

Vertically scaled

①

cost implications

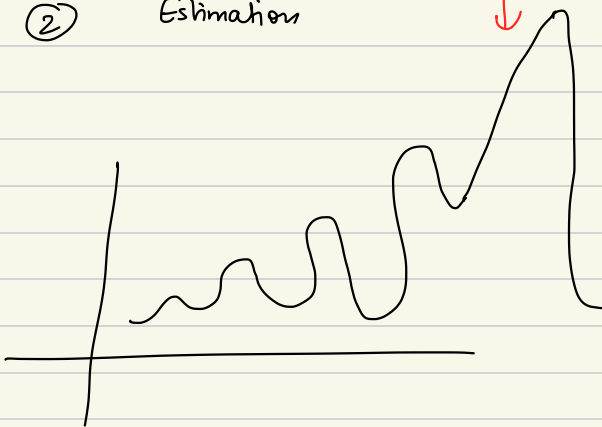


②

Estimation



(scale flexibility is lost)



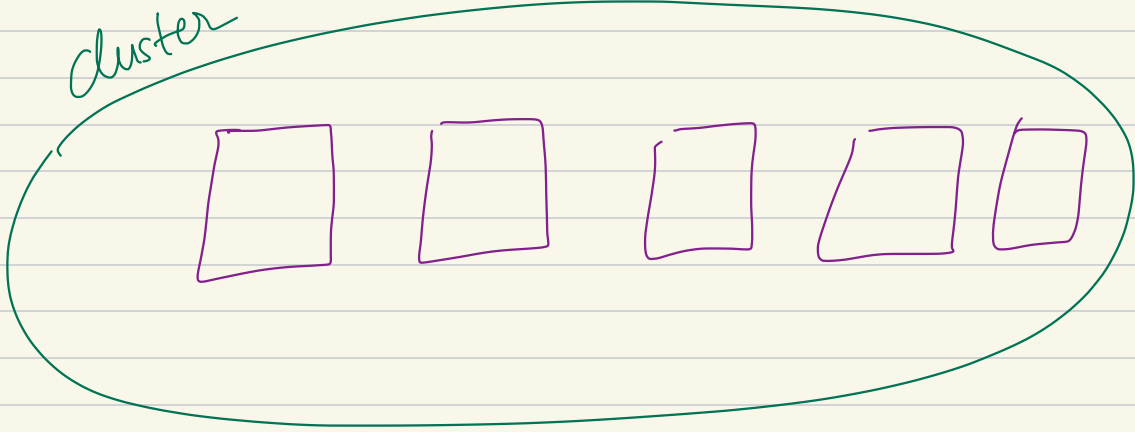
③

SPOF

④



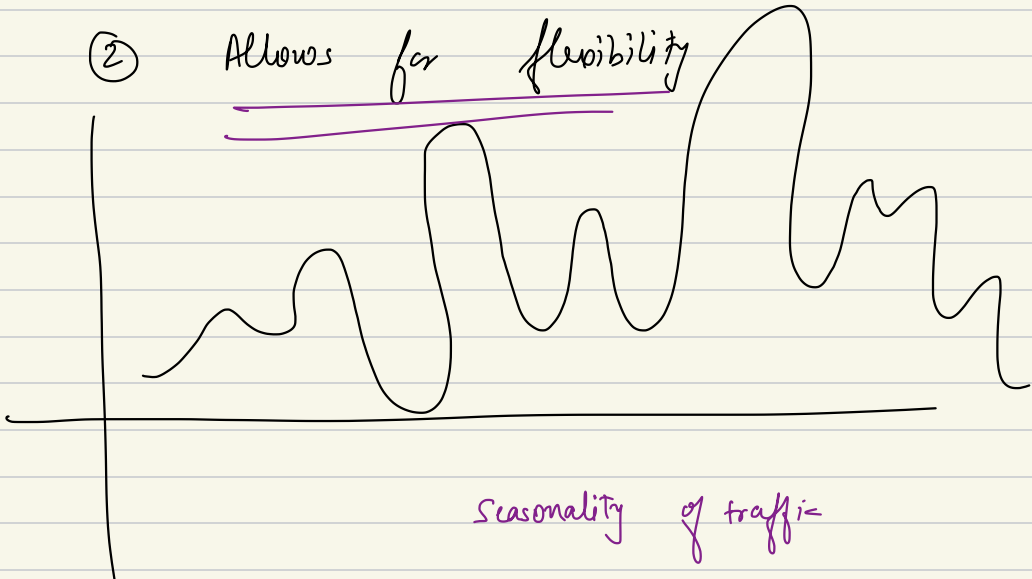
