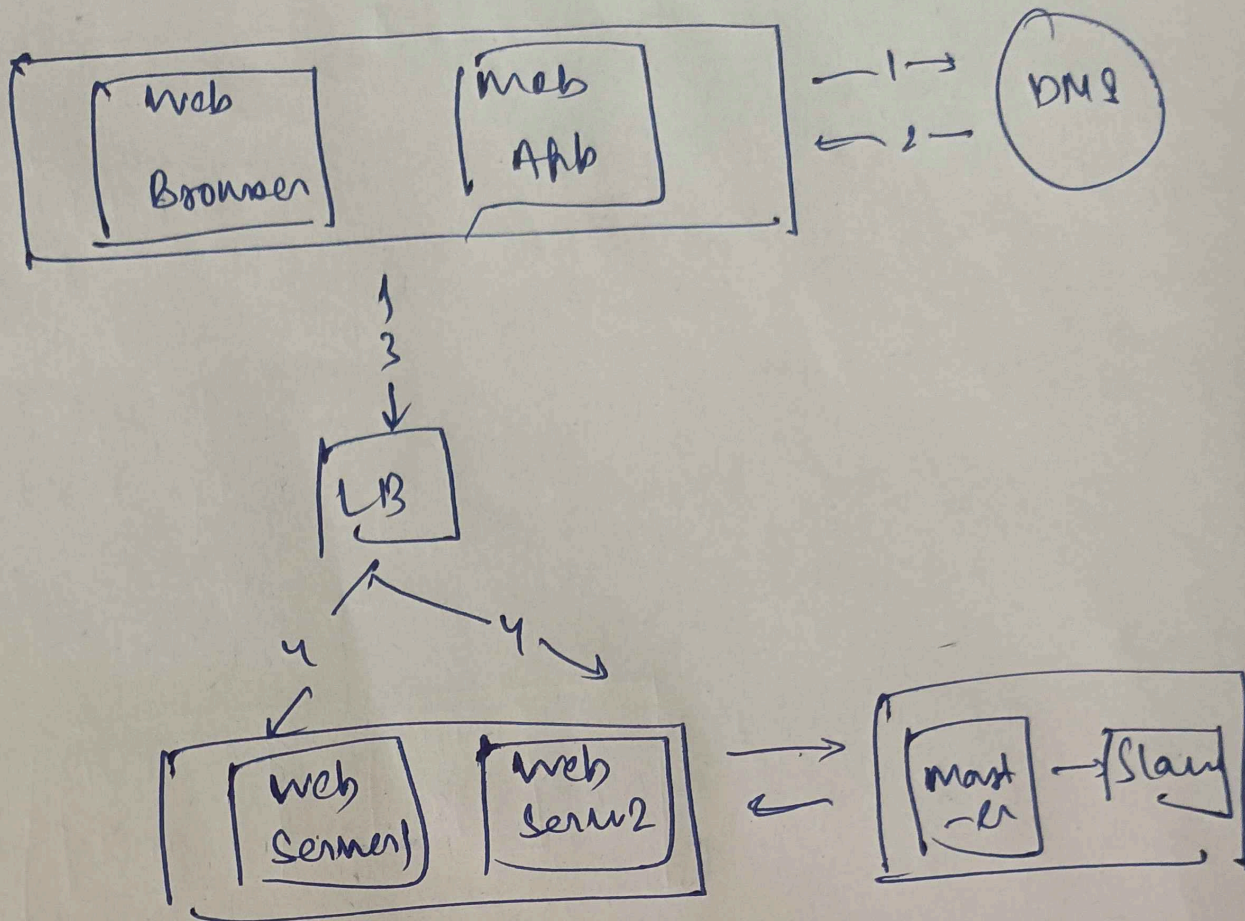
To SPOR se bachne ke liye.



Slaves master ki copy h,

\*) agar master kharab hua to kuch samay ke liye koi slave aa jayega as master.

