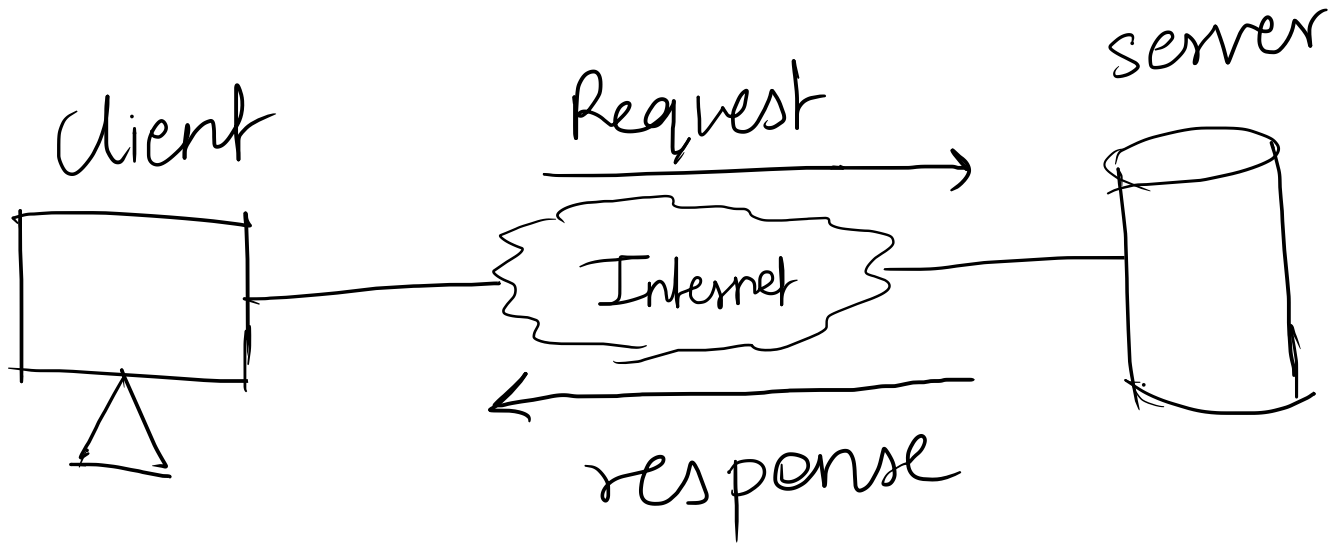


# Client - Server Model

Client : Service requester

Server : Service provider



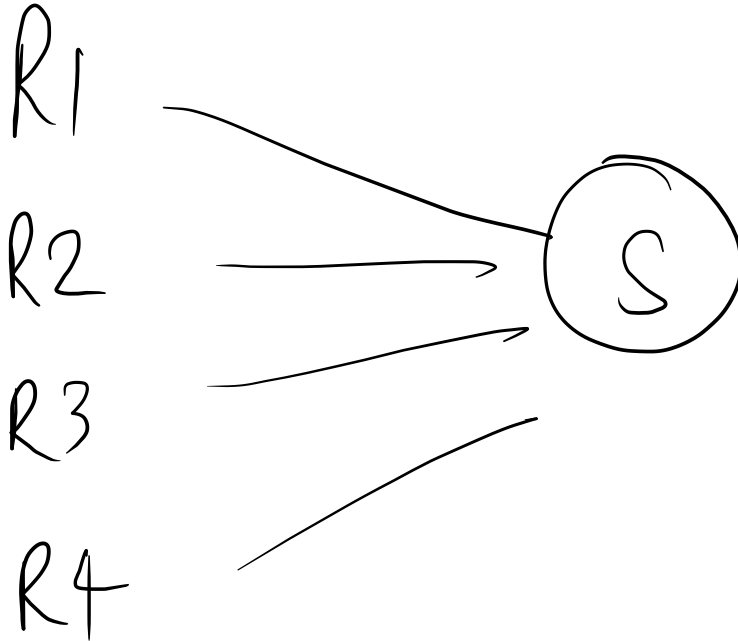
Client - Server model

# DNS (Domain name system)

- User enters URL
- DNS server lookup for the address
- HTTP / HTTPS request
- Server sends the files of website

Requests

Server



# Popular Website

Requests ↑

(heavy traffic)

Server

(load ↑)

