

# Vertical vs Horizontal Scaling

## Master System Design

Progress 35/130 chapters

Search topics...

API Fundamentals 5/10

Databases & Storage 4/12

Database Scaling Techniques 2/8

Caching 5/6

Asynchronous Communications 3/4

Tradeoffs 0/9

Vertical vs Horizontal Scaling

Concurrency vs Parallelism

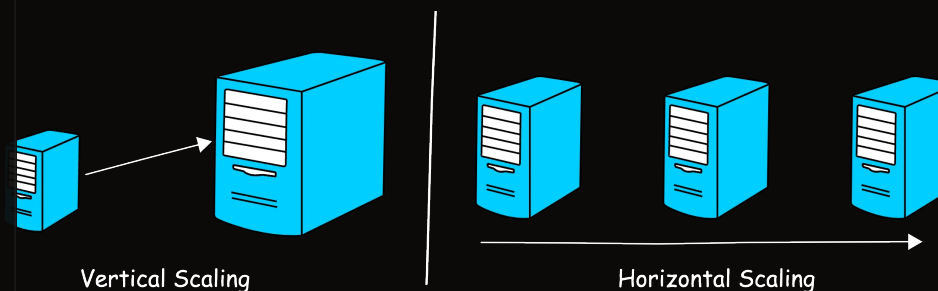
Ashish Pratap Singh



4 min read

When your application gets bigger, it needs more resources.

To handle this growth, two common approaches are **vertical** and **horizontal** scaling.



In this article, we will explore the pros and cons of these two scaling strategies and discuss when to use one over the other.

## 1. Vertical Scaling (Scaling Up)

★ Get Premium  
Subscribe to unlock full access to all premium content

Subscribe Now

Reading Progress 0%

### On this page

1. Vertical Scaling (Scaling Up)
2. Horizontal Scaling (Scaling Out)
3. When to Choose Vertical vs Horizontal scaling...
4. Combining Vertical and Horizontal Scaling

- Long Polling vs WebSockets
- Stateful vs Stateless Architecture
- Strong vs Eventual Consistency
- Push vs Pull Architecture
- Monolith vs Microservices
- Synchronous vs Asynchronous Communications
- REST vs GraphQL



Distributed System Concepts

0/10 ▾



Architectural Patterns

0/5 ▾



Microservices

0/6 ▾

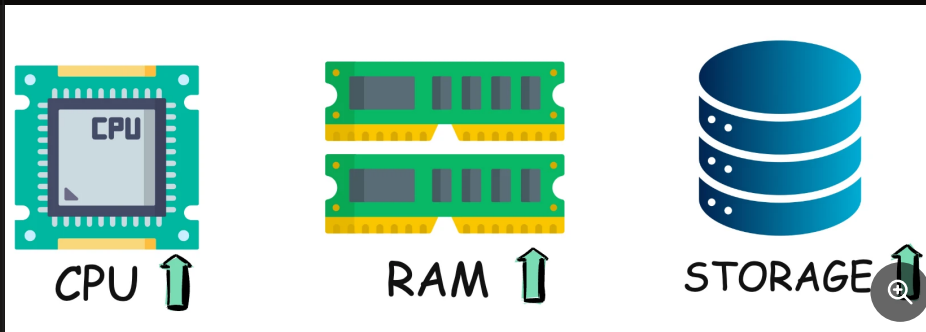


Big Data Processing

0/5 ▾

**Vertical scaling**, also known as "**scaling up**" involves boosting the power of an existing machine within your system to handle increased loads.

This can mean upgrading the **CPU**, **RAM**, **Storage**, or other hardware components to boost the server's capacity.



- **Upgrading CPU:** Replacing your server's processor with a more powerful one.
- **Increasing RAM:** Adding more memory to handle larger datasets and reduce reliance on slower storage.
- **Enhancing Storage:** Switching to faster storage (like SSDs) or increasing overall storage capacity.

## Pros of Vertical Scaling

1. **Simplicity:** Vertical scaling is relatively simple to implement as it doesn't require changes to the application architecture.
2. **Lower latency:** Since all the resources are located on a

single machine, vertical scaling can eliminate the need for inter-server communication thus lowering the latency.