HORIZONTAL SEALING



commodity hardware

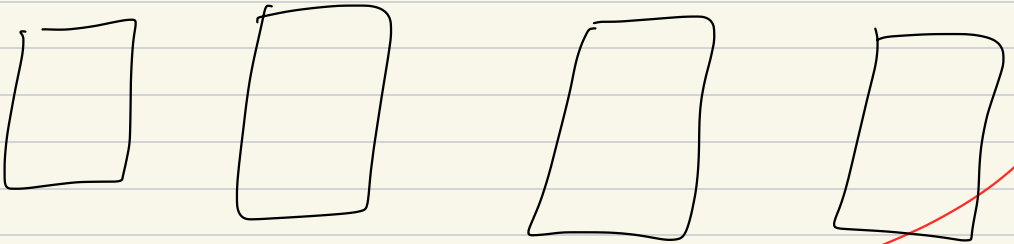
① cheaper

② Allows for flexibility





③ No SPOF

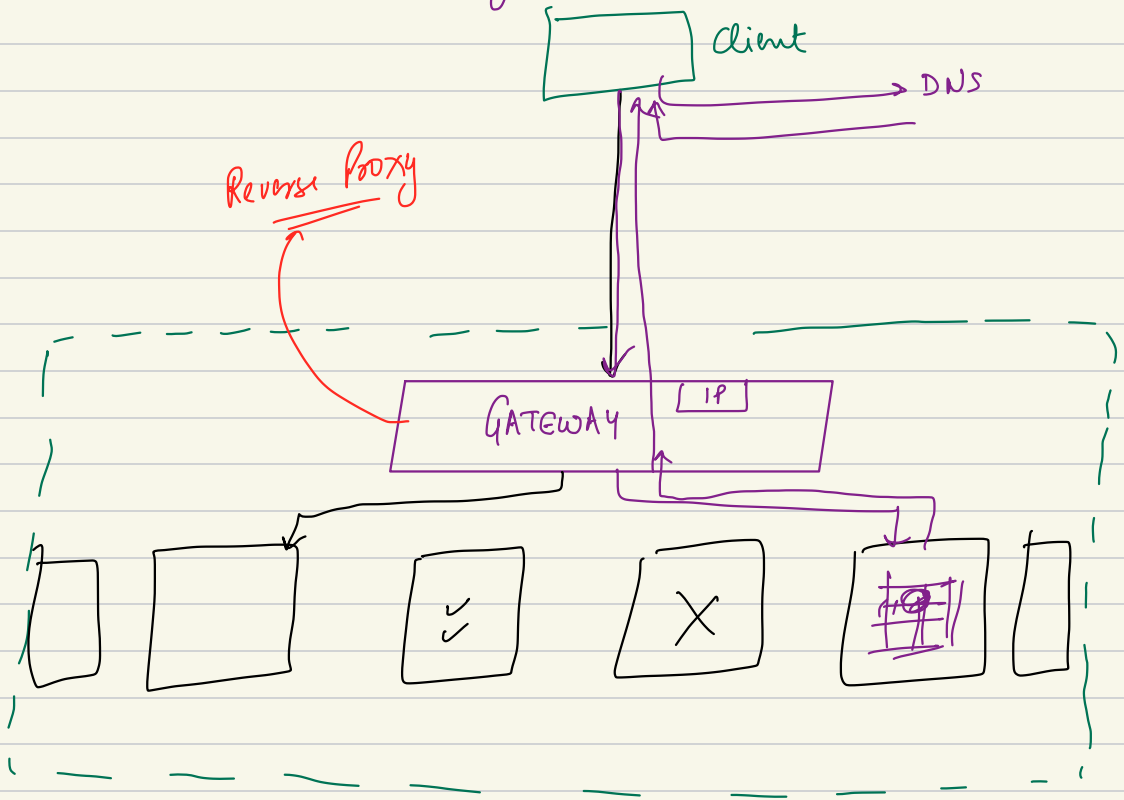


These machines have to communicate with each other over network..

