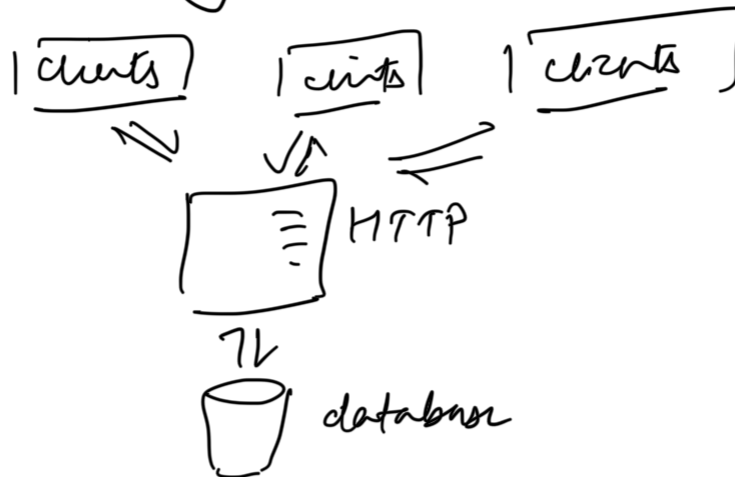# System Design Overview

## Horizontal vs Vertical scaling

→ Single server system
- is a single point of failure
- budget friendly and still has a place
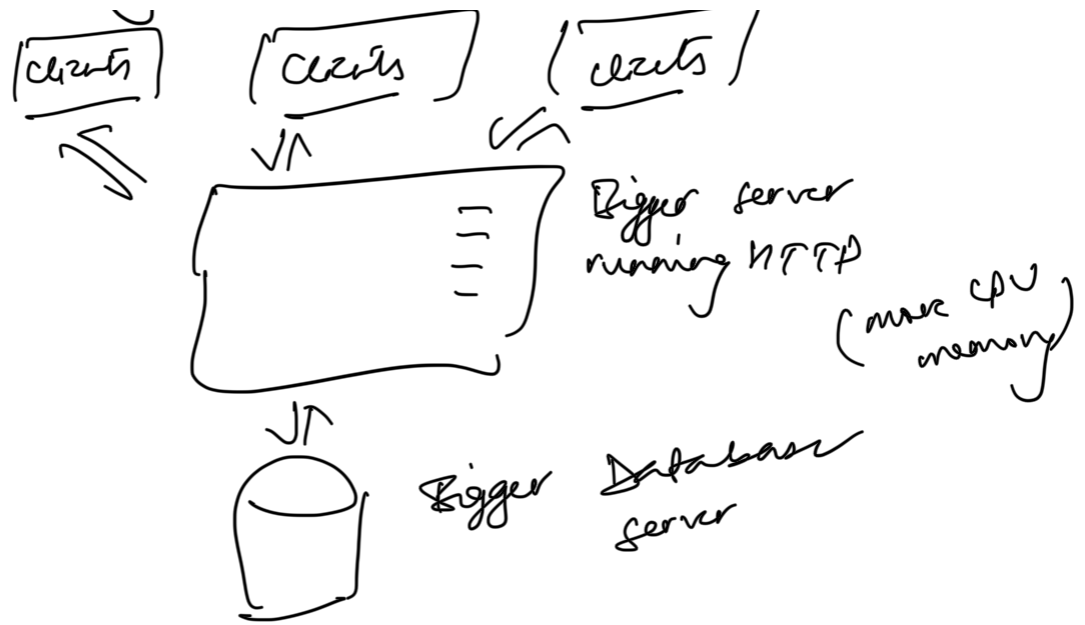
but less tolerant of failures


HTTP Server
Database

let us make the design a little bit more
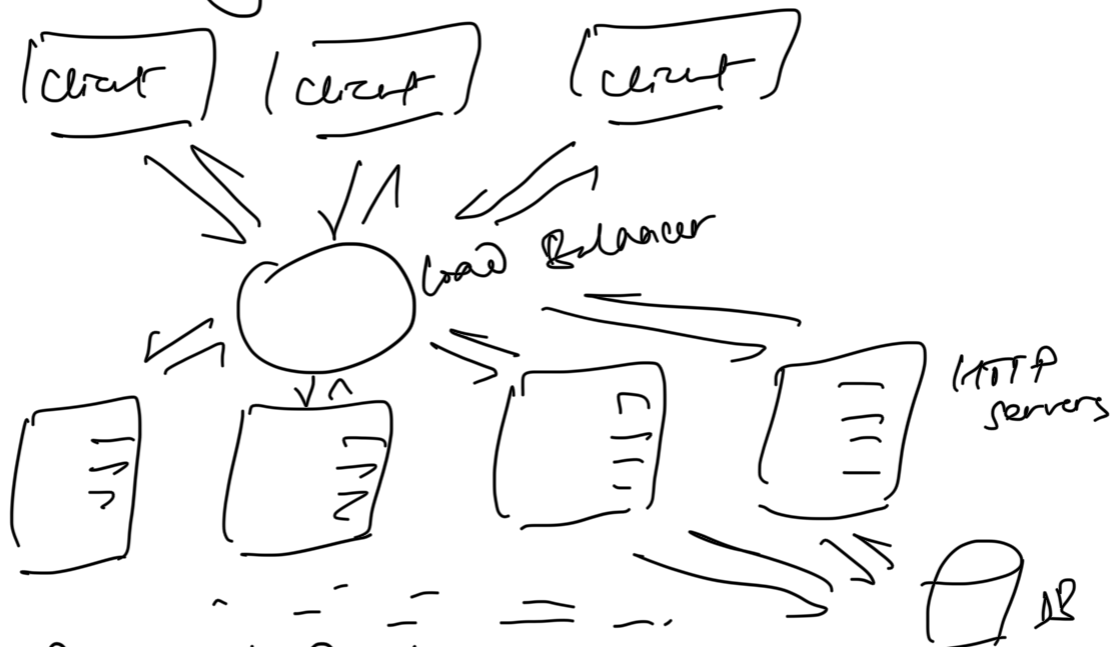tolerant so that HTTP server and Database
are hosted separately