3. **Reduced software costs:** In the initial phase, vertical scaling may be more cost-effective than horizontal scaling, especially when dealing with moderate increases in demand.
4. **No Major Code Changes:** Often requires little to no adjustments to your application's codebase.

### Cons of Vertical Scaling

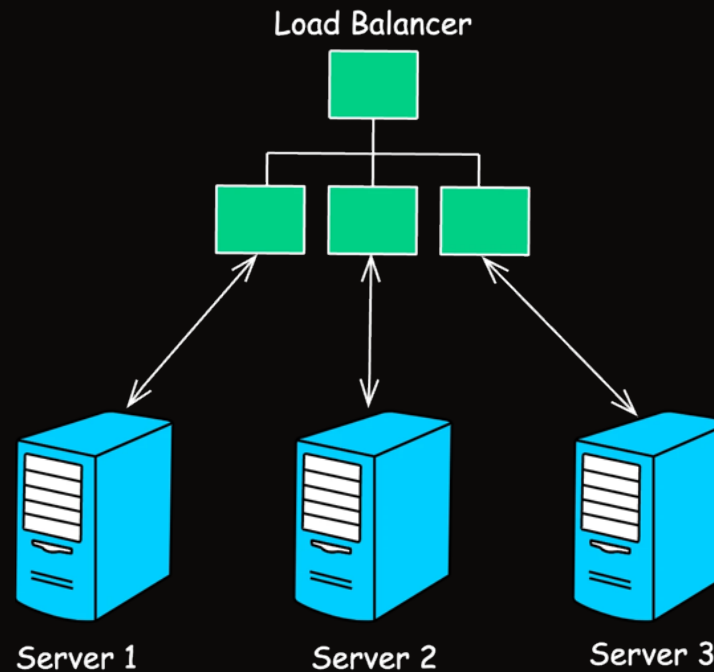
1. **Limited scalability:** There is a limit to how much a single machine can be upgraded.
2. **Single point of failure:** With all resources on one server, any hardware failure can bring down the entire system.
3. **Downtime:** Upgrading hardware often requires taking the server offline, which can be a significant disadvantage.
4. **Higher Costs in the Long Run:** High-end servers with powerful CPUs and large amounts of RAM can get very expensive as you scale.

---

## 2. Horizontal Scaling (Scaling Out)

**Horizontal scaling**, or **scaling out**, involves adding more servers or nodes to the system to distribute the load across multiple machines.

Each server runs a copy of the application, and the load is balanced among them often using a load balancer.



## Pros of Horizontal Scaling

1. **Near-Limitless Scalability:** You can continue to add nodes as long as your architecture supports it, providing the ability to handle larger loads.
2. **Improved fault tolerance:** The failure of one node does not bring down the entire system, minimizing downtime.

3. **Cost-effective:** Horizontal scaling can be more cost-effective as it uses commodity hardware instead of expensive high-end servers.

## Cons of Horizontal Scaling

1. **Complexity:** Distributing the application across multiple servers introduces complexity in terms of data consistency, load balancing, and inter-server communication.
2. **Increased latency:** Communication between servers can introduce additional latency compared to a single machine.
3. **Cost:** Initial setup and maintenance costs can be higher due to the complexity of the infrastructure.
4. **Application Compatibility:** Your application's code might need adjustments to work effectively in a distributed environment.

---

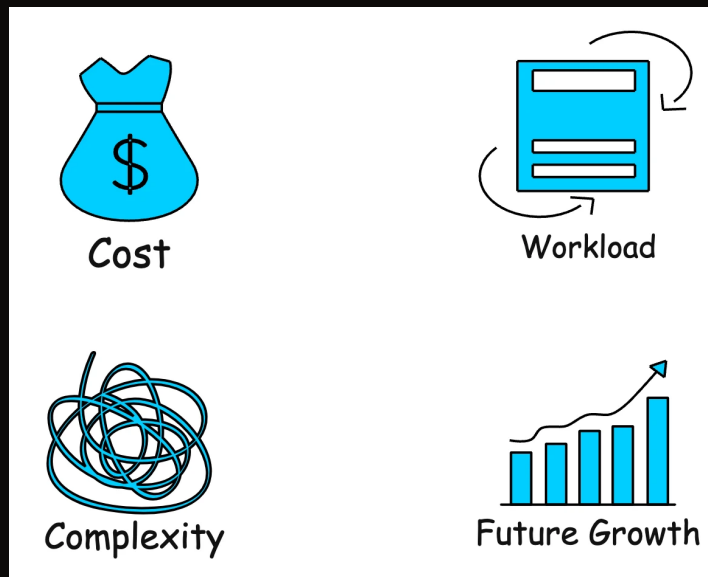
## 3. When to Choose Vertical vs Horizontal scaling