Horizontal scaling becomes tougher as you
have to navigate the challenges of
machines talking to each other over
network

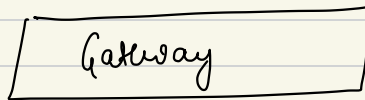
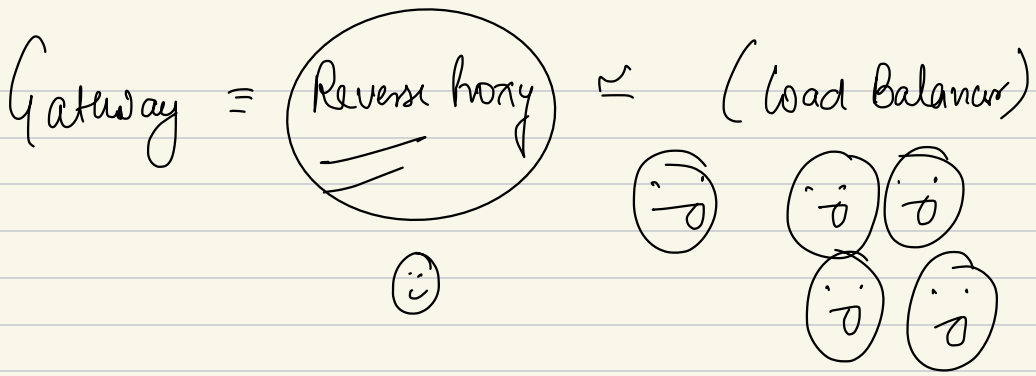


- ① Buy a domain
- ② static public IP address
- ③ DNS entry (domain \leftrightarrow IP)



✓ Gateways usually follow

Active Passive Gatewaying



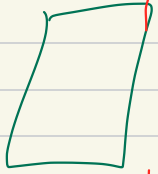
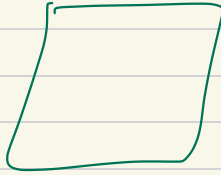
Gateway machine can (usually does) also perform the role of Load Balancer

and hence some people also call this machine a load balancer 😊

Gateway \rightarrow First point of contact for the outside world.

Load Balancer \rightarrow When you have many incoming requests and you want a machine to distribute these requests amongst different machines of the cluster.

