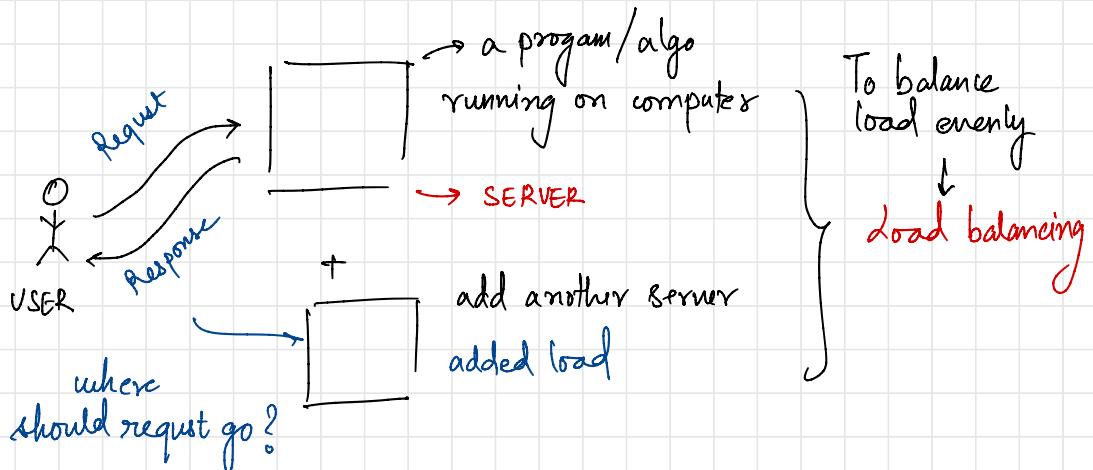


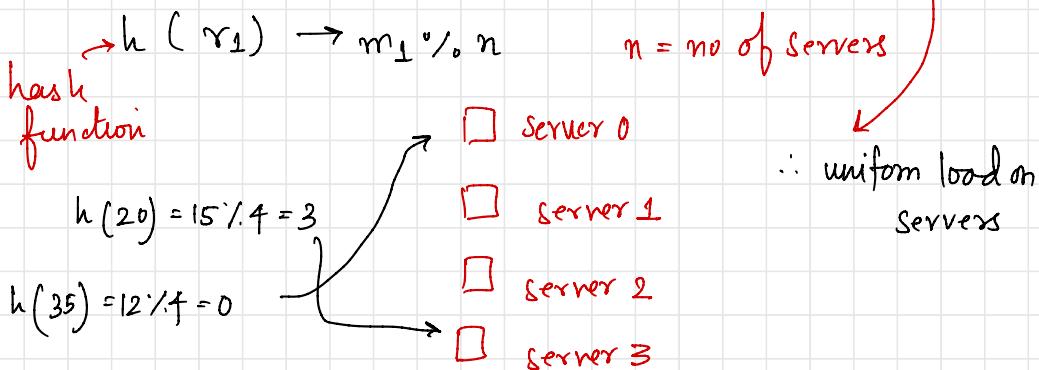
# LOAD BALANCING ?



Concept of consistent Hashing for load Balancing

Request ID for each request (ideal random IDs)

$$0 \rightarrow (M-1)$$



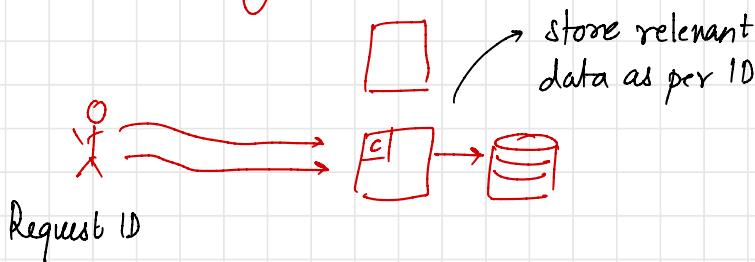
In actual scenario IDs are not random

one ID repeat



sent to same server  
again & again

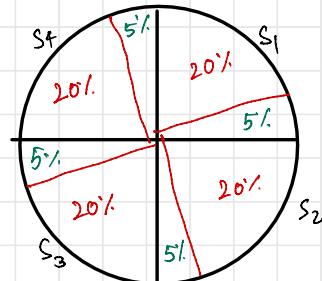
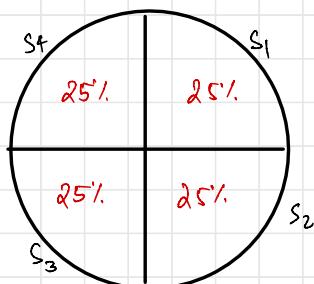
Using a cache memory at server (local data)



