

Load Balancers

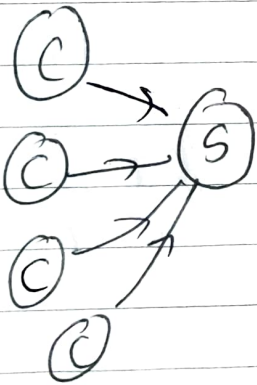
balances workload

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

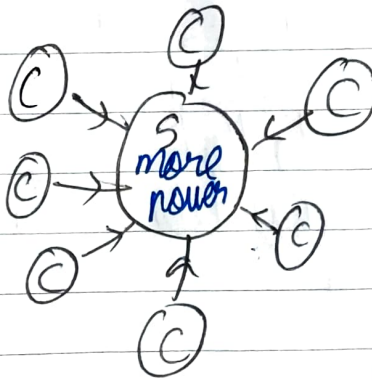
(ii) This prevents 1 server from getting too overloaded by too many client reqs.

(iii) Also clients do not know abt load balancers' existence so load balancer = reverse proxy

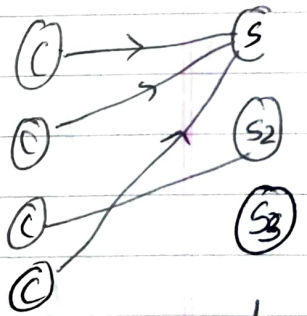
Case 1: Server overload due to lot of reqs.



Case 2: Vertical scaling: (power of server ↑ but still old power, so overload)

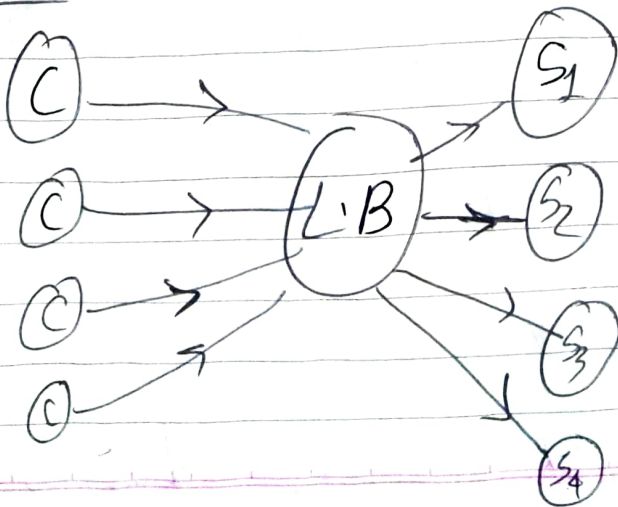


Case 3: Horizontal scaling: Multiple servers but still most client reqs to 1 server.



S1: overloaded

Case 4: Load Balancer



reqs. evenly distrib.
① No ~~too~~ overload of server

② Throughput ↑

③ Latency ↓ (no server clogged)

Other places where load balancing can be done:

- 1) b/w servers and dB
2) at D.N.S layer i.e. same ~~or~~ service → multiple servers so LB directs you to multiple servers.
eg:- 'amazon' → single domain gets multiple IP addresses ($\frac{1}{1}$ for .com) for multiple servers.
- amazon.com amazon.in
- ↓ ↓
- Amazon D.N.S. Amazon D.N.S.
in the US. in India
- ↓ ↓
- US server Indian server

eg2: Netflix US VS ~~Netflix~~ Netflix India

eg 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 L.B. → physical machine dedicated to load balancing.

Software L.B. → ~~more~~ not physical machines
(eg:- a cloud services L.B. → cloud server)
use: more power, customizaⁿ

SOFTWARE L.B.

- 1) How does L.B. know which servers it can access?