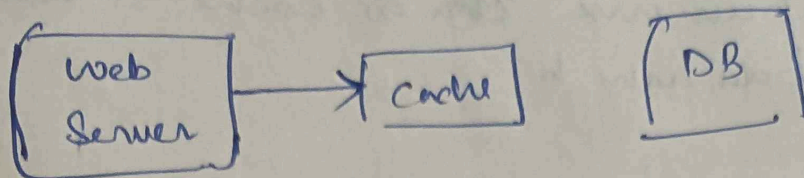




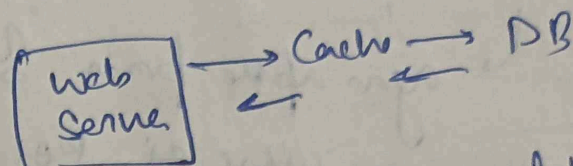
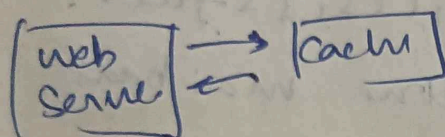


# Cache

→ Faster than DB



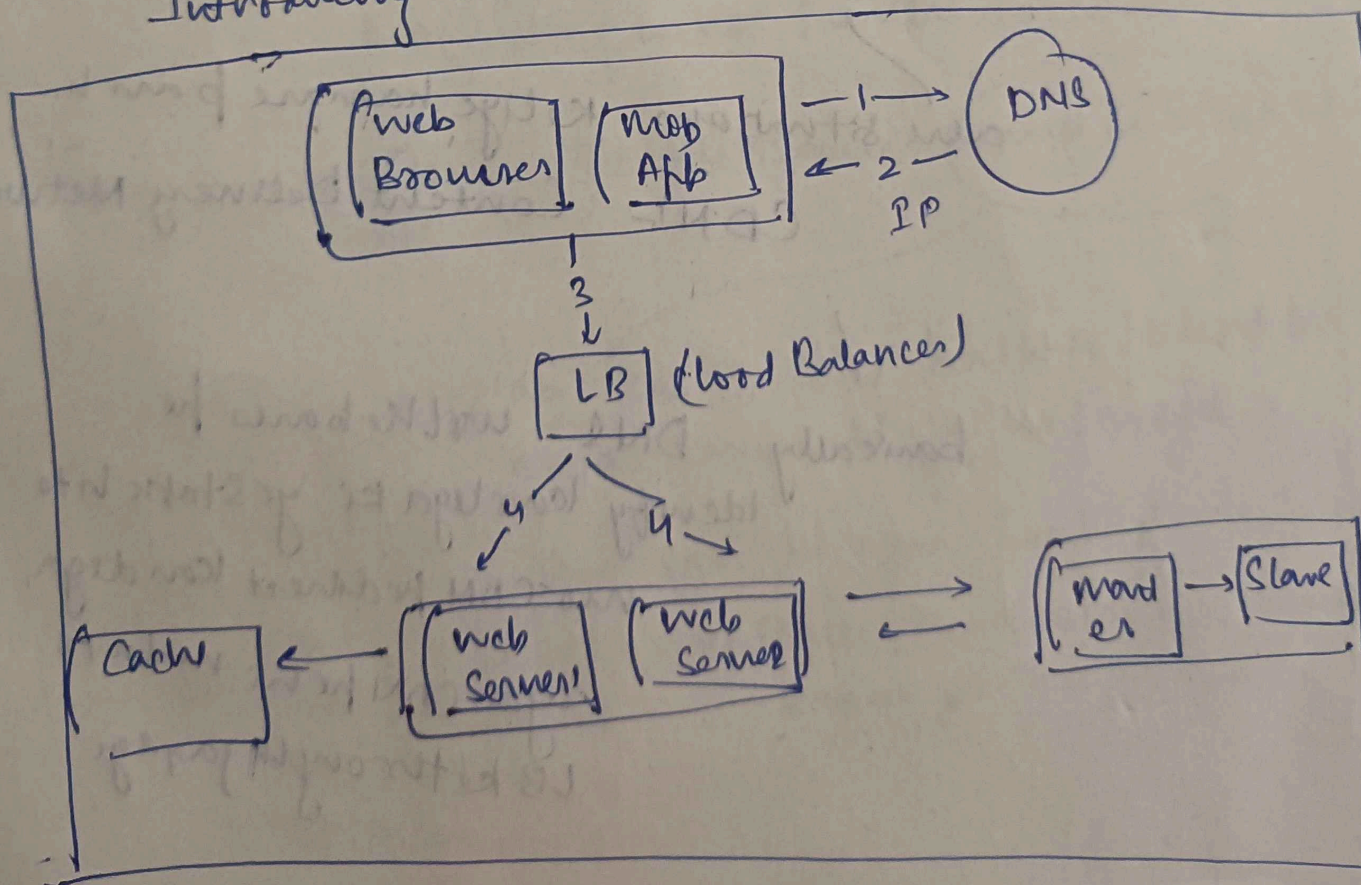
agar  
Cache ko hum  
req ka res nahi  
Cache me



