# Scaling

| Horizontal scaling | Vertical scaling |
|---|---|
| □ □ □ □ □ | □ |
| → Load balancing required | → N/A |
| → Resilient | → single point of failure |
| → network calls (RPC) | → Interprocess Communicator |
| → data inconsistency | → consistent |
| → scales well as users increase | → hardware limit |

# Load Balancing

→ N-servers have load on them. Balancing load is known as load balancing.
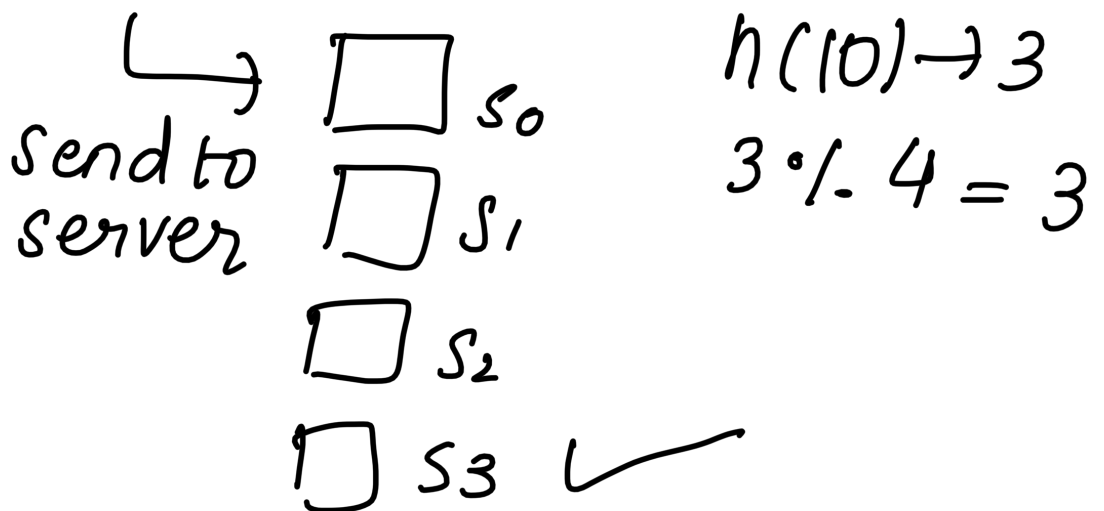
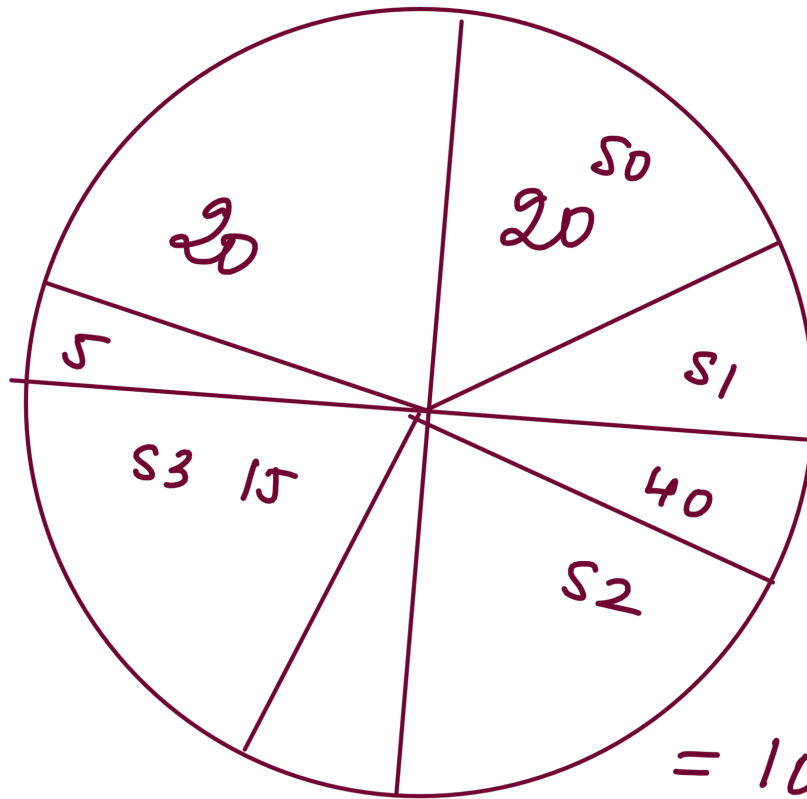To do that : – consistent hashing