(i) Sys. designers config. L-B and server to know
abt each other

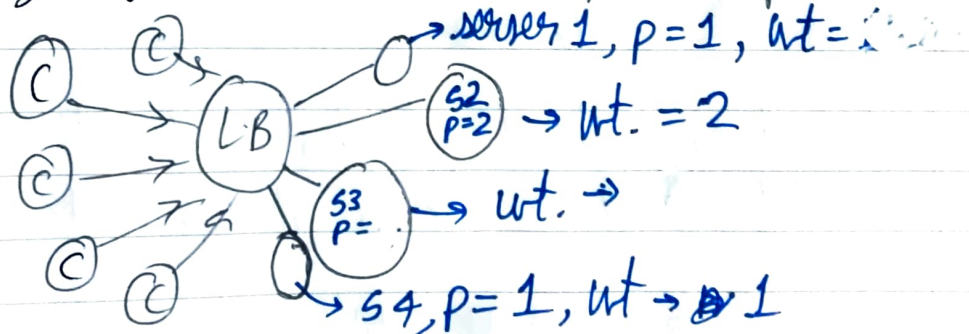
(ii) On adding / removing old server, it registers / deregisters itself w.r.t. LB

2) Ways to select servers to direct load to

(1) Random router \rightarrow ~~No~~ No logic used. Just random selectⁿ of server from available servers
 \downarrow
 by chance, one server may get overloaded

(2) Round Robin \rightarrow L.B. 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 new batch of reqs. \rightarrow it'll again start from 1st server
 \downarrow
 server's power not optimally used.

Let size of server \propto its power (p)



(3) Weighted Round Robin : L.B. goes as per round robin approach but if a weight is added on 2nd server, L.B. will send it a couple more reqs.
 \therefore Power of each server \Rightarrow optimized

Request Batch	Round Robin (RR)	Wt'd. R-R.
1	S1	S1
2	S2	S2
3	S3	S2
4	S4	S3
5	S1	S3
6	S2	S3
7	S3	S4

$\left. \begin{matrix} S2 \\ S2 \end{matrix} \right\} \text{avg wt} = 2$
 $\left. \begin{matrix} S3 \\ S3 \\ S3 \end{matrix} \right\} \text{wt} = 3$

③ Performance/Load Based: L.B. 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 ^{is the} server using
- etc.

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

④ IP based server select: 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

eg:- You're caching results of a certain req.

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

- eg:-
- i) Client 1 researches for "cheese pizza"
 - ii) ~~Client 1 knows that the response to this request is in server 1~~
 - iii) L.B. 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) L.B. redirects client 1's req. to server S3
 - vi) Client 1 gets his/her response super fast.

⑤ Path-based → L.B. 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

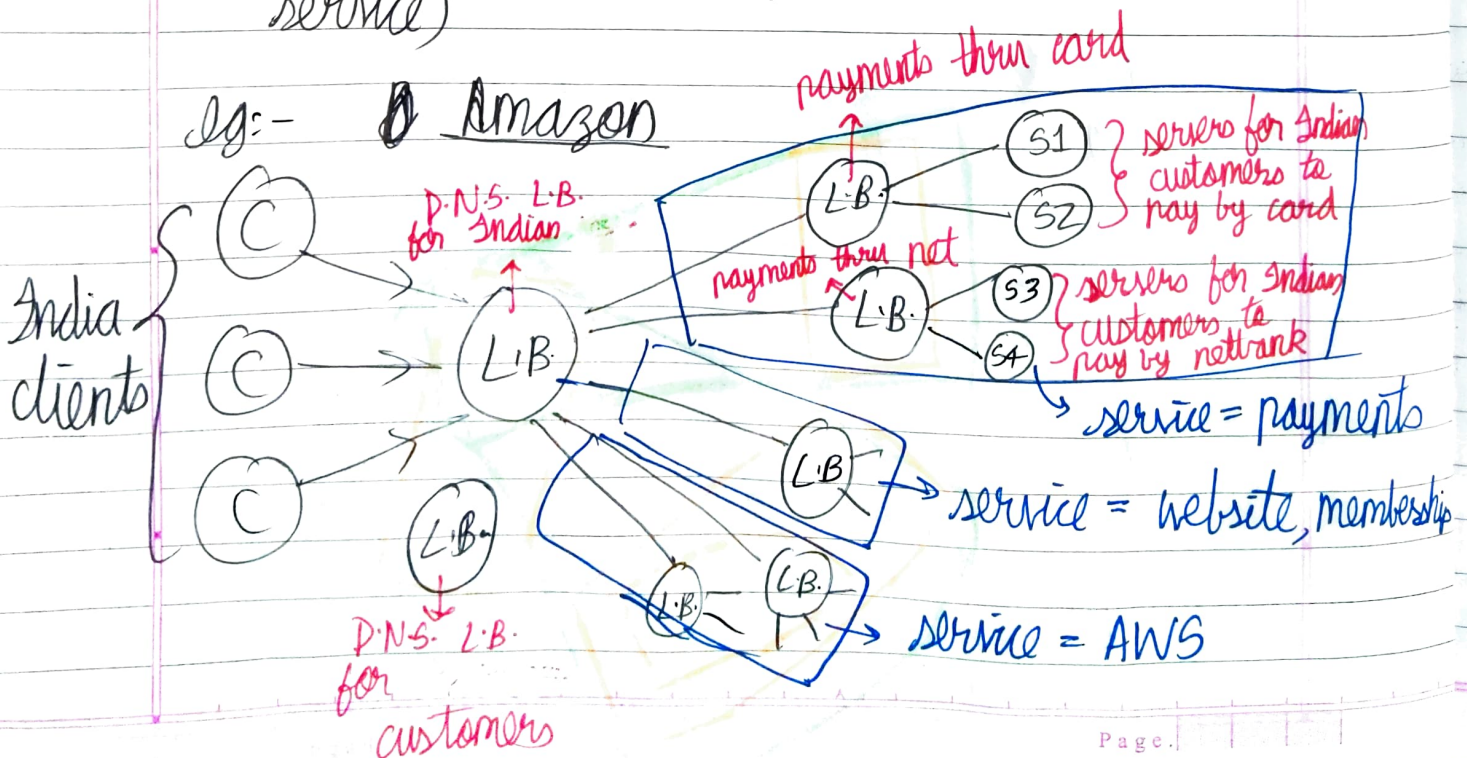
Use:

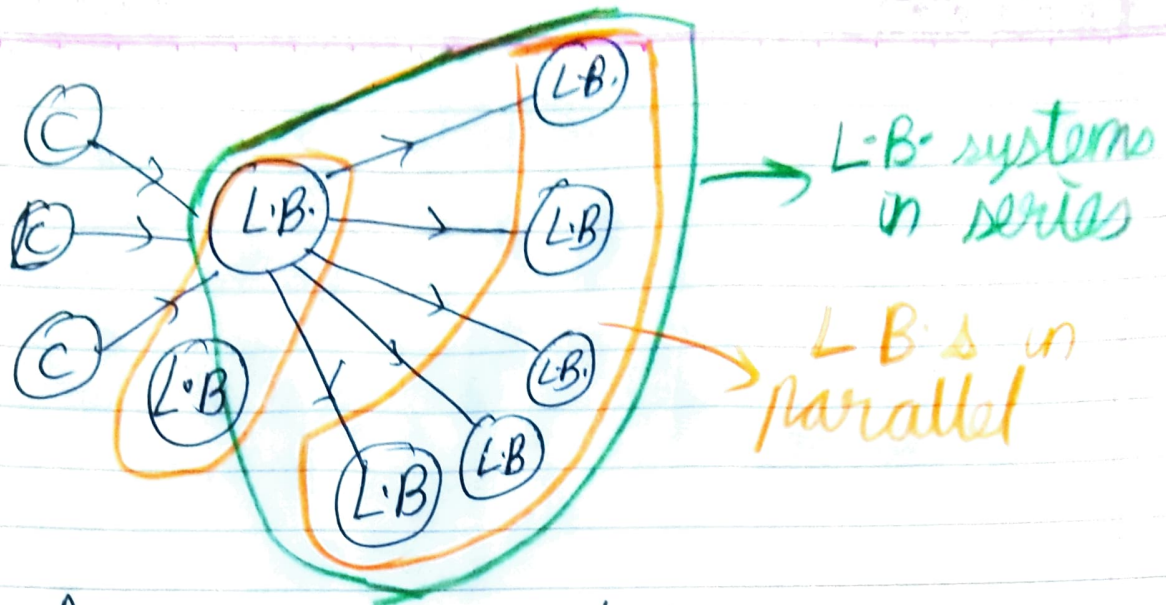
- 1) ~~Any~~ Deployment of big change to any service (eg:- running code service), affects only the servers allotted to this service.
- 2) If one set of servers for a service start failing, atleast other services are still up and running. (unaffected)

MULTIPLE LOAD BALANCERS

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

eg:- Amazon





Amazon eg: - Just L.B.s

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