(Cache me un data  
database se aayega,  
Cache me store,  
aur res bhejega)

"Cache Eviction Algo"

## Introducing Cache in Architecture



Abhi dikkat to bot hi aa rahi h, hamari application  
cache se chal rahi h, bahut saare users ko handle  
bhi kar paa rahi h. aur agar user ko jaldi data  
chahiye to hum unke ~~data~~ cache se lake  
de bhi paa rahi h.

lekin

↳ agar apke paas "Static Content" bahut jayda,  
jaise ki koi static HTML, CSS, JS files.  
to wahan mai is architecture mai baar  
baar cache ya DB ko karna sirf  
latency badhega.

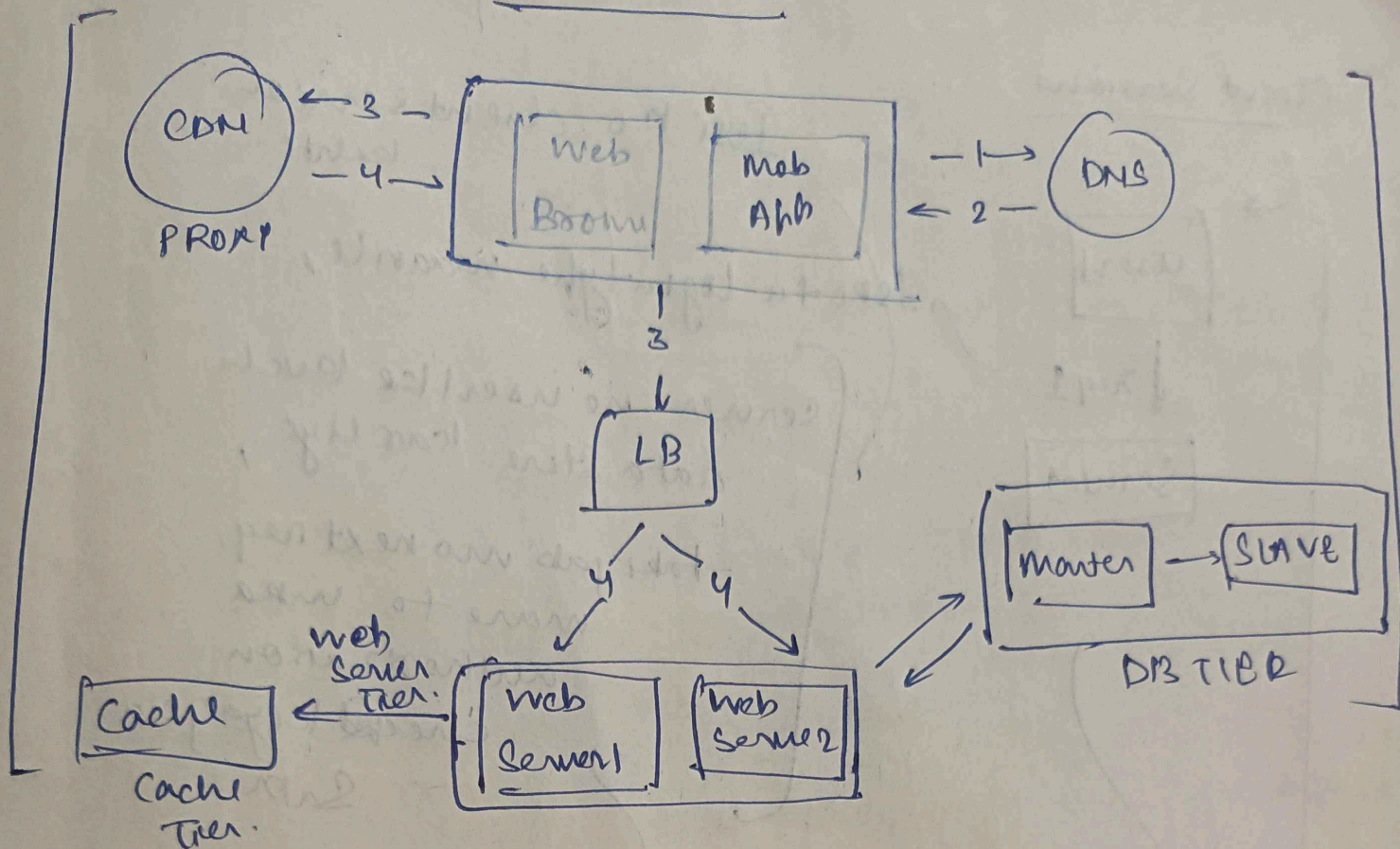
aur situations ke liye hamare paas h.