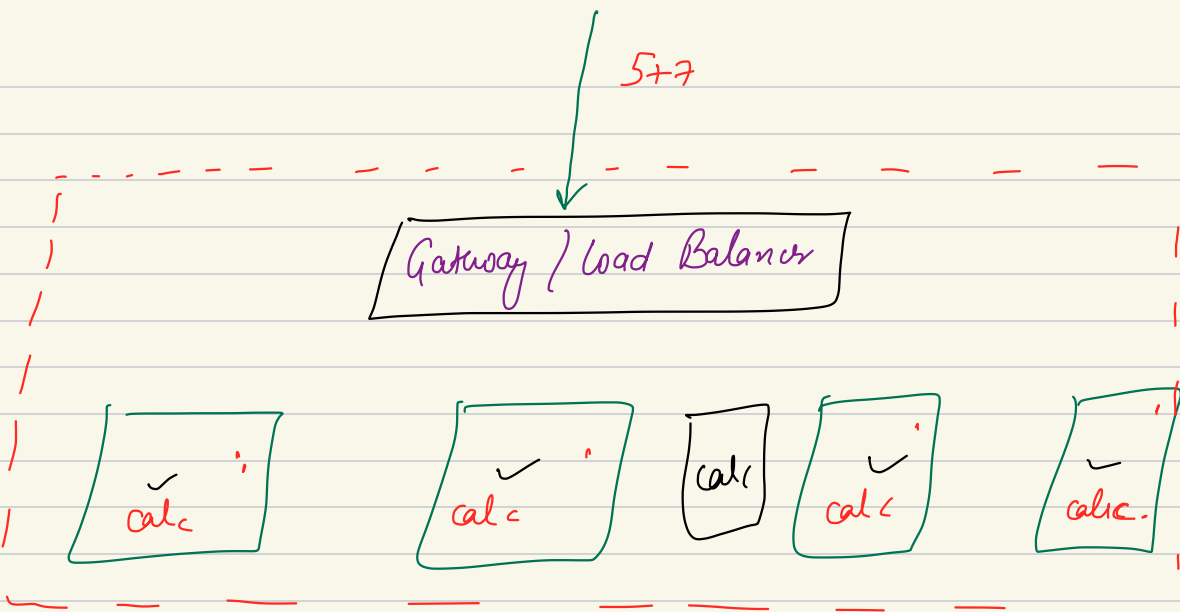
Gateway / Load Balancer



2 kinds of Load Balancing

① stateless Load Balancing

② stateful Load Balancing



if each of the machines in the backend cluster are equally well equipped to handle incoming requests.

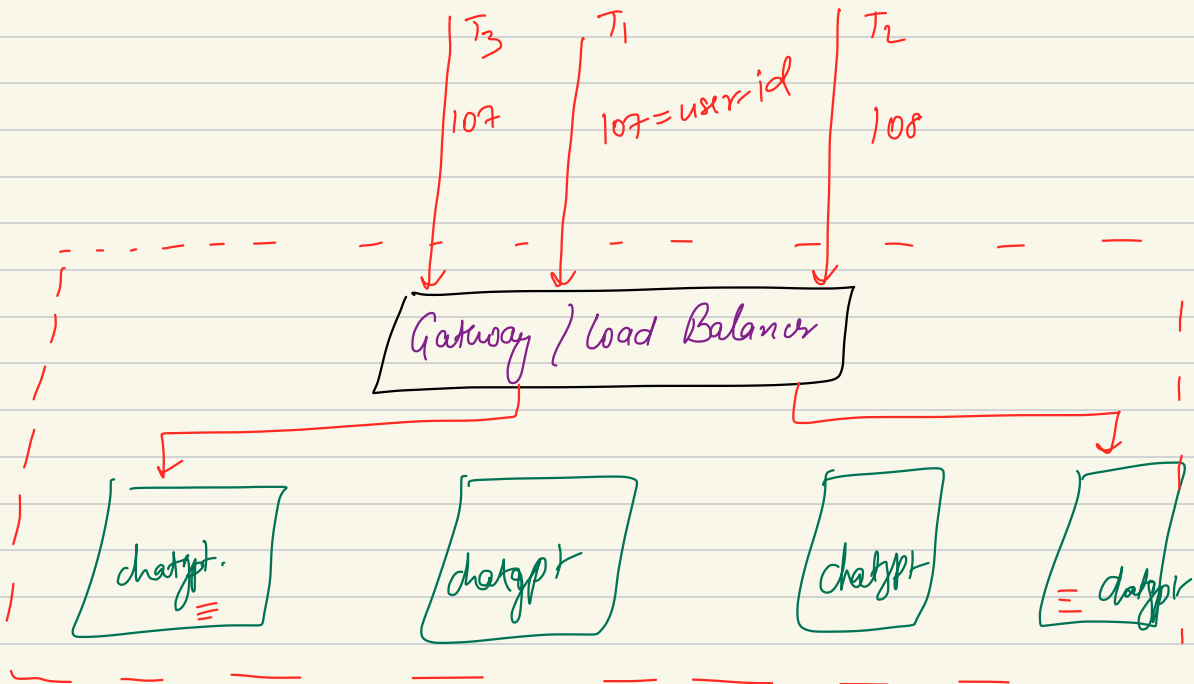
"stateless load balancer"

Algo:

