



AI Model Deployment

Scalable AI Deployment Strategies

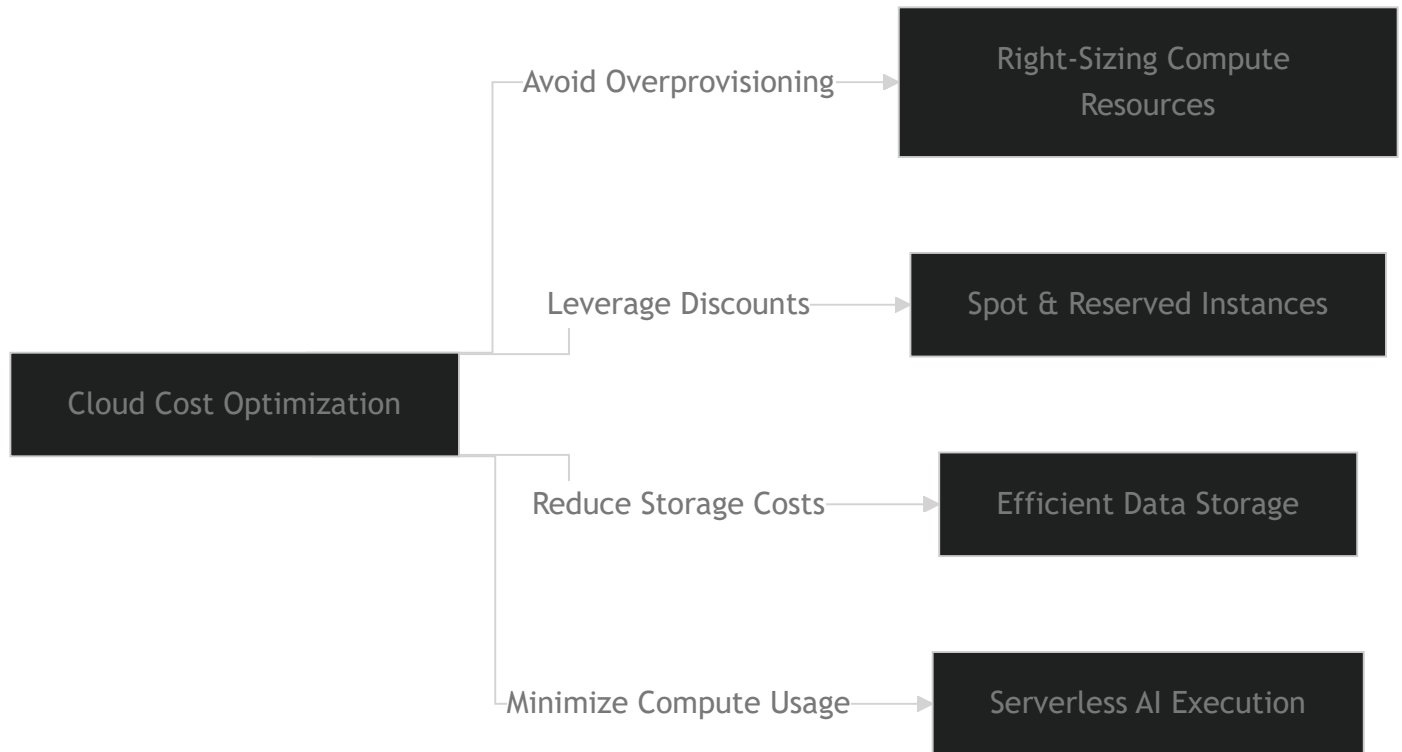
Scalable AI Deployment Strategies

Cost Considerations in AI Deployment

Efficient AI deployment requires balancing cost and performance. Cloud costs are primarily influenced by:

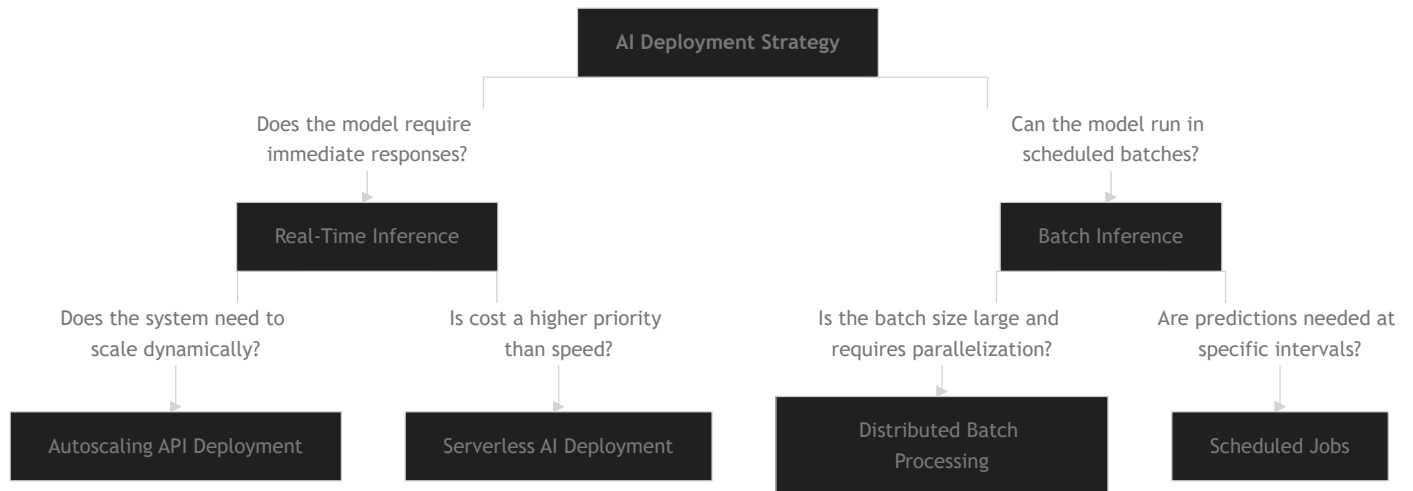
- **Compute:** GPU, TPU, or CPU resources required for inference.
- **Storage:** Model size, data retention, and real-time access needs.
- **Networking:** Data transfer between cloud services and external endpoints.
- **Operational Scaling:** Autoscaling policies, serverless pricing, and reserved instances.

Cloud Cost Optimization Strategies:

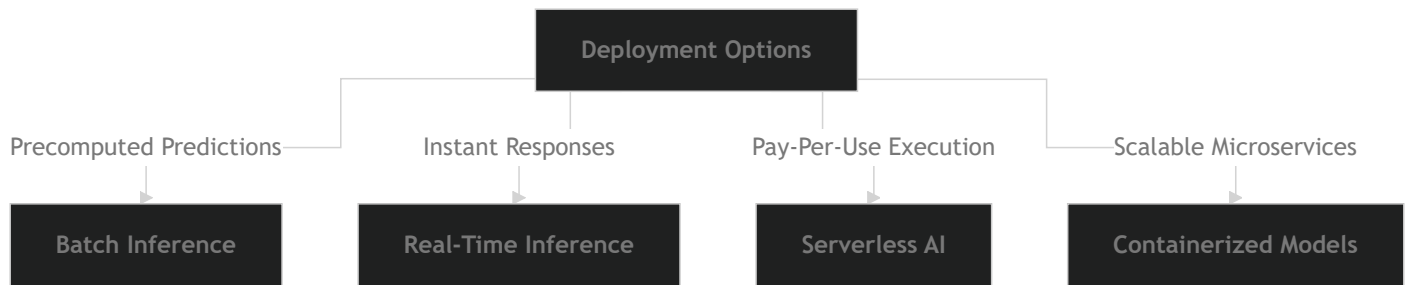


Scalable AI Deployment Strategies

Choosing the right deployment strategy depends on performance needs, cost, and operational complexity. Below is a **decision-making framework** for selecting an AI deployment strategy:



Deployment Options:



Hands-on Coding: Deployment & Cost Benchmarking

To make informed deployment decisions, AI architects must evaluate both **inference performance** and **cost efficiency**.

The following coding exercises will show some ways you can benchmark model inference times and estimate cloud costs, providing practical insights into optimizing AI deployment.

1. Benchmarking Model Inference Times

Copy

```

import time
import tensorflow as tf

# Load a sample model (pre-trained MobileNetV2)
model = tf.keras.applications.MobileNetV2()
input_data = tf.random.normal([1, 224, 224, 3])

# Measure inference time
start_time = time.time()
prediction = model(input_data)
inference_time = time.time() - start_time

print(f"Inference Time: {inference_time:.4f} seconds")

```

Think about it: How does inference time change with different model sizes and hardware?

2. Estimating Cloud Costs Using an API

Copy

```
import requests
```

```
# Example: Query AWS Pricing API for EC2 GPU instances
```

```
response = requests.get("https://pricing.us-east-1.amazonaws.com/example-pricing-  
data = response.json()
```

```
print("Sample Cost Estimate:", data['price'])
```



Think about it: How do different instance types affect cost-performance trade-offs?

Key Takeaway

By understanding **cost optimization**, **scalable deployment strategies**, and **benchmarking AI performance**, you can design AI architectures that balance cost, performance, and operational complexity.

[< Previous](#)

© 2025 General Assembly
[Attributions](#)

[Next >](#)