CDN = Content Delivery Network

basically DNS work ke baad hi  
identify kar lega ki ye static h to  
us CDN ko direct kar dega,  
agar CDN hi nahi milta to  
LB ko through jayegi



## CDN Introduced



AutoScaling is not possible for server.

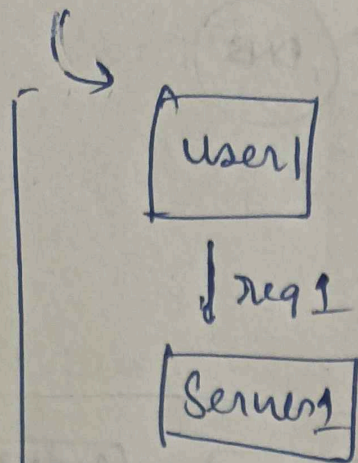
→ Dikkat ye h ki Me WebServer khud se Scale ni ho Sakta.

matlab kal ko load badh to khud ba khud naya server add ni ho jayega,

aur iske hi agar load dekhna to web server ka count khud ba khud kaam ni ho jayega



## Client Sessions



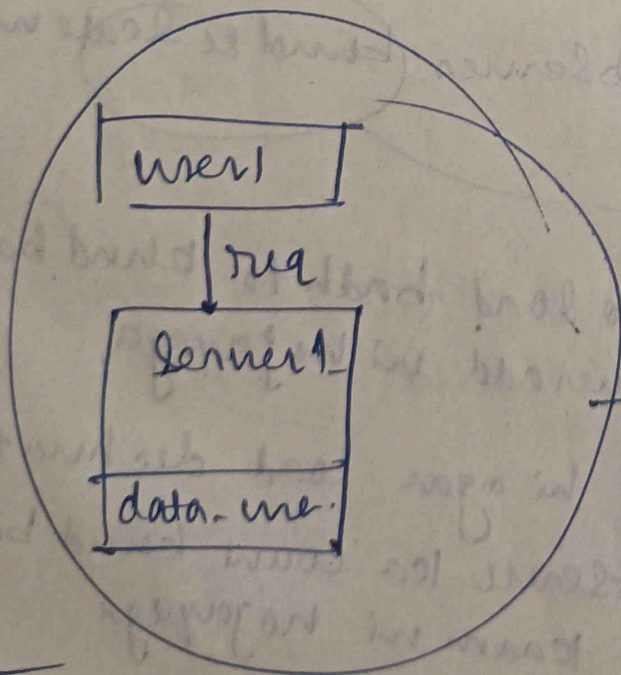
Yesi ko client session  
kehte.

Isko tu login type maante,

Server ne user1 ke kuch  
data store kar liye,

tabi jab wo next req  
mare to uska  
authorization  
check krna pade  
Same

Let say us data ko  
"data-user"  
keh sake.



