



## ISDM (INDEPENDENT SKILL DEVELOPMENT MISSION)



# MODEL DEPLOYMENT STRATEGIES & API CREATION

## ◆ CHAPTER 1: INTRODUCTION TO MODEL DEPLOYMENT

### 1.1 What is Model Deployment?

Model deployment is the process of integrating a trained machine learning (ML) model into a **production environment** where it can **serve real-world predictions**.

### 1.2 Why is Model Deployment Important?

- **Enables real-time predictions** for applications.
- **Automates decision-making** based on data inputs.
- **Allows scalability** for handling multiple requests.
- **Bridges the gap between ML models and real-world applications.**

### 1.3 Common Deployment Scenarios

- ✓ **Web applications** (e.g., Chatbots, Recommendation Systems)
- ✓ **Mobile applications** (e.g., Face recognition, Speech-to-text)
- ✓ **Embedded systems** (e.g., IoT-based smart devices)
- ✓ **Enterprise software** (e.g., Fraud detection, Credit scoring)

## ◆ CHAPTER 2: MODEL DEPLOYMENT STRATEGIES

There are several ways to deploy ML models depending on the use case.

### 2.1 Batch Processing Deployment

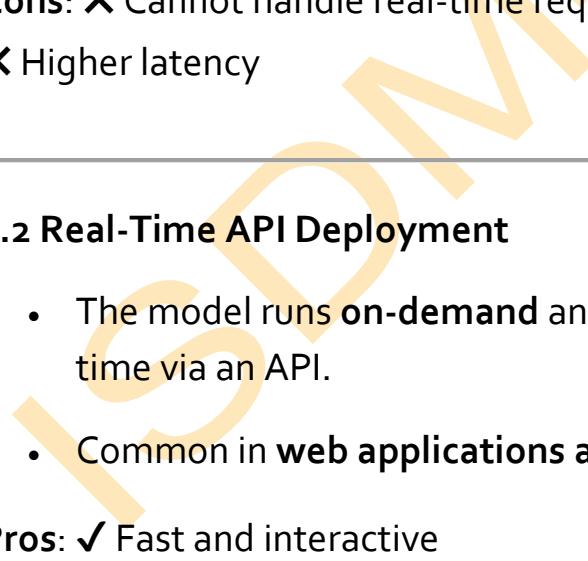
- Predictions are made **offline** and stored in a database for later use.
- Used when **real-time inference is not needed**.
- **Example:** Generating customer purchase recommendations once a day.

**Pros:** ✓ Efficient for large datasets

✓ Easier to scale

**Cons:** ✗ Cannot handle real-time requests

✗ Higher latency



---

### 2.2 Real-Time API Deployment

- The model runs **on-demand** and provides predictions in real time via an API.
- Common in **web applications and mobile apps**.

**Pros:** ✓ Fast and interactive

✓ Can handle multiple requests

**Cons:** ✗ Requires **server maintenance**

✗ **Latency issues** for complex models

---

## 2.3 Edge Deployment

- The model is deployed **locally** on a device rather than a server.
- Used in **IoT devices, mobile apps, and self-driving cars.**

**Pros:** ✓ Low latency

✓ Works without an internet connection

**Cons:** ✗ Limited computing power

✗ Difficult to update models

## 2.4 Cloud Deployment

- The model is deployed using **cloud services** (AWS, Google Cloud, Azure).
- Can handle large-scale **real-time inference.**

**Pros:** ✓ Highly scalable

✓ Secure and managed infrastructure

**Cons:** ✗ Costs depend on usage

✗ Requires **cloud expertise**

## 2.5 Containerized Deployment (Using Docker & Kubernetes)

- The model is packaged in a **Docker container** and managed using **Kubernetes**.
- Used for **scalable and portable deployments.**