clients  clients  clients
HTTP
database

little bit better because we can scale independently

(bigger)
Vertical Scaling: use beefier machine (CPU↑, mem↑)

clients     clients     clients

Bigger server running HTTP
(more CPU memory)

Bigger Database server

- Still each machine is not tolerant to failures. If it goes down, service goes down
- Each machine can only be scaled so much ie there is a max of CPU and max of mem that you cannot scale beyond

Horizontal Scaling: Add more servers



client     client     client

Load Balancer

HTTP servers

- needs a load balancer
- fault tolerant : one server goes down, request can be served still by other available servers
- can scale infinitely
- Downside : now we have more stuff to maintain

- Horizontal scaling is easier if the servers are **Stateless**, meaning any subsequent request coming into the server should not rely on the server's state. (changes can be and will be made in database).
  - This means any next incoming request should be independent of earlier request
  - Because if they are dependent (and not reflected on the database), then there is no way to be certain which HTTP server the earlier request was served by and which server the next request will be routed to
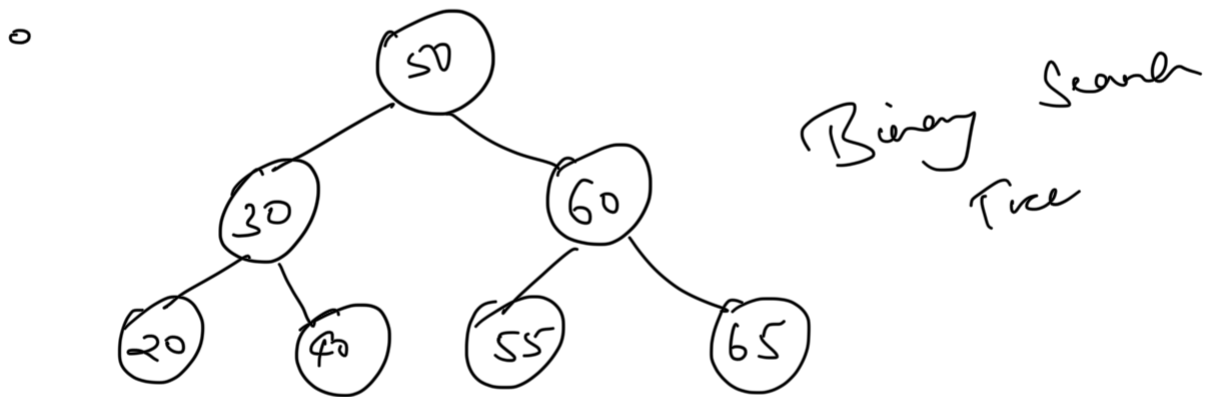
# Binary Search Tree

- Binary Tree: Each member has a left and a right child. Each node has a most 2 children.
- Binary Search Tree: Each left is less than

~right child and parent. And right
child is greater than parent.
- Rule applies to the entire tree, meaning
any value to the left must be less
and any value to the right must be higher

- Access in a BST is $O(\log n)$
  worst case is $O(n)$

○

```
            50
          /    \
        30      60
       /  \    /  \
     20   40  55   65
```

Binary Search
Tree

# Fail Over Strategies : Database

→ Servers can be
  • hosted within your own company
  • using cloud servers
  • using fully managed "Serverless" services

→ How to scale DB
    → using "cold stand by server"
        • create a periodic backup in other
  storage
        ○ upsides :
          ○ the one

- ...
- Downside
  - possible data loss
  - downtime could be significant
→ using "warm standby server"
  - create another database server that is a replication, and is always ready to handle request
  - many databases allow replication. As a data is written/updated/deleted, operation is duplicated in both the places; may be after a delay
  - there could still be a tiny window where you could loose data; but still a better option
  - if the primary goes down, request will be directed to the backup/replication
- How is this handled in practise??

→ using "Hot standby approach"
  - replication in real time
  - thus it's hot

→ using "Multi Primary"
  - writing simultaneously both both servers rather that relying on replication

# Horizontal Scaling : Sharding