ab to me

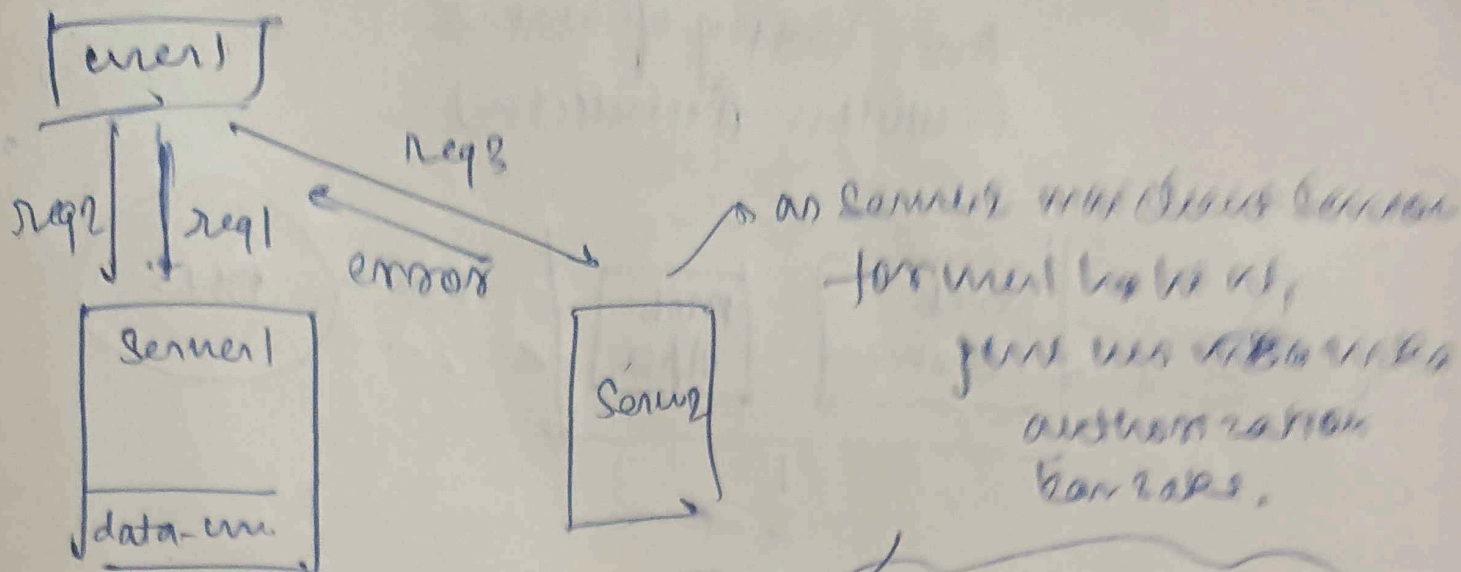
dikhaiye h ki,

agar alga req  
user1 ne Server2  
pe maare to wo  
error milega or

Server2 main  
"data-user" h ki hi

"Statefull"  
architecture

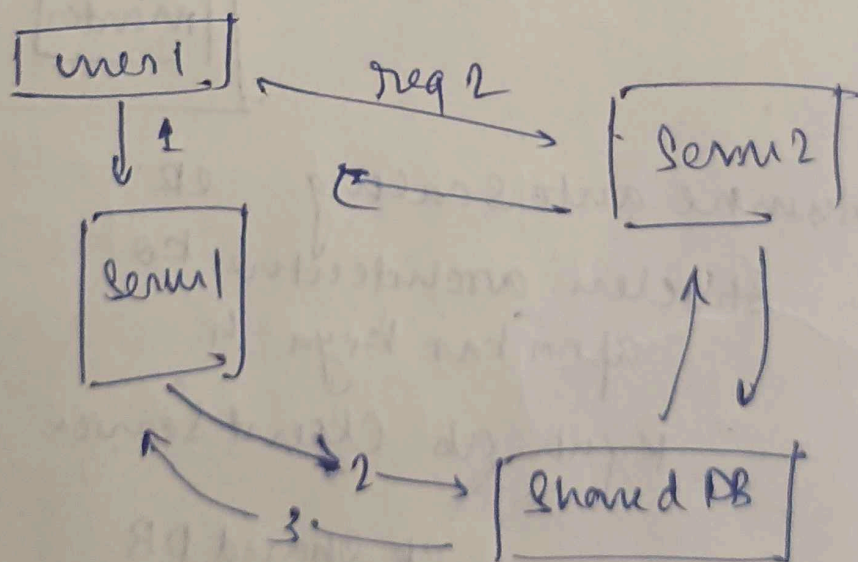




to matlab load Balancer ko ye yaad rakhna padiga ki wo user ka req server pehla bheje.