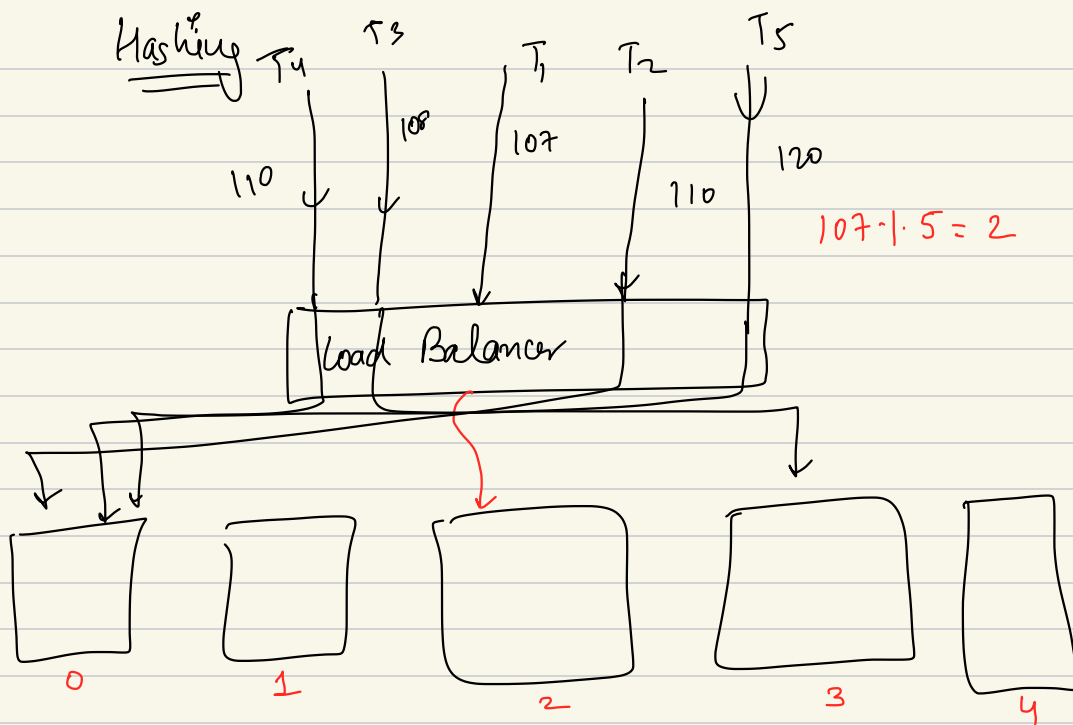
① Round Robin ✓

② Random

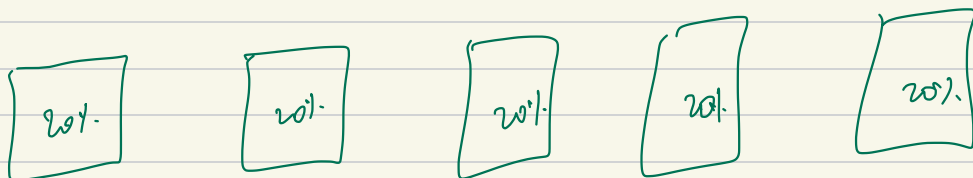
Stateful Load Balancing

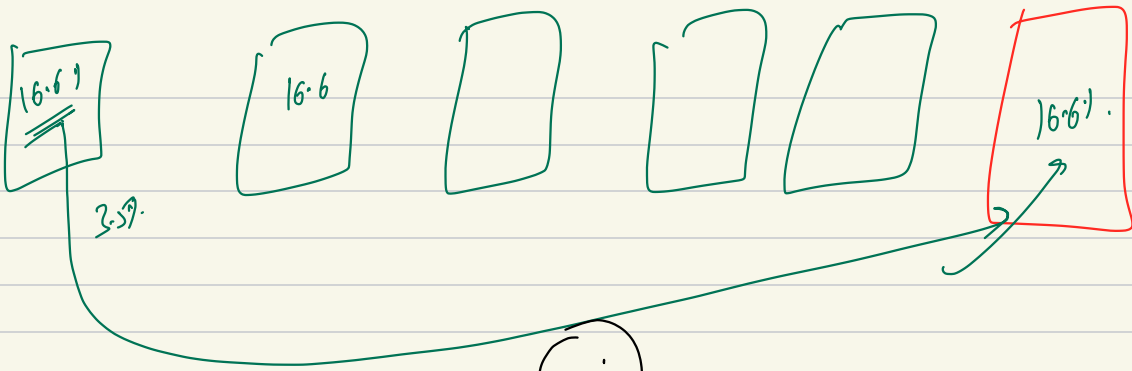


~~Round Robin~~



Simple Hashing 😊 😊





Hashing

Diagram showing a sequence of circles. The first two circles are outlined in red and contain a dot and a slash. The next three circles are outlined in purple and contain a dot and a slash. The word "Hashing" is written in purple between the red and purple circles.

