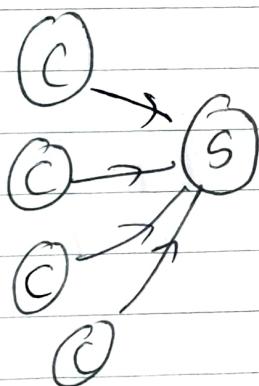


Load Balancers

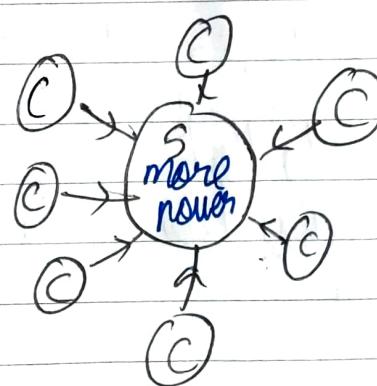
Load Balancer → (i) A server that sits b/w a set of clients and set of servers which distributes load / routes traffic evenly across all servers in our set.

- (ii) This prevents 1 server from getting too overloaded by too many client reqs.
- (iii) Also clients do not know abt load balancers' existence so load balancer = reverse proxy

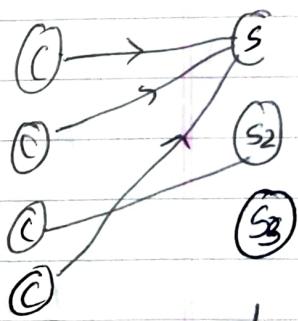
Case 1: Server overload due to lots of reqs.



Case 2: Vertical scaling: power of servers ↑ but still lots of reqs, so overlod

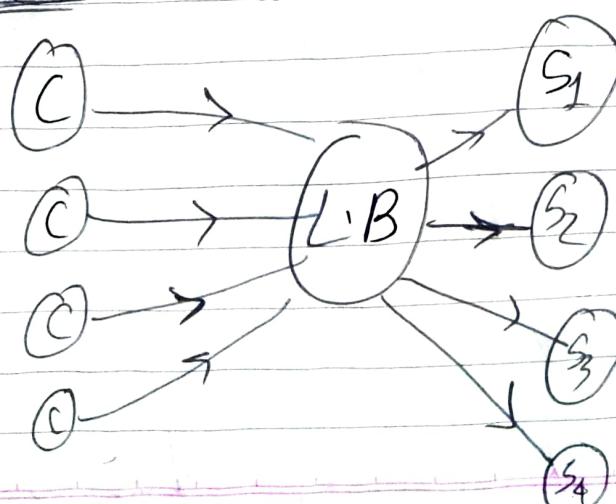


Case 3: Horizontal scaling
Multiple servers but still most client reqs to 1 server.



S1: overloaded

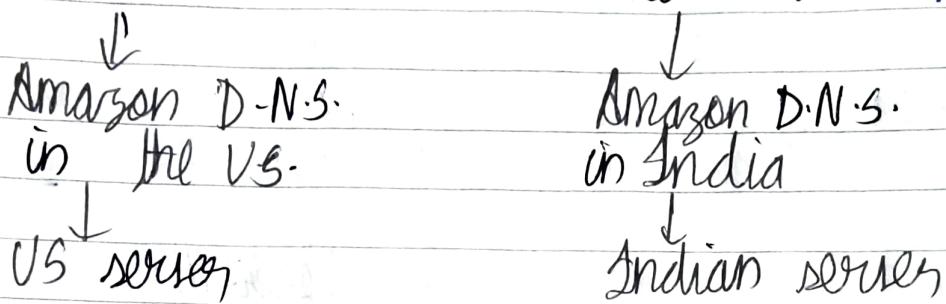
Case 4: Load Balancer



① Req's evenly distri.
② No overload of servers

③ Throughput ↑
④ Latency ↓ (no server clogged)

Other places where load balancing can be done:



Q2: Netflix VS ~~VS~~ ~~Notif~~ Netflix India

e.g. On opp page \rightarrow we are accessing diff. IP addresses of google.com i.e. diff. servers and we get allotted the server related to our specific IP address (abbr. IPA) by the load balancer.