Iska soln h Stateless Architecture

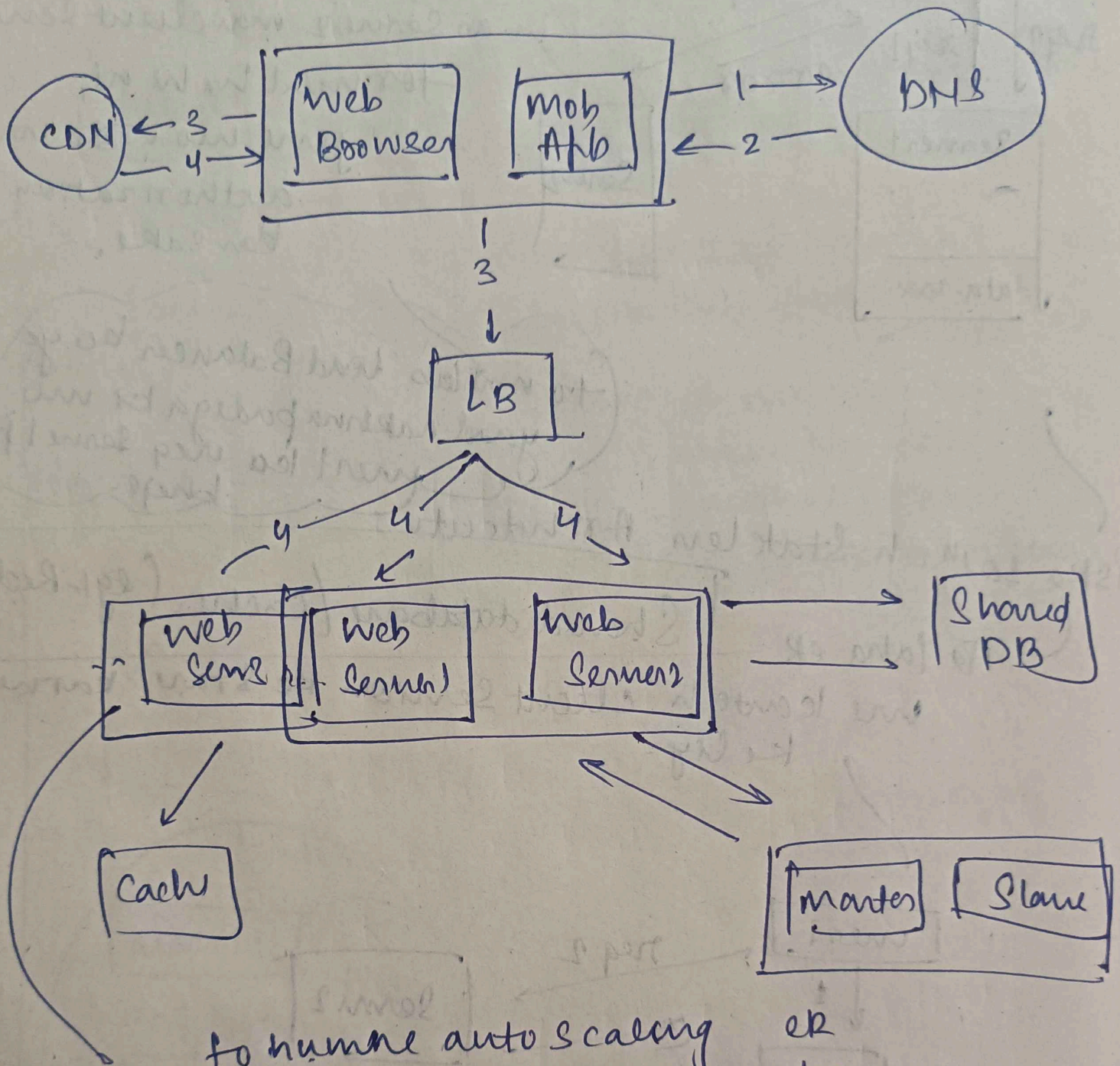
Yaha ek Shared database / cache (eg. Redis) me karte h client session ko store karne ke liye.



Yaha for client session store honge.



# Auto Scaling of Server. (Stateless Architecture)



to humne auto scaling ek  
stateless architecture ko  
apna kar liya h.