# Application Scaling

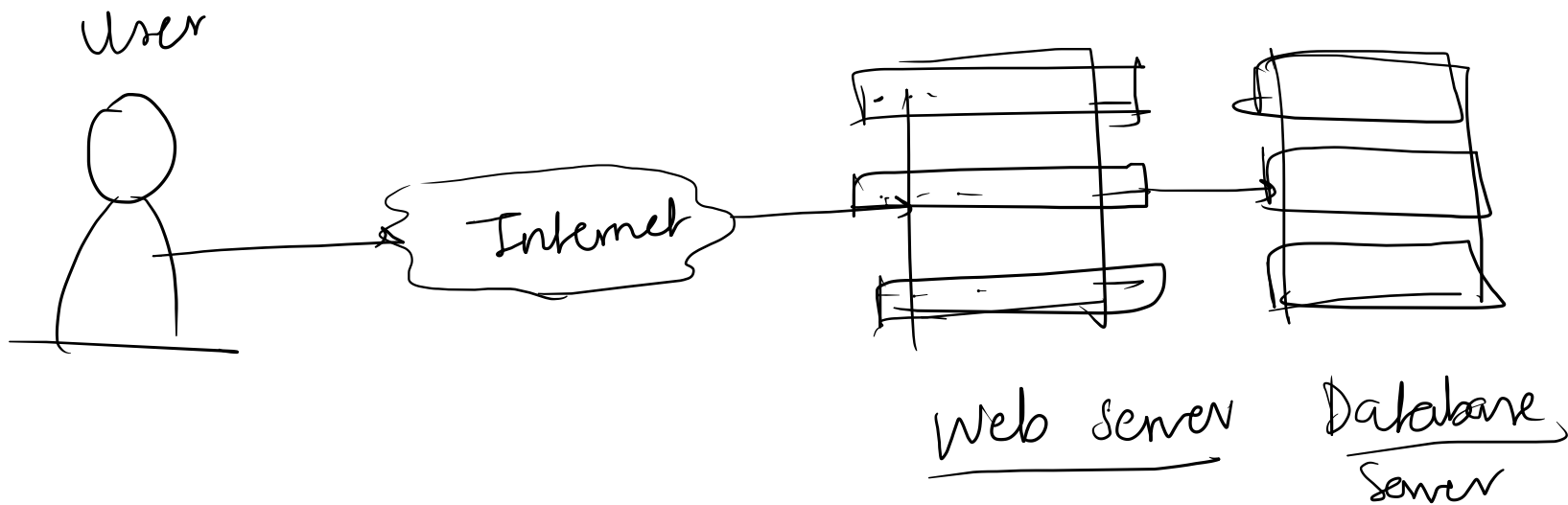
```
graph TD; A[Application Scaling] --> B[Horizontal Scaling]; A --> C[vertical Scaling]; B --> D["↑ no. of servers"]; C --> E["↑ size of server"]
```