NOTE

Hardware LB → physical machine dedicated to load balancing

Software LB → ~~more~~ not physical machines
(e.g. - a cloud services LB → ~~cloud server~~)
use: More power, customization

SOFTWARE L.B.

- 1) How does LB know which servers it can access?

i) sys. designer config. LB and servers to know abt each other

i) sys. designers config. L-B and servers to know abt each other

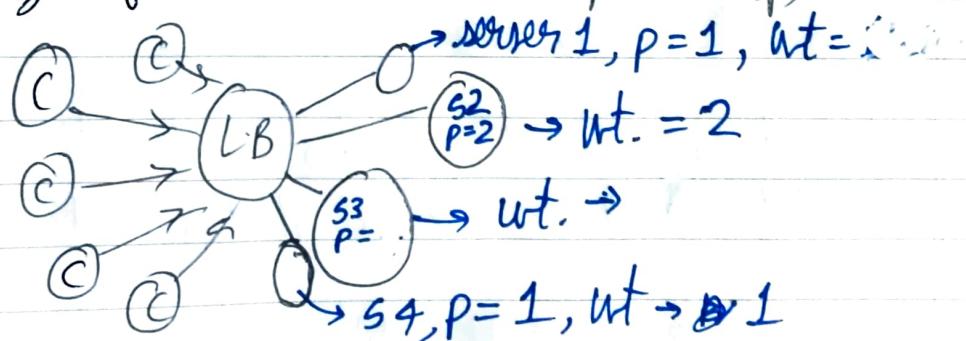
ii) On adding / removing old servers, it registers / deregisters itself w.r.t. LB

2) Ways to select servers to direct load to

① Random redirect → ~~No logic used. Just random selecⁿ of server from available servers~~
by chance, one server may get overloaded

② Round Robin → LB goes thru all servers in order and allot a batch of reqs. to a server and so on until all servers have equal load so, for next batch of reqs., it'll again start from 1st server
↓
server's power not optimally used.

Let size of server & its power (P)



③ Weighted Round Robin : LB goes as per round robin approach but if a weight is added on 2nd server, LB will send it a couple more reqs.
∴ Power of each server → optimized

Request Batch Round Robin(RR) Wtd. R.R.

1	S1	S1
2	S2	S2 } avg wt = 2
3	S3	S2 }
4	S4	S3 } wt = 3
5	S1	S3 }
6	S2	
7	S3	S4

(3) Performance/Load Based: LB will do performance/health checks on servers

- how much traffic a server is handling at any given time
- how long a server is taking to respond to traffic
- how many resources server uses
- etc.

Based on these checks, it allocates reqs. (more reqs. to low : high performance servers)

(4) IP based server selec: On getting requests from clients, L.B. hashes the IP address of the client and depending on this hashed value, the client's req. gets directed to a specific server

Use :- Caching

e.g.: - You're caching results of certain reqs.

in your servers, it's helpful to redirect the client to a series in which, the ~~req~~ response to the request ~~is~~ has already been cached

- eg:-
- i) Client 1 searches for "cheese pizza"
 - ii) ~~LB~~ knows that, the response to this request is ~~is~~ in server 1
 - iii) LB hashes down IPA of client to say "1" → basically, a VID
 - iv) Now, it sees that in server S3 a result for ~~req~~ a prior req. made with VID = 1 is cached and ready
 - v) LB redirects client 1's req. to server S3
 - v) Client 1 gets his/her response super fast.