Kyun ki ab Client Server

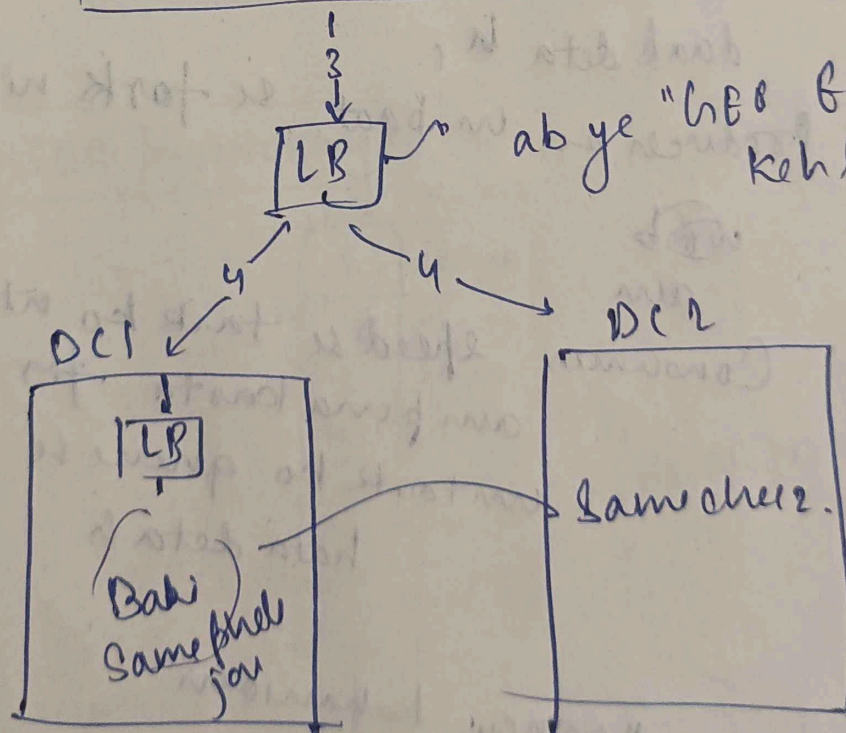
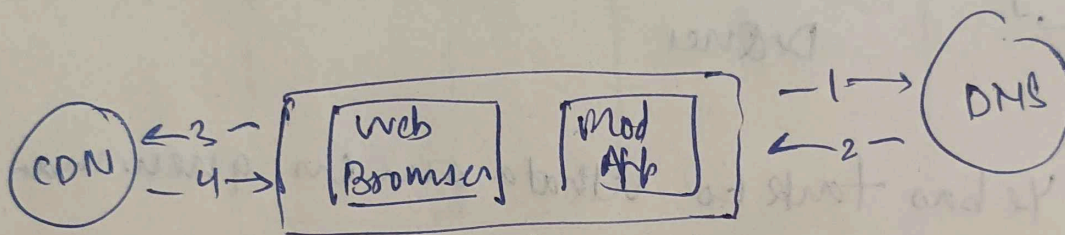
ek Shared DB  
mai store hoga  
naa hi har ek  
server mai



Statku ka matlab h ki DIB ko koi matlab n  
 kaun mer aaga h, unko aur  
 authenticate kar de raha h.

kyon ye international ko server karne mai kafi time  
 lagayega, iseliye hum data regions banayenge.