	1/5	1/6
101	1	5
102	2	0
103	3	1
104	4	2

↑
.
.
.

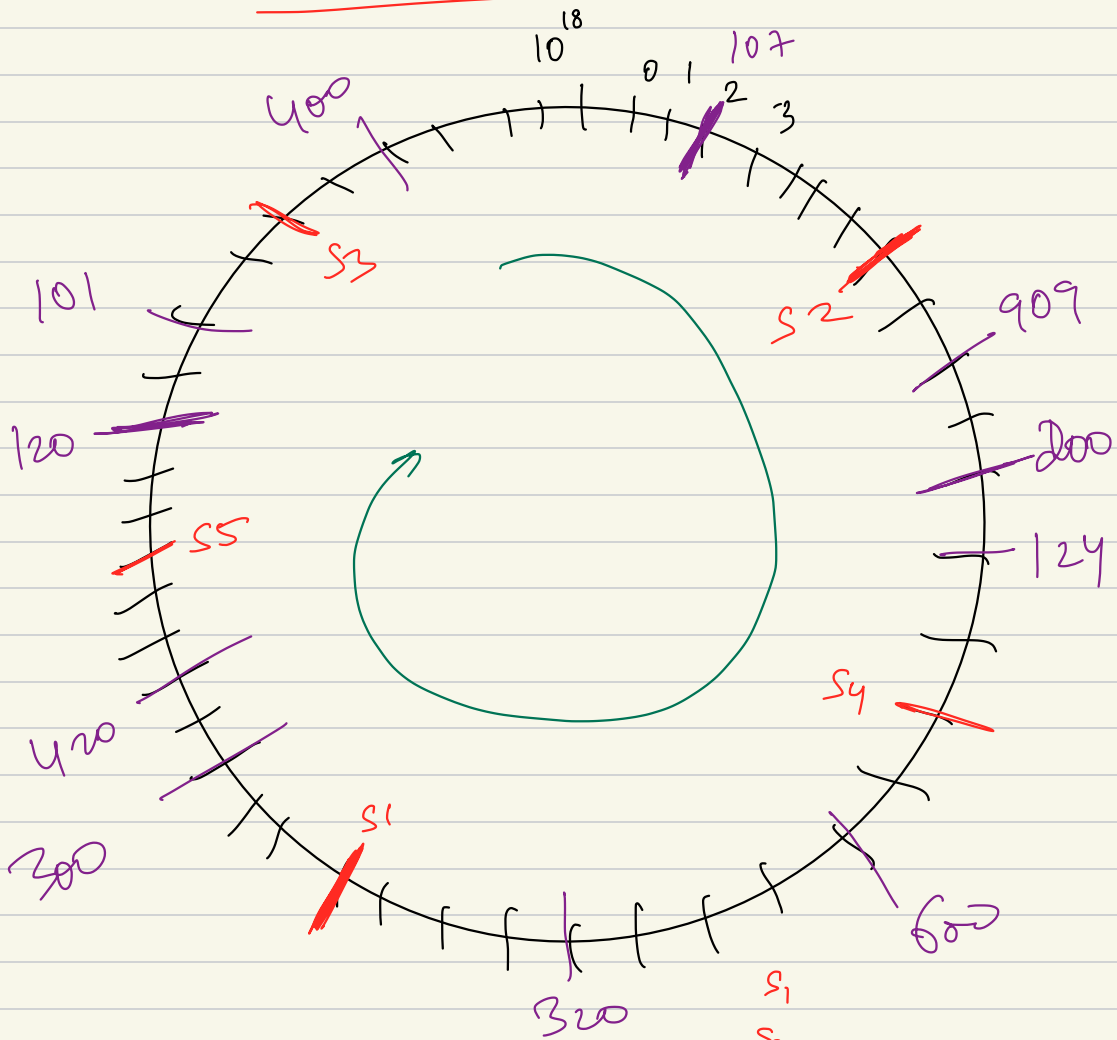
11

a

Hashing at first seemed a promising approach
but on closer look it becomes clear that

Simple hashing will NOT be able
to optimally handle increase/decrease in
machines 😊