⑤ Path-based → LB selects ^{server or set of servers} based on path of ~~the reqs~~.

eg:- AlgoExpert has servers S1, S2, S3, ..., S25

All reqs related to payments } → S1, S2, ..., S9

All reqs related to running code → S10 - S18

All reqs related to watching vids → S19 - S21, S22, S24, S25

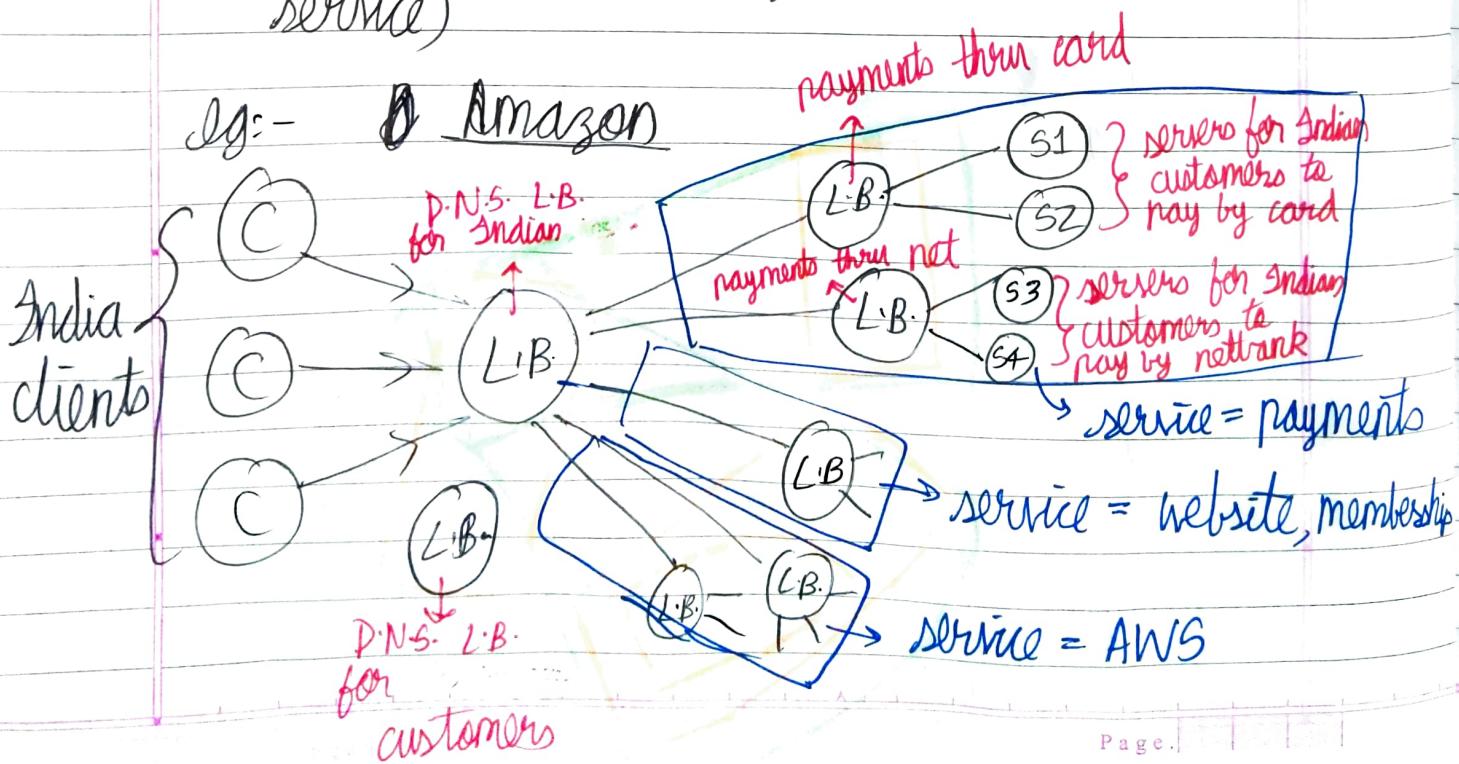
All reqs related to "about" page → S23

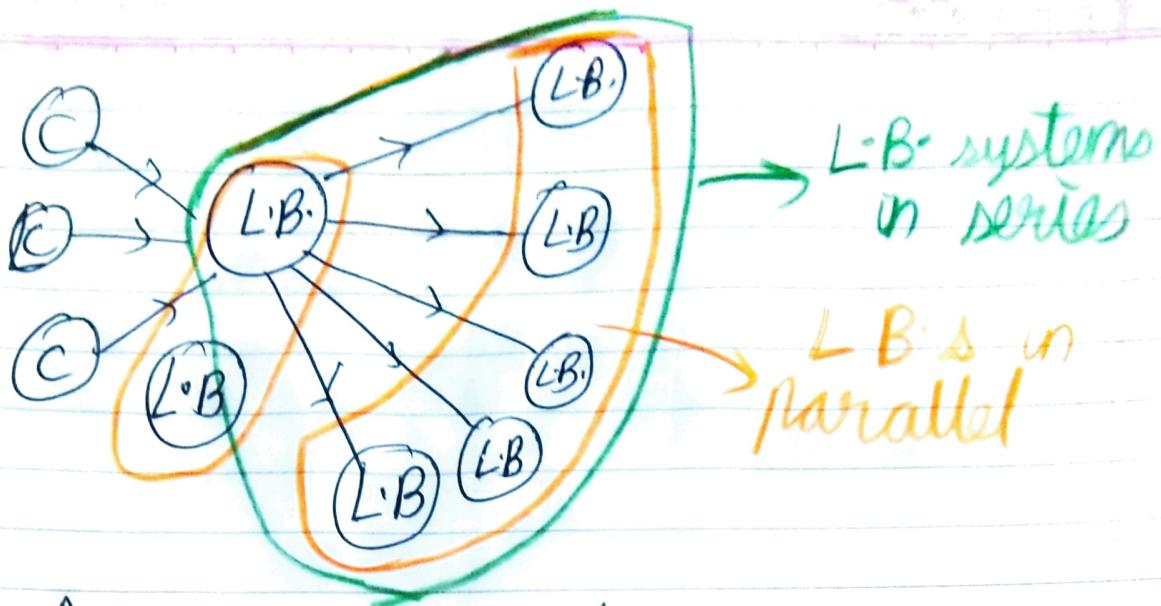
Use:

- 1) Any deployment of big change to any service (eg:- running code service) affects only the servers allotted to this service.
- 2) If one set of servers for a service start failing, atleast other services are still up and running. (unaffected)

MULTIPLE LOAD BALANCERS

- 1) In series → To keep divide a batch of reqs into smaller batches by based on some criteria (path eg:- path) on the req.
- 2) In parallel → To ① just prevent 1 L.B. itself from getting overloaded or ② Redirect to diff. servers based on req. criteria or server criterial (eg:- performance, diff kinda service)





Amazon eg: - Just L-B-S

code on opp. page → shows used R.R. selecⁿ meth.

```
nginx.conf — load_balancing

❶ nginx.conf × JS server.js

❷ nginx.conf
1 events { }
2
3 http {
4     upstream nodejs-backend {
5         server localhost:3000 weight=3;
6         server localhost:3001;
7     }
8
9     server {
10        listen 8081;
11
12        location / {
13            proxy_set_header systemexpert-tutorial true;
14            proxy_pass http://nodejs-backend;
15        }
16    }
17 }
```

```
elements-MBP:load_balancing clementmihailescu$ PORT=3001 node server.js
stening on port 3001.
'systemsexpert-tutorial': 'true',
host: 'nodejs-backend',
connection: 'close',
'user-agent': 'curl/7.54.0',
accept: '*/*' }
'systemsexpert-tutorial': 'true',
host: 'nodejs-backend',
connection: 'close',
'user-agent': 'curl/7.54.0',
accept: '*/*' }
'systemsexpert-tutorial': 'true',
host: 'nodejs-backend',
connection: 'close',
'user-agent': 'curl/7.54.0',
accept: '*/*' }
```

You can see that the requests get distributed