↓

evenly distribute load

Request Id → 0 to M-1 → sent to server

Take the ID, say $r_1$ then hash it

$h(r_1) \to m_1 \% n \to$ server

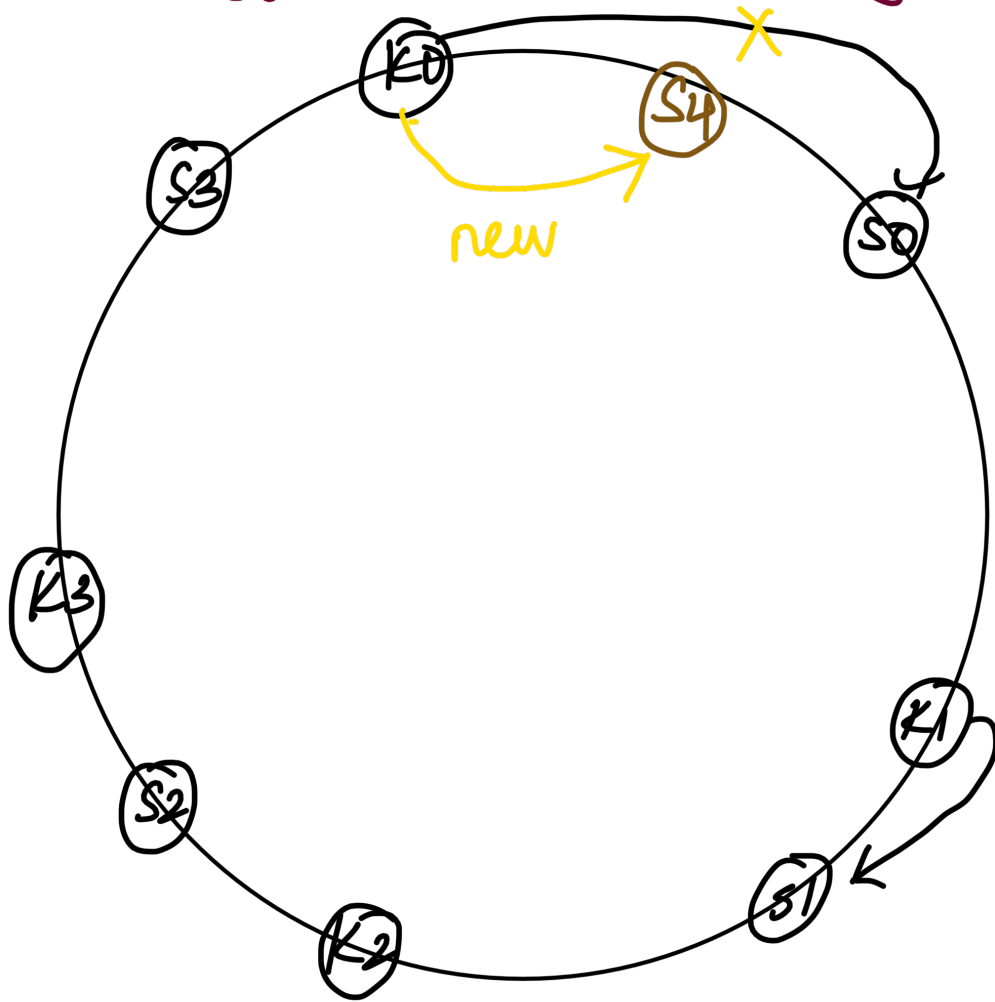Send to server □ $S_0$

□ $S_1$

□ $S_2$

□ $S_3$ ✓

$h(10) \to 3$

$3 \% 4 = 3$

on increasing servers, hash generated server number changes. → consistent hashing

S0 20

S1

40

S2

S3 15

5

20

$+5 + 5 + 10$
$+10 + 15$
$+15 + 20$
$+20$
$\downarrow$
$S4$

$= 100 = M$

# Consistent Hashing



In simple hashing, when new server added, almost all keys need to be remapped.

With consistent hashing, on adding new servers, only fraction of keys are relocated