

Chapter 1: Scaling from Zero to Millions of Users – Interview Guide

This chapter explains how to **design a system from scratch** and gradually scale it to **millions of users**. It covers **every concept**, explaining problems, solutions, trade-offs, and best practices.

Problem Statement: Why Do We Need Scaling?

Imagine you launch a **small website** that serves **100 users per day**. Initially, everything runs smoothly on a **single server**:

✅ The **web app, database, caching, and file storage** all operate **from one machine**.

But what happens when your website grows to **millions of users**?

- ❌ **Increased load slows down responses.**
- ❌ **Database queries take longer.**
- ❌ **Server crashes when overloaded.**
- ❌ **Users experience downtime.**

To **prevent failures** and ensure **high availability**, we need to **scale the system step by step**.

1. Starting Simple – Single Server Setup

What Is It?

A **single server** handles everything:

- **Frontend + Backend** (handles user requests).
- **Database** (stores user accounts, posts, transactions).
- **File Storage** (stores images, videos, documents).
- **Cache** (stores frequently used data).

Users send **HTTP requests**, and the server **processes and returns responses**.

❌ Problems (Why This Fails with Growth?)

- ❶ **Performance Bottleneck:** CPU, RAM, and disk I/O have limits.
- ❷ **Single Point of Failure:** If the server crashes, everything stops.
- ❸ **Limited Scalability:** Cannot handle a sudden increase in traffic.
- ❹ **Slow Response Times:** More data = slower database queries.

First Scaling Step: Separate Web & Database Servers

- ✓ **Web Server:** Handles user requests.
- ✓ **Database Server:** Stores and retrieves data separately.

💡 **Why? Web & DB can now scale independently.**

2. Choosing the Right Database – SQL vs NoSQL

The **database** is a **major bottleneck** when scaling. You must decide between **SQL (Relational DB)** and **NoSQL (Non-Relational DB)**.

SQL (Relational Databases)

- ✓ **Structured data (tables, rows, columns).**
- ✓ **Supports complex queries (joins, ACID transactions).**
- ✓ **Good for financial, banking, and e-commerce applications.**
- ✗ **Doesn't scale well horizontally – requires vertical scaling.**

Examples: MySQL, PostgreSQL, Oracle.

NoSQL (Non-Relational Databases)

- ✓ **Designed for large-scale applications.**
- ✓ **Scales horizontally (adds more servers when needed).**
- ✓ **Great for real-time apps (social media, analytics, logs).**
- ✗ **Weaker consistency (eventual consistency).**

Examples: MongoDB (document-based), Redis (key-value), Cassandra (column-based).

💡 **When to use NoSQL?**

- **High write traffic** (e.g., logging, real-time analytics).
 - **Flexible data models** (e.g., social media posts).
-

3. Scaling Strategies – Vertical vs Horizontal

♦ Vertical Scaling (Scaling Up)

- **Upgrade the existing server** (add more RAM, CPU, SSD).
- **Pros:** Simple, no code changes needed.
- **Cons:** Expensive, hardware limits, single point of failure.

♦ Horizontal Scaling (Scaling Out)


- **Add more machines** to handle requests.
- **Requires load balancer & distributed data management.**
- **Pros:** Fault tolerance, cheaper, scalable.
- **Cons:** Requires architectural changes.

 **Best Practice:** Horizontal scaling is preferred for large-scale systems.

4. Load Balancer – Distribute Traffic Across Multiple Servers

Why Use a Load Balancer?

Without a load balancer:

 All users hit **one web server**, which slows down and crashes.

With a load balancer:

- ✓ **Traffic is evenly distributed** across multiple servers.
- ✓ **Prevents overloading of any single server.**
- ✓ **Ensures high availability** (if one server fails, others take over).

How It Works?

- 1 Users request a website (e.g., **mysite.com**).
- 2 DNS resolves the domain to a load balancer's IP.
- 3 The load balancer forwards the request to one of the web servers.

Load Balancing Strategies

- **Round Robin:** Requests are sent to servers in rotation.
- **Least Connections:** Directs traffic to the server with the fewest active connections.

- **IP Hashing:** Routes users to the same server based on their IP.

Popular Load Balancers:

- **Nginx** (most used for web applications).
 - **AWS ELB** (Elastic Load Balancer for cloud).
 - **HAProxy** (high-performance open-source LB).
-

5. Database Replication – Handling High Read Traffic

Problem: Single Database Becomes a Bottleneck

- More users = More queries = Slower performance.

Solution: Master-Slave Replication

- **Master DB** → Handles **writes** (INSERT, UPDATE, DELETE).
- **Slave DBs** → Handle **reads** (SELECT queries).
- Load balancer directs read queries to slave databases.

Pros:


- Improves read performance by distributing load.
- Ensures redundancy (if master crashes, a slave can take over).


Cons:

- **Replication Lag** (slaves might not have the latest data).
 - **Writes are still limited** to one master.
-

6. Caching – Reduce Database Load & Speed Up Responses

Why Use Cache?

 Database queries are slow – fetching the same data repeatedly wastes resources.

 Cache stores frequently accessed data in memory, reducing DB load.

Types of Caching

① Application Cache (Local Cache)

- Stores data in **RAM of the web server**.
- **Pros:** Fastest, low latency.
- **Cons:** Not shared between multiple servers.

② Database Cache (Redis, Memcached)

- **Stores query results** for quick retrieval.
- **Pros:** Fast reads, reduces DB load.
- **Cons:** Needs cache invalidation strategy.

③ Content Cache (CDN – Content Delivery Network)


- **Caches static files (JS, CSS, images, videos, etc.).**
- **Pros:** Reduces bandwidth usage & speeds up content delivery.

Cache Strategies

- **Write-through:** Write to cache & DB at the same time.
 - **Write-back:** Write to cache first, sync to DB later (fast but risky).
 - **Eviction Policies:**
 - **LRU (Least Recently Used)** – removes the least accessed items.
 - **LFU (Least Frequently Used)** – removes the least requested items.
-

7. Content Delivery Network (CDN) – Faster Static Content Delivery

Why Use a CDN?

-  Users far from the server experience slow load times.
-  CDN stores static files across multiple locations worldwide.

How It Works?

- ① User requests a file (image, CSS, JS).
- ② CDN serves it from the nearest location.
- ③ If not in cache, CDN fetches it from the origin server and stores it.

Popular CDNs: Cloudflare, AWS CloudFront, Akamai.



Key Takeaways for Interviews

- Start simple, scale step by step.
- Use horizontal scaling (scale-out) over vertical scaling.
- Load balancer + Replicated DB + Cache = High Performance.
- CDN reduces latency for global users.



"Scale out before scaling up" – Always prefer horizontal scaling! 