Horizontal Scaling

↑ no. of servers

vertical Scaling

↑ size of server



# Problems

## ① Single Point of Failure

If server goes down, service unavailable

- Bad experience for users
- Unacceptable for service providers



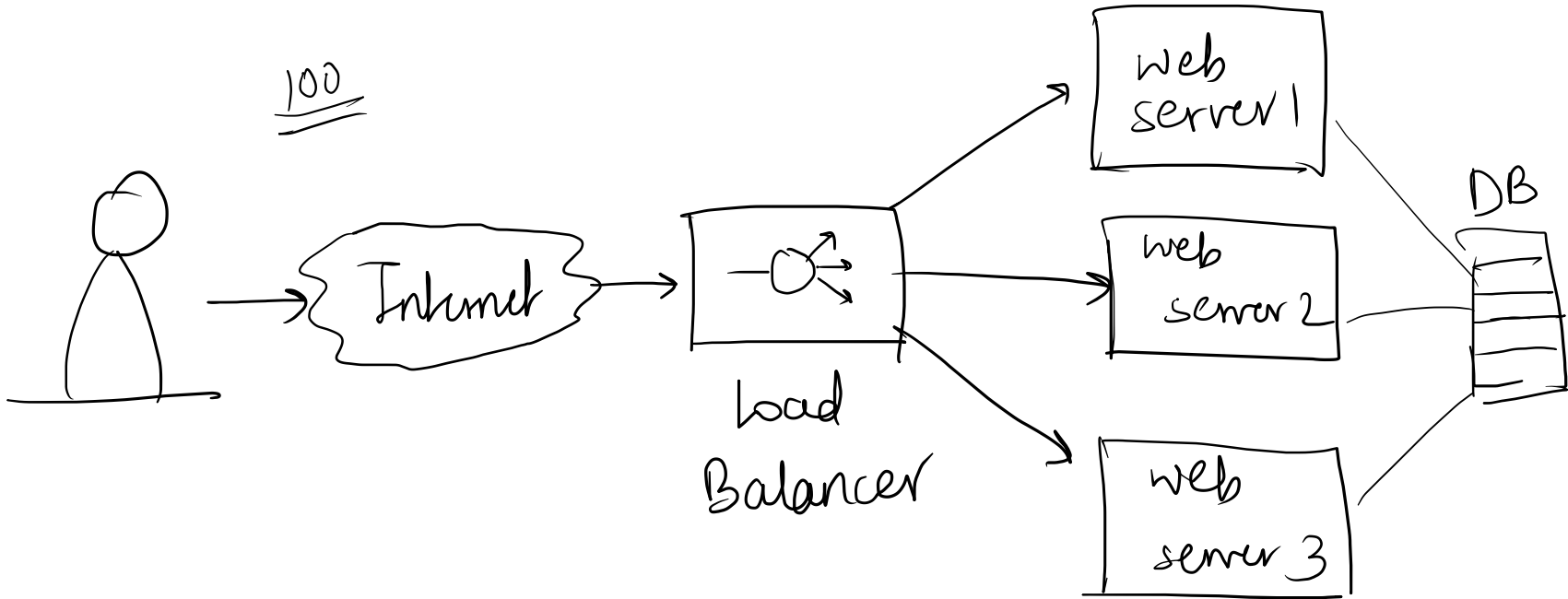
## ② Overloaded servers

→ limitation of no. of requests  
that a server can handle

→ Requests ↑ server overloaded

# LOAD BALANCER

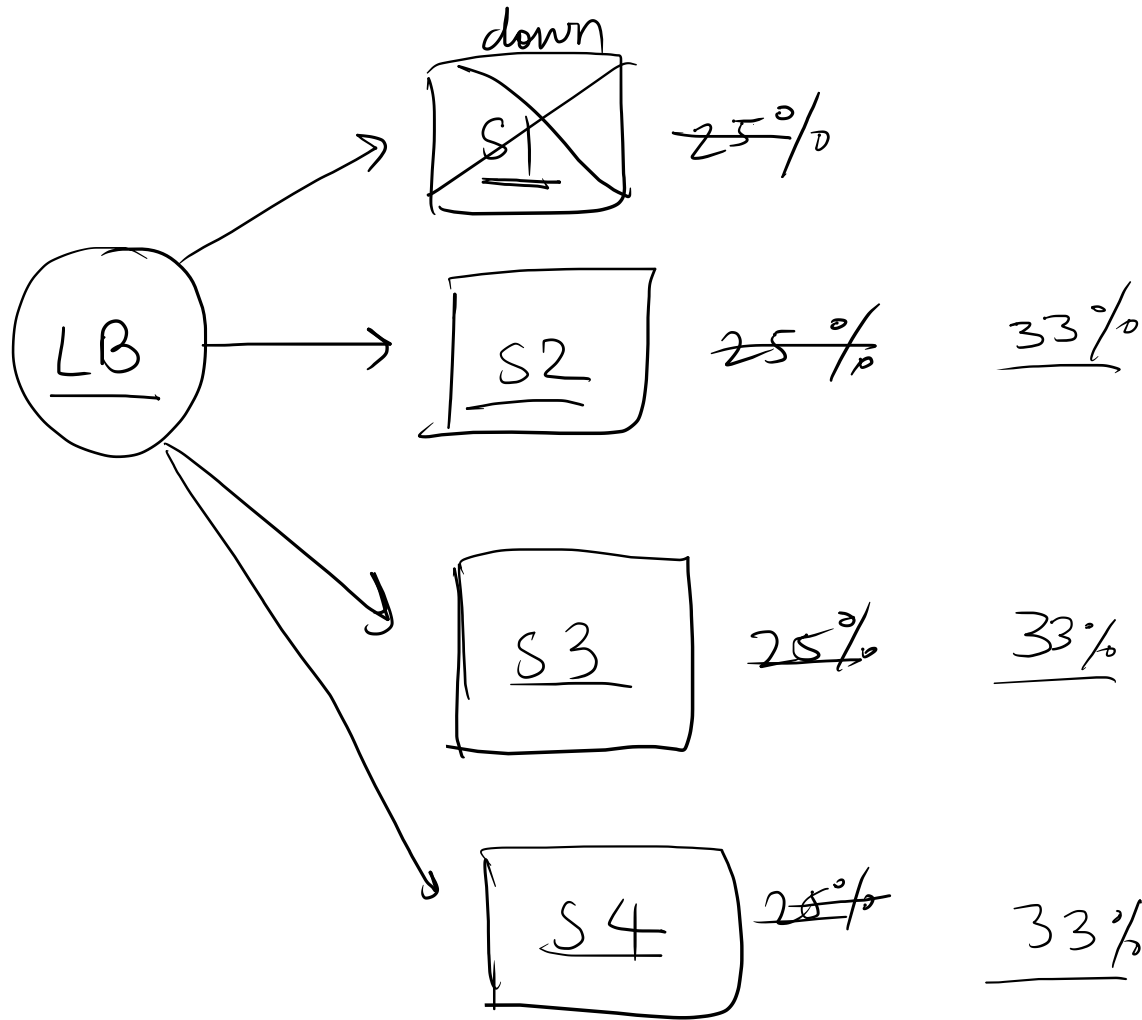
---



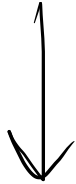
## Solution :