Instead of increasing more servers → great change in system

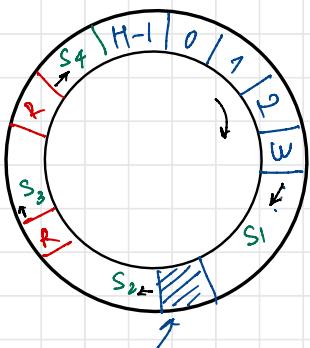
small changes in complete system

more efficient



## CONSISTENT HASHING

Request ID  $\rightarrow h(r_i)$



request point to  
server

S<sub>0</sub> S<sub>1</sub> S<sub>2</sub> S<sub>3</sub> S<sub>f</sub>

hashing server IDs

$h(0) \mod M$

M = Search Space

Going clockwise finding nearest server  
server next in clockwise direction  
will serve request

Expected Load factor =  $\frac{1}{\text{no. of servers}}$  (ideal)

Why use this architecture over simple hash?

- + Adding a new server will not change system load drastically
- \* No issues for adding new servers.

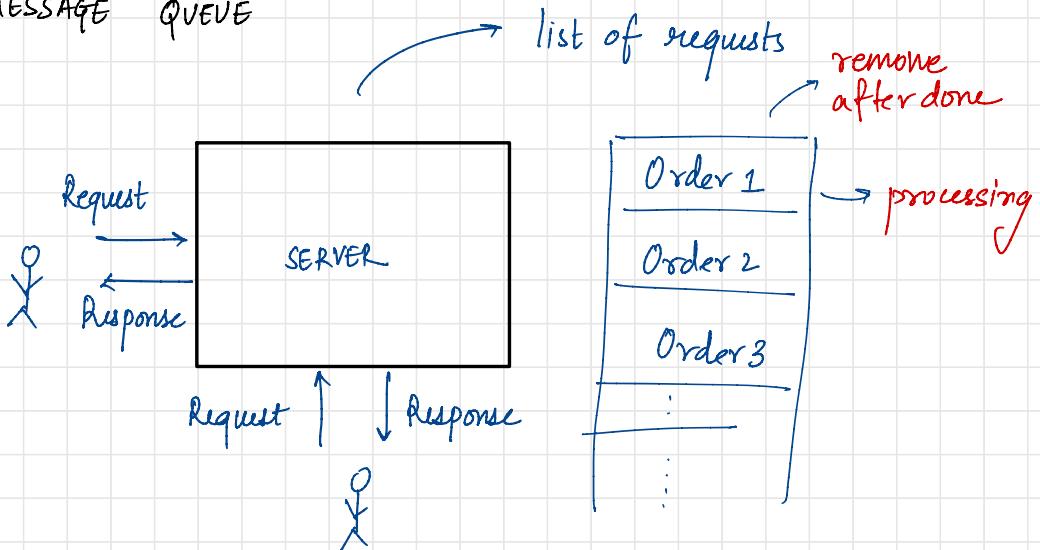
Disadvantage?

- Skewed load balancing in practical

Solution ; Making Virtual Servers

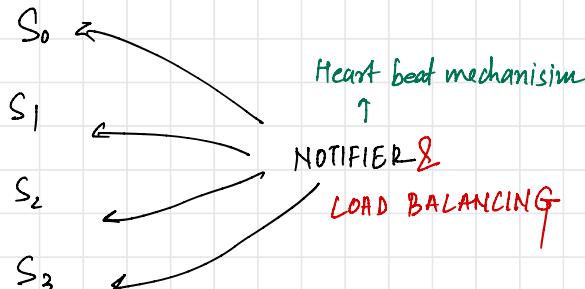
- ↳ using multiple hash functions
- ↳ k hash functions

## MESSAGE QUEUE



## ASYNCHRONOUS PROCESSING

\* In case of fault other servers taking requests



ID	TYPE REQUEST	DONE?
1	P1	Y
2	P2	N
3	P3	N



IF FAULT with server it will stop responding to notifier

DIVIDE NOT DONE REQUESTS → USING LOAD BALANCER  
to ensure no DUPLICACY

## MESSAGE / TASK QUEUE

Take request → assign to server → If server takes long  
(assuming fault)  
↓  
Send to next server

Examples of Task Queues (Encapsulating  
complexity on server side)

Rabbit MQ

Zero MQ

Java Messaging Service