**Pros:** ✓ Works across different environments

✓ Supports **auto-scaling**

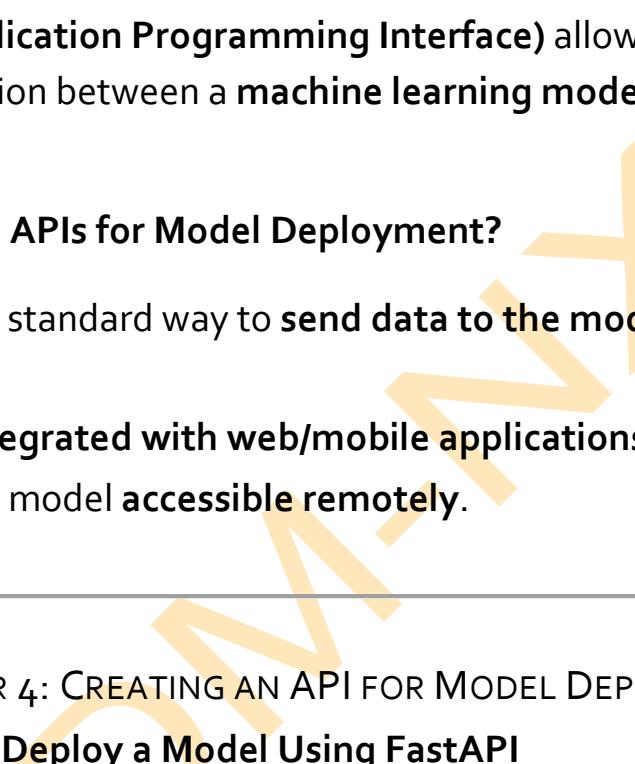
- Cons:**
- ✗ Complex setup
  - ✗ Requires container orchestration knowledge
- 

- ◆ CHAPTER 3: API CREATION FOR MODEL DEPLOYMENT

### 3.1 What is an API?

An **API (Application Programming Interface)** allows communication between a **machine learning model** and other applications.

### 3.2 Why Use APIs for Model Deployment?

- ✓ Provides a standard way to **send data to the model and get predictions.**
  - ✓ Can be **integrated with web/mobile applications.**
  - ✓ Makes the model **accessible remotely.**
- 

---

- ◆ CHAPTER 4: CREATING AN API FOR MODEL DEPLOYMENT

### 4.1 Steps to Deploy a Model Using FastAPI

FastAPI is a **lightweight, high-performance** Python framework used for creating APIs.

#### Step 1: Install FastAPI and Uvicorn

```
pip install fastapi uvicorn
```

#### Step 2: Train a Sample Model

```
import pickle
```

```
import numpy as np
```

```
from sklearn.linear_model import LogisticRegression
```

```
from sklearn.datasets import load_iris
```

```
# Load dataset
```

```
data = load_iris()
```

```
X, y = data.data, data.target
```

```
# Train model
```

```
model = LogisticRegression()
```

```
model.fit(X, y)
```

```
# Save model
```

```
with open("iris_model.pkl", "wb") as file:
```

```
    pickle.dump(model, file)
```

---

### Step 3: Create an API Using FastAPI

```
from fastapi import FastAPI
```

```
import pickle
```

```
import numpy as np
```

```
# Load trained model
```

```
with open("iris_model.pkl", "rb") as file:
```

```
    model = pickle.load(file)
```

```
# Create FastAPI instance
```

```
app = FastAPI()
```

```
# Define API endpoint
```

```
@app.post("/predict/")
```

```
def predict(features: list):
```

```
    features = np.array(features).reshape(1, -1)
```

```
    prediction = model.predict(features)[0]
```

```
    return {"prediction": int(prediction)}
```

---

#### Step 4: Run the API

```
uvicorn main:app --reload
```

- The API is now accessible at <http://127.0.0.1:8000>.
- 

#### Step 5: Test the API

```
import requests
```

```
url = "http://127.0.0.1:8000/predict/"
```

```
data = {"features": [5.1, 3.5, 1.4, 0.2]}
```

```
response = requests.post(url, json=data)
```

```
print(response.json())
```

---

## ◆ CHAPTER 5: DEPLOYING THE API ON A CLOUD SERVER

### 5.1 Deploying on AWS EC2

1. **Launch an EC2 instance (Ubuntu).**
2. **Install dependencies:**
3. sudo apt update
4. sudo apt install python3-pip
5. pip install fastapi uvicorn scikit-learn numpy
6. **Upload model files.**
7. **Run the FastAPI application.**
8. **Expose the API using a public IP.**

### 5.2 Deploying with Docker

#### 1. Create a Dockerfile

```
FROM python:3.8
```

```
WORKDIR /app
```

```
COPY . /app
```

```
RUN pip install fastapi uvicorn numpy scikit-learn
```

```
CMD ["uvicorn", "main:app", "--host", "0.0.0.0", "--port", "8000"]
```

#### 2. Build and Run Docker Container

```
docker build -t model-api .
```

```
docker run -p 8000:8000 model-api
```

---

#### ◆ CHAPTER 6: MONITORING AND SCALING

##### 6.1 API Performance Monitoring

- Use **Prometheus** and **Grafana** for monitoring.
- Track **response time, latency, and error rates**.

##### 6.2 Scaling the API

- Use **Kubernetes** for auto-scaling.
- Deploy on **AWS Lambda** for **serverless execution**.

#### ◆ CHAPTER 7: SUMMARY

Topic	Key Takeaways
Model Deployment	Process of integrating ML models into production.
Deployment Strategies	Batch, real-time, edge, cloud, and containerized deployment.
API Creation	FastAPI is a lightweight framework for deploying ML models.
Cloud & Docker Deployment	Models can be deployed on <b>AWS EC2, Docker, and Kubernetes</b> .
Monitoring & Scaling	Use <b>Prometheus, Grafana, and Kubernetes</b> to optimize API performance.

---

◆ CHAPTER 8: NEXT STEPS

- Deploy a model on Google Cloud or Azure.
- Use Flask for API creation instead of FastAPI.
- Learn how to scale APIs using Kubernetes.

ISDM-NxT



# DEPLOYING ML MODELS USING FLASK & FASTAPI

## 📌 CHAPTER 1: INTRODUCTION TO ML MODEL DEPLOYMENT

### 1.1 What is Model Deployment?

Model deployment is the process of integrating a trained machine learning model into a **production environment** where it can provide predictions on new data.

### 1.2 Why is Deployment Important?

- Allows real-time predictions for end-users.
- Makes ML models accessible via **APIs** for web or mobile applications.
- Enables **scalability**, so models can handle large volumes of requests.

### 1.3 Deployment Options

- **Local Deployment** (Flask, FastAPI, Streamlit)
- **Cloud Deployment** (AWS, GCP, Azure)
- **Containerized Deployment** (Docker, Kubernetes)

This guide focuses on **Flask and FastAPI** for local deployment.



## 📌 CHAPTER 2: DEPLOYING ML MODELS USING FLASK

### 2.1 What is Flask?

Flask is a **lightweight Python web framework** used to create RESTful APIs.

## 2.2 Installing Flask

```
pip install flask
```

## 2.3 Creating a Basic Flask API

```
from flask import Flask, request, jsonify
```

```
app = Flask(__name__)
```

```
@app.route("/")
```

```
def home():
```

```
    return "Welcome to Flask API!"
```

```
if __name__ == "__main__":
```

```
    app.run(debug=True)
```

Run the script and access <http://127.0.0.1:5000/> in your browser.

---

## CHAPTER 3: DEPLOYING AN ML MODEL WITH FLASK

### 3.1 Train and Save a Machine Learning Model

We will use **Scikit-learn** to train a simple **Iris classification model**.

```
import pickle
```

```
from sklearn.datasets import load_iris
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
# Load dataset  
  
iris = load_iris()  
  
X_train, X_test, y_train, y_test = train_test_split(iris.data, iris.target,  
test_size=0.2, random_state=42)  
  
# Train model  
  
model = RandomForestClassifier(n_estimators=100)  
model.fit(X_train, y_train)  
  
# Save model using pickle  
  
pickle.dump(model, open("model.pkl", "wb"))
```

---

### 3.2 Create Flask API for Model Deployment

```
import pickle  
  
from flask import Flask, request, jsonify  
  
# Load trained model  
  
model = pickle.load(open("model.pkl", "rb"))  
  
# Initialize Flask app  
  
app = Flask(__name__)
```

```
@app.route("/")
def home():

    return "ML Model API Running!"

@app.route("/predict", methods=["POST"])
def predict():

    data = request.json # Get data from request

    prediction = model.predict([data['features']]) # Make prediction

    return jsonify({"prediction": int(prediction[0])})

if __name__ == "__main__":
    app.run(debug=True)
```

---

### 3.3 Test the API Using Postman or cURL

Use **Postman** or **cURL** to test the endpoint.

#### Sample Request (POST)

URL: <http://127.0.0.1:5000/predict>

JSON Body:

```
{
    "features": [5.1, 3.5, 1.4, 0.2]
}
```

Expected Response:

```
{
```

```
    "prediction": o  
}
```

---

## 📌 CHAPTER 4: DEPLOYING ML MODELS USING FASTAPI

### 4.1 What is FastAPI?

FastAPI is a **modern, high-performance** Python framework for building APIs.

### 4.2 Installing FastAPI

```
pip install fastapi uvicorn
```

uvicorn is an ASGI server for running FastAPI applications.

---

## 📌 CHAPTER 5: DEPLOYING AN ML MODEL WITH FASTAPI

### 5.1 Create FastAPI ML Model Deployment

```
import pickle  
from fastapi import FastAPI  
from pydantic import BaseModel
```

```
# Load trained model
```

```
model = pickle.load(open("model.pkl", "rb"))
```

```
# Initialize FastAPI app
```

```
app = FastAPI()
```

```
class InputData(BaseModel):
    features: list

@app.get("/")
def home():
    return {"message": "FastAPI ML Model Running!"}

@app.post("/predict")
def predict(data: InputData):
    prediction = model.predict([data.features])
    return {"prediction": int(prediction[0])}
```

---

## 5.2 Run the FastAPI Server

Run the FastAPI server using:

uvicorn filename:app --reload

Example:

uvicorn fastapi\_app:app --reload

- Access API at: <http://127.0.0.1:8000/>
  - Access documentation at: <http://127.0.0.1:8000/docs>
- 

## 5.3 Test FastAPI Using Swagger UI

FastAPI automatically generates an **interactive API documentation**.

1. Open <http://127.0.0.1:8000/docs>.
2. Click on **POST /predict**.
3. Enter JSON input:

```
{
  "features": [5.1, 3.5, 1.4, 0.2]
}
```

4. Click **Execute** and check the response.

## CHAPTER 6: COMPARING FLASK AND FASTAPI

Feature	Flask	FastAPI
Speed	Slower	<input checked="" type="checkbox"/> Fast
Ease of Use	<input checked="" type="checkbox"/> Easy	<input checked="" type="checkbox"/> Easy
Async Support	<input checked="" type="checkbox"/> No	<input checked="" type="checkbox"/> Yes
Automatic Documentation	<input checked="" type="checkbox"/> No	<input checked="" type="checkbox"/> Yes (Swagger UI)

## CHAPTER 7: DEPLOYING FLASK/FASTAPI ON CLOUD

### 7.1 Deploy Using Docker

```
# Use Python image
```

```
FROM python:3.9
```

```
# Set working directory
```

```
WORKDIR /app
```

```
# Copy files  
COPY ..  
  
# Install dependencies  
RUN pip install -r requirements.txt
```

```
# Run Flask app  
CMD ["python", "app.py"]
```

Build and run the Docker container:  
docker build -t ml-api .  
docker run -p 5000:5000 ml-api

## 7.2 Deploy Using Heroku

1. Install **Heroku CLI**.

2. Login:

```
heroku login
```

3. Create a **requirements.txt** file:

```
pip freeze > requirements.txt
```

4. Deploy to Heroku:

```
git init
```

```
git add .
```

```
git commit -m "Deploy ML model"
```

```
heroku create my-ml-api
```

```
git push heroku master
```

---

## 📌 CHAPTER 8: SUMMARY & CONCLUSION

### 8.1 Key Takeaways

- ✓ **Flask** is great for **simple ML model deployment**.
- ✓ **FastAPI** is **faster** and supports **async requests**.
- ✓ APIs can be deployed on **Docker, Heroku, AWS, or GCP**.

### 8.2 Next Steps

- Use Docker to containerize ML models.
  - Deploy FastAPI on cloud platforms (AWS Lambda, Google Cloud).
  - Integrate ML models into web applications using React or Flask.
-

---

# CLOUD DEPLOYMENT ON AWS, GOOGLE CLOUD (GCP), AND AZURE

---

## CHAPTER 1: INTRODUCTION TO CLOUD DEPLOYMENT

### 1.1 What is Cloud Deployment?

Cloud Deployment refers to **deploying applications, services, or machine learning models** on cloud platforms like **Amazon Web Services (AWS)**, **Google Cloud Platform (GCP)**, and **Microsoft Azure**. Cloud platforms provide **scalability, security, and cost-effectiveness** without requiring on-premise hardware.

### 1.2 Why Use Cloud Deployment?

- Scalability** – Scale resources up or down based on demand.
- Cost Efficiency** – Pay only for the resources used.
- Reliability** – Cloud services ensure high availability and disaster recovery.
- Security** – Cloud providers offer **secure networking, authentication, and compliance**.
- Automation** – CI/CD pipelines streamline deployment and updates.

### 1.3 Key Cloud Deployment Models

- **Public Cloud** – Fully hosted on cloud providers (AWS, GCP, Azure).
- **Private Cloud** – Dedicated infrastructure for organizations.
- **Hybrid Cloud** – Combination of on-premise and cloud services.

## CHAPTER 2: CLOUD DEPLOYMENT ON AWS

### 2.1 What is AWS?

**Amazon Web Services (AWS)** is the leading cloud platform offering **compute, storage, networking, and AI services.**

### 2.2 AWS Services for Deployment

- **EC2 (Elastic Compute Cloud)** – Deploy virtual machines.
- **S3 (Simple Storage Service)** – Store files and static assets.
- **RDS (Relational Database Service)** – Managed databases (MySQL, PostgreSQL).
- **Lambda** – Serverless computing for event-driven applications.
- **Elastic Beanstalk** – Automated deployment for web applications.
- **ECS & EKS** – Deploy Docker containers using Kubernetes.

### 2.3 Steps to Deploy a Web App on AWS EC2

#### 1. Create an EC2 Instance

- Go to AWS EC2 Dashboard.
- Click **Launch Instance** and select **Ubuntu or Amazon Linux**.
- Choose instance type (**t2.micro for free tier**).
- Configure networking and security groups (allow SSH, HTTP/HTTPS).
- Add key pair (for SSH access) and launch the instance.

#### 2. Connect to EC2

```
ssh -i your-key.pem ec2-user@your-ec2-public-ip
```

### 3. Install Dependencies

```
sudo apt update
```

```
sudo apt install -y python3-pip nginx
```

### 4. Deploy a Flask Web App

```
pip3 install flask
```

```
echo "from flask import Flask"
```

```
app = Flask(__name__)
```

```
@app.route('/')
```

```
def home():
```

```
    return 'Hello, AWS!'
```

```
if __name__ == '__main__':
```

```
    app.run(host='0.0.0.0', port=5000)" > app.py
```

Run the app:

```
python3 app.py
```

### 5. Configure Security Group to Allow Traffic

- Allow port 5000 for external access.

### 6. Use Nginx as Reverse Proxy

```
sudo nano /etc/nginx/sites-available/default
```

Modify the server block:

```
server {
```

```
    listen 80;
```

```
    location / {
```

```
proxy_pass http://127.0.0.1:5000;  
}  
}
```

Restart Nginx:

```
sudo systemctl restart nginx
```

7. **Access the App via <http://your-ec2-public-ip>**

## 2.4 AWS Lambda for Serverless Deployment

1. **Go to AWS Lambda Console.**
2. **Create a new function (Choose Python).**
3. **Upload your application code.**
4. **Set up API Gateway for public access.**
5. **Invoke function using API Gateway URL.**

📌 **CHAPTER 3: CLOUD DEPLOYMENT ON GOOGLE CLOUD (GCP)**

### 3.1 What is GCP?

**Google Cloud Platform (GCP)** provides cloud computing services like **Compute Engine**, **Cloud Storage**, **Cloud Functions**, and **App Engine**.

### 3.2 GCP Services for Deployment

- **Compute Engine** – Deploy virtual machines.
- **App Engine** – Fully managed PaaS for deploying web apps.
- **Cloud Run** – Serverless containers for microservices.
- **Cloud Functions** – Event-driven serverless computing.

- **Kubernetes Engine (GKE)** – Deploy Kubernetes clusters.

### 3.3 Steps to Deploy a Web App on GCP Compute Engine

#### 1. Create a VM Instance

- Go to **GCP Console → Compute Engine**.
- Click **Create Instance** and choose **Ubuntu**.
- Set firewall rules to allow **HTTP & HTTPS**.

#### 2. Connect to VM

```
gcloud compute ssh instance-name
```

#### 3. Install Web Server

```
sudo apt update
```

```
sudo apt install apache2 -y
```

#### 4. Deploy Flask App

```
sudo apt install python3-pip
```

```
pip3 install flask
```

Run a Flask app similar to AWS steps.

#### 5. Allow Traffic Using Firewall Rules

```
gcloud compute firewall-rules create allow-http --allow tcp:80
```

#### 6. Access the Web App Use `http://your-vm-ip` in the browser.

### 3.4 Deploying on GCP App Engine (PaaS)

#### 1. Enable App Engine

```
gcloud app create --region=us-central
```

#### 2. Deploy the App

gcloud app deploy

### 3. View the App

gcloud app browse

---

## 📌 CHAPTER 4: CLOUD DEPLOYMENT ON MICROSOFT AZURE

### 4.1 What is Azure?

**Microsoft Azure** provides cloud services like **Azure Virtual Machines, App Services, Azure Functions, and Kubernetes Service.**

### 4.2 Azure Services for Deployment

- **Azure Virtual Machines** – Deploy Linux/Windows servers.
- **Azure App Service** – Managed hosting for web applications.
- **Azure Functions** – Serverless event-driven applications.
- **Azure Kubernetes Service (AKS)** – Managed Kubernetes for containerized apps.

### 4.3 Steps to Deploy on Azure Virtual Machines

#### 1. Create a Virtual Machine

- Go to **Azure Portal** → **Virtual Machines**.
- Choose **Ubuntu Server**.
- Configure firewall to allow HTTP (Port 80).

#### 2. Connect to VM

ssh azureuser@your-vm-ip

#### 3. Install Web Server

```
sudo apt update
```

```
sudo apt install nginx -y
```

#### 4. Deploy Flask App

```
pip3 install flask
```

Run the app and configure Nginx as a reverse proxy (same as AWS).

#### 5. Access the App Go to <http://your-vm-ip>

### 4.4 Deploying Web App on Azure App Service

#### 1. Create an App Service

- o Go to **Azure Portal → App Services.**
- o Choose **Python Flask App template.**

#### 2. Deploy the Application

- o Use **GitHub Actions or Azure CLI** to deploy.

#### 3. Access the App Use the generated **Azure App URL.**

## 📌 CHAPTER 5: COMPARING AWS, GCP, AND AZURE

Feature	AWS	GCP	Azure
Compute	EC2	Compute Engine	Virtual Machines
Serverless	Lambda	Cloud Functions	Azure Functions
Managed App Hosting	Elastic Beanstalk	App Engine	App Service

Container Deployment	EKS	GKE	AKS
Pricing	Pay-as-you-go	Pay-as-you-go	Pay-as-you-go
Best For	Enterprise & Startups	AI & Big Data	Enterprises using Microsoft Stack

## 📌 CHAPTER 6: SUMMARY & NEXT STEPS

### 6.1 Key Takeaways

- ✓ AWS is the most mature cloud platform with extensive services.
- ✓ GCP is best for AI & data-heavy applications.
- ✓ Azure is ideal for businesses using Microsoft technologies.
- ✓ Serverless platforms (Lambda, Cloud Functions) reduce infrastructure management.
- ✓ CI/CD pipelines automate deployments for scalability.

### 6.2 Next Steps

- Implement Docker & Kubernetes for containerized apps.
- Use Terraform for Infrastructure as Code (IaC).
- Set up CI/CD pipelines with GitHub Actions for automated deployment.



# EDGE AI & IOT INTEGRATION

## 📌 CHAPTER 1: INTRODUCTION TO EDGE AI & IOT INTEGRATION

### 1.1 What is Edge AI?

**Edge AI** refers to running **artificial intelligence (AI) models directly on edge devices**, such as **smart sensors, IoT devices, mobile phones, and