Things to consider to decide between vertical and horizontal scaling:

- **Cost:** Analyze initial hardware costs vs. long-term opera-

tional expenses.

- **Workload:** Is your application CPU bound, memory bound, or does it lend itself to distribution?
- **Architectural Complexity:** Can your application code handle distributed workloads?
- **Future Growth:** How much scaling do you realistically anticipate?



## When to Choose Vertical Scaling

Vertical scaling is a good fit in the following scenarios:

1. **Limited Scalability:** Small to medium-sized applications with a limited growth forecast and your needs are easily met by hardware upgrades.

2. **Legacy applications:** When there is a tight coupling between components, making it difficult to distribute across multiple servers.
3. **Low Latency:** When low latency is a critical requirement, and inter-server communication overhead is unacceptable.
4. **Cost-sensitive projects:** When the budget does not allow for a complex infrastructure and the cost of scaling horizontally outweighs the benefits, such as in the case of expensive software licenses.

## When to Choose Horizontal Scaling

Horizontal scaling is a good fit in the following situations:

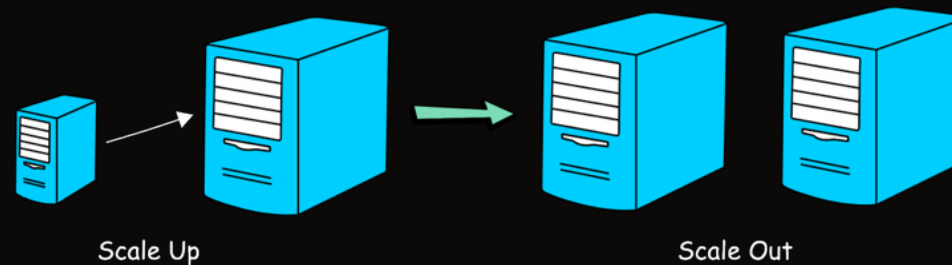
1. **Rapid Growth:** When experiencing rapid growth and requiring the ability to handle increasing traffic.
2. **High availability needs:** When the application needs to be highly available and resilient to node failures.
3. **Easily Distributable:** When the application can be easily distributed across multiple servers without significant modifications.
4. **Microservices architectures:** When applications are designed around microservices, which naturally lend themselves to horizontal scaling.
5. **Cost Effectiveness:** When cost-effectiveness is a priori-

ty, and the use of commodity hardware is preferred.

## 4. Combining Vertical and Horizontal Scaling

In some cases, a combination of vertical and horizontal scaling can be used to optimize system performance and cost-effectiveness.

**Example:** A system can initially scale vertically until it reaches the practical limits of a single machine, and then switch to horizontal scaling to accommodate further growth.



Many successful systems use a combination of both:

- **Vertically Scaled Clusters:** Powerful individual machines form the nodes of a horizontally scaled cluster.
- **Database Sharding:** Data is partitioned across multiple




database servers (horizontal), while each database server might be vertically scaled.


Choosing between vertical and horizontal scaling depends heavily on the specific needs of the application, the expected scale of growth, budget, and how critical uptime is to the business.


Often, the best approach involves a combination of both strategies, starting perhaps with vertical scaling for immediate needs and planning for horizontal scaling as the long-term solution.

[< Prev: Kafka Use Cases](#)

 Take Notes

 Star

 Mark as Complete

 Ask AI

[Next: Concurrency vs Parallelis... >](#)