CONSISTENT HASHING



5 servers

s_1
 s_2
 s_3
 s_4
 s_5

$$H_S(\text{server-id}) = \left[\text{Output} \right]_{\underline{0 - 10^{18}}}$$

$$H_S(s_1) = 10^{14}$$

$$H_S(s_2) = 100$$

$$H_S(s_3)$$

$$H_S(s_4)$$

$$H_S(s_5)$$

$$H_R(\text{request-id}) = \left[0 - 10^{18} \right]$$

\swarrow
 sharding-id

$$T_1 \quad \underline{107}$$

$$T_2 \quad 120$$

$$T_3 \quad 200$$

$$T_{100} \quad 107$$

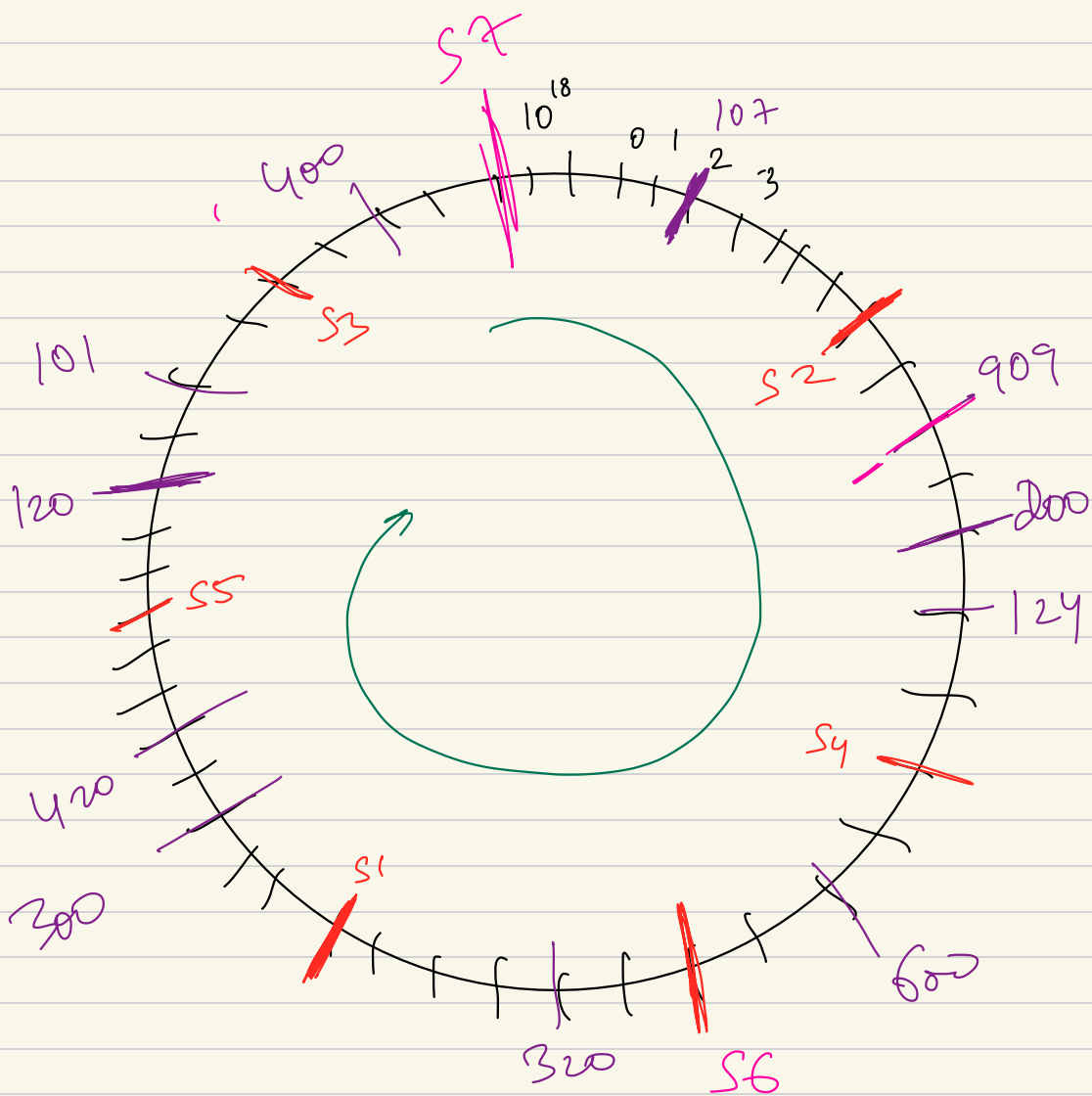
$$H_R(\text{request-id}) = \text{value}$$

