

Week 1 - Lesson 3 - Scalability: vertical vs horizontal

- 1 What does "scaling" really mean?
- Scaling = handling more load without breaking the system

Load can be:

- more users
- more requests
- more data

2 Vertical scaling (Scale Up)

What it is

- make one machine stronger
- more CPU, RAM, Disk

Example

- Upgrade from 8GB RAM → 64GB RAM

Pros

- Simple
- No code changes

Cons (very important)

- Hardware limits
- Very expensive
- Single point of failure

If this machine dies → system dies

3 Horizontal Scaling (Scale Out)

What it is

- Add more machines
- Distribute the load

Example

- 1 server \rightarrow 10 servers \rightarrow 100 servers

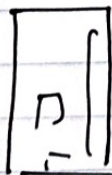
Pros

- Practically unlimited scale
- Fault tolerance
- cheaper

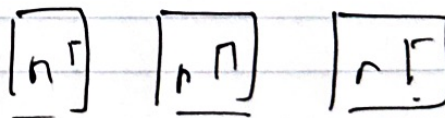
Cons

- Coordination complexity
- Data consistency issues
- Needs load balancing

4 Visual comparison



vertical



horizontal

5 why big systems prefer horizontal scaling

Let's say:

- One server handles 1000 req/sec

Then:

- 10 servers \rightarrow 10,000 req/sec
- One server fails \rightarrow system still works

This is how:

- Google
- Netflix
- Instagram

state

6 The hidden problem of horizontal scaling
This is where system design becomes interesting

Problem 1: State

If user logs in on server A:

- what if next request goes to server B?

\rightarrow State must move out of server

Problem 2: Data

If 10 servers write to same DB:

- DB becomes bottleneck

\rightarrow DB must also scale

7 Important rule

Servers must be stateless to scale horizontally.

We'll revisit this many times

8 Mini Design Exercise

Task 3

you have:

- 5 application servers
- 1 database

Traffic increase 10x

Answer

1. What breaks first?
2. Why can't you just add more app servers?
3. What does this tell you about?