- 1.) Add a L.B. in front of servers
- 2.) Handle multiple requests (heavy volume)
- 3.) Share the requests across multiple servers
- 4.) Even if one of the servers go down, service will continue.
- 5.) faster response time, low latency

- L.B. minimize response time
- maximise throughput
- High availability and reliability
- continuous health checks to monitor servers capability of handling request



# Horizontal Scaling

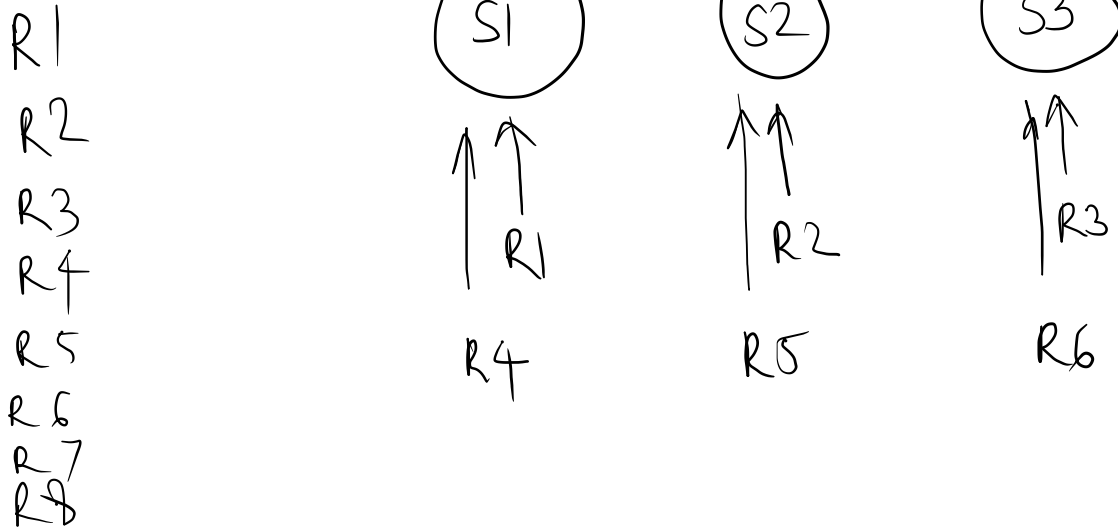


## Load Balancer ✓

# Load Balancing Algorithms

## ① Round Robin Algorithm

Requests are distributed in a sequential/rotational manner



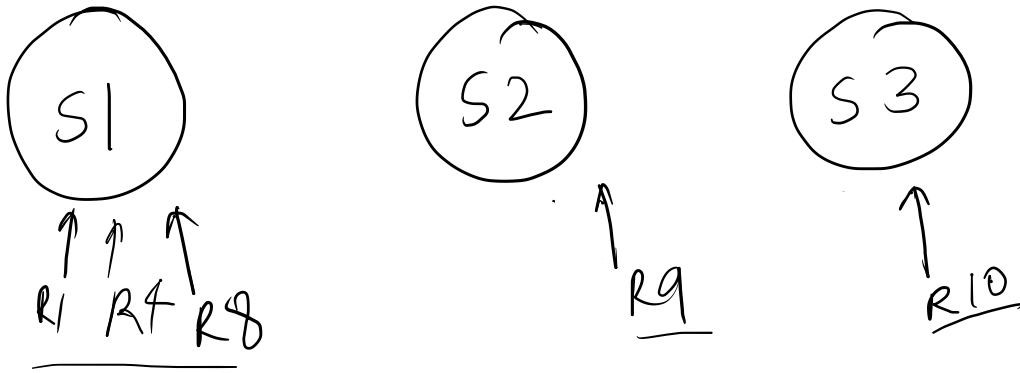
$$\frac{R_i \% N}{}$$

## Round Robin limitation:

- doesn't consider the load already on the server
- risk of single server overloading



## ② Least Connection Method



- \* Requests directed to the server with fewest no. of requests / active connections
- \* expensive, as L.B. needs to compute, to identify the server with least no. of requests

