

# Horizontal Scaling vs Vertical Scaling

Tuesday, October 29, 2024 6:08 PM

## Vertical Scaling

---

--> Vertical scaling, also known as "scaling up," involves adding more power to an existing server.

--> This includes increasing storage, RAM, CPU, and network capacity to enhance the server's overall performance.

--> Unlike horizontal scaling, where multiple servers are added to share the load, vertical scaling focuses on making a single machine more powerful.

### Advantages

---

- **Easier Hardware Upgrades:** Upgrading the existing hardware is simpler than setting up a new server, as you only need to enhance the current machine.
- **Cost-Effective Resource Use:** You pay only for the additional resources you need, avoiding the cost of a completely new setup.
- **Simplified Maintenance:** Since everything runs on a single machine, maintenance and upgrades are generally easier to manage.
- **Better for Applications with High Data Consistency:** For applications requiring strict data consistency, vertical scaling is often preferable as all data is managed in one place, avoiding the complexities of data distribution across multiple servers.

### Disadvantages

---

- **Single Point of Failure:** If the server goes down, all services hosted on it are affected, creating a significant risk for critical applications.
- **Physical Limitations:** There is a maximum capacity for how powerful a single server can be, meaning there's an upper bound on scalability.
- **Expensive High-End Hardware:** Upgrading to top-tier hardware can be costly, and expenses can escalate quickly as you reach the physical limits of vertical scaling.
- **Limited Elasticity for Demand Spikes:** Vertical scaling can handle gradual growth, but it's less suited for sudden traffic spikes or unpredictable demand, as adding resources may require downtime.

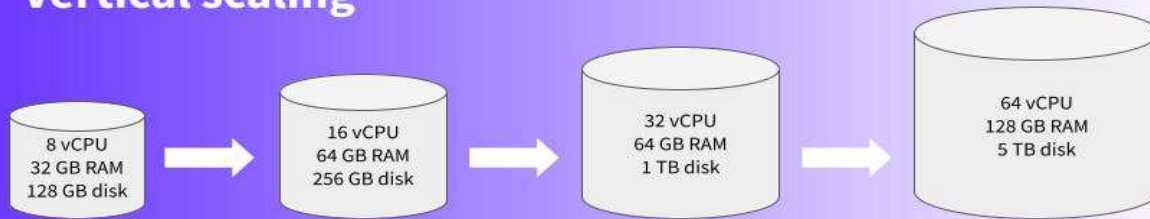
### When to Use Vertical Scaling

---

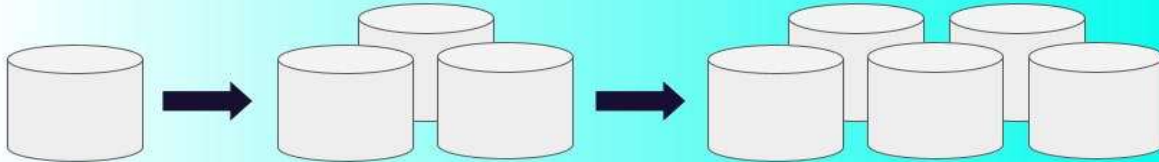
Vertical scaling is ideal in scenarios where:

- **High Data Consistency** is essential, such as in financial systems or databases requiring synchronous transactions.
- **Resource Requirements Are Predictable**, allowing the organization to plan upgrades and allocate budgets effectively.
- **Legacy Systems** are in use that aren't designed for horizontal scaling, making vertical scaling the only viable option.

## Vertical scaling



## Horizontal scaling



### Horizontal Scaling

→ Horizontal scaling, also known as "scaling out," involves adding more servers to your infrastructure to distribute the workload across multiple machines.

→ Unlike vertical scaling, where resources are added to a single server, horizontal scaling allows you to expand capacity by simply adding more servers as needed.

#### Advantages

- **High Availability:** By spreading workloads across multiple servers, horizontal scaling helps ensure that a single point of failure doesn't take down the entire system. If one server goes down, others can continue to operate.
- **Flexible Capacity Growth:** You can add more servers as demand grows, allowing your infrastructure to scale dynamically with minimal downtime.
- **Improved Performance:** Distributing the workload across multiple servers can lead to better performance, especially during peak loads, as different tasks or requests are handled by different servers.
- **Cost Efficiency at Scale:** In some cases, it may be cheaper to add multiple lower-spec servers rather than upgrading a single high-spec server to handle all the workload.

#### Disadvantages

- **Complex Setup and Management:** Managing a distributed system is more complex than managing a single server. It requires skills in load balancing, orchestration, and distributed system design.
- **Data Consistency Challenges:** Maintaining data consistency across multiple servers requires data replication and synchronization mechanisms, which can be complex to implement and may impact performance.
- **Network Latency:** In a horizontally scaled environment, data and requests may need to travel across the network, potentially causing latency issues, especially in applications with high-frequency data synchronization.

#### When to Use Horizontal Scaling

Horizontal scaling is ideal in scenarios where:

- **High Availability** is crucial, such as in e-commerce sites, social media platforms, and other applications requiring minimal downtime.

- **Unpredictable Traffic Patterns** demand flexibility, such as applications experiencing spikes in traffic or seasonal demands.
- **Distributed Data Processing** is needed, as in large databases or big data systems, where different nodes can process data in parallel to speed up tasks.