or  
 data center.  
 e.g. CR India, eb America



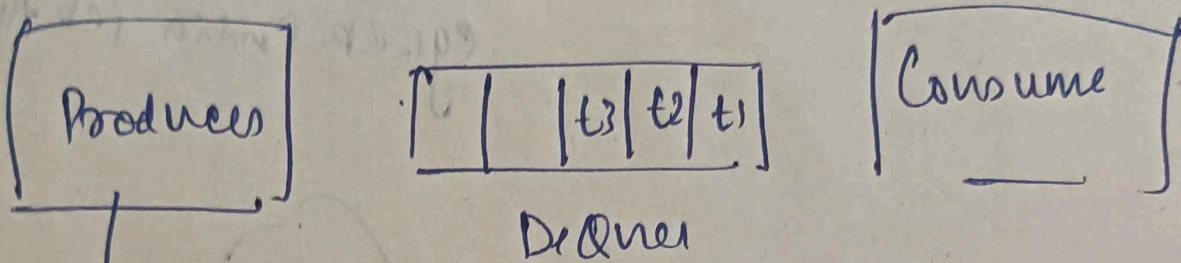
ab ye "Geo Enabled DNS"  
 kehlayega



# Managing Queue (or PUB-SUB model)

eg:- Kafka, RabbitMQ

Isme do parts hote h  
→ Producer  
→ Consumer



Ye bas task ko uthata aur us queue mein  
daal deta h,

Producer ko in baat se fark ni pata

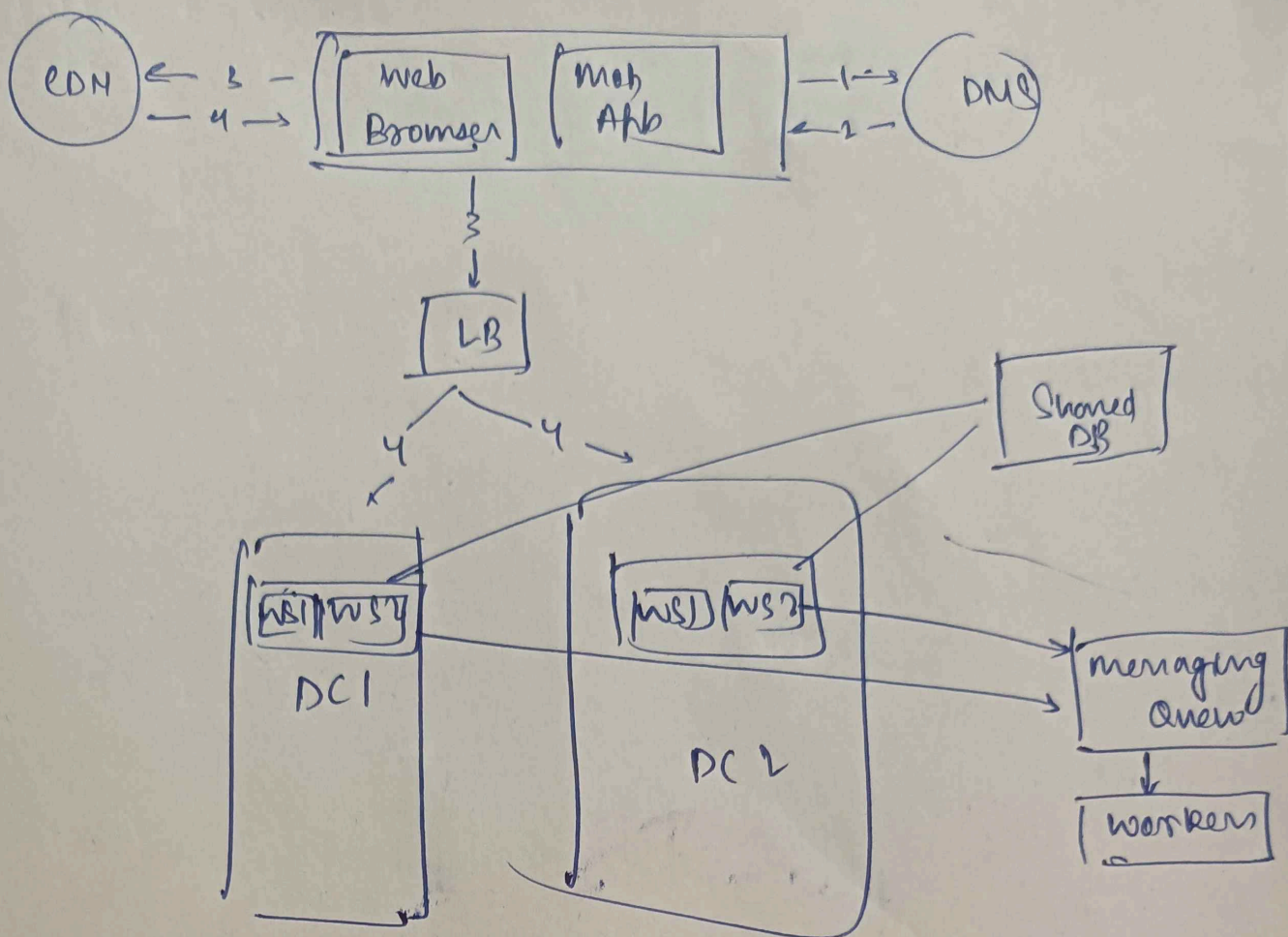
~~work~~

aur

Consumer speed se task ko uthate  
aur pura karta fir  
untask ko queue se  
hata deta h

-> async behaviour.





### Sharding (DB partitioning)