\downarrow
 request-id

$$H_R(120) = \text{value 2}$$

$$H_R(200) = \text{value 3}$$

$$H_R(107) = \text{value 1}$$

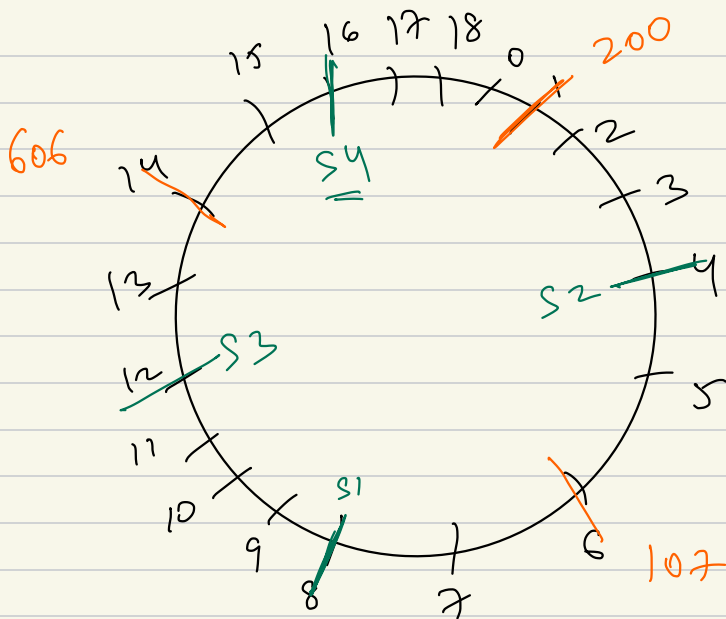


S6 gets added

$$\underline{\underline{H_5}}(S_6) = \text{value}$$

$$H_5 \equiv \left[7x^2 + 19x + 104x^3 - x^{24} \right] \cdot 10^{18}$$

$$\underline{\underline{H_R}} \equiv \left[\quad \quad \quad \right] \cdot 10^{18}$$



$$H_s(\text{server 1}) = 8$$

$$H_s(\text{server 2}) = 4$$

$$H_s(\text{server 3}) = 12$$

add S4

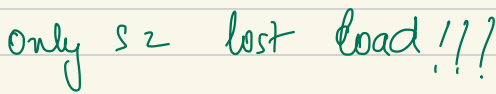
$$H_s(\text{server 4}) = 16$$

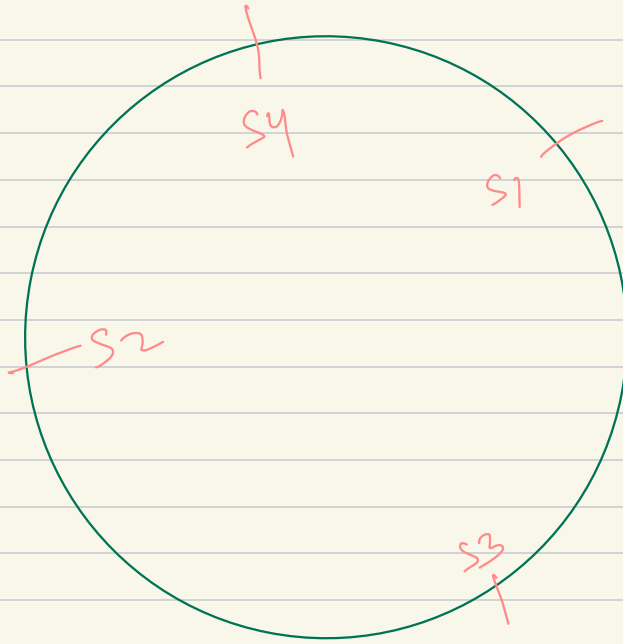
$$H_R(107) = 6$$

$$H_R(200) = 1$$

$$H_R(606) = 14$$

only 606
needs to
move!!!
 $S_2 \rightarrow S_4$





remove S4

S2 same

S3 same

S1 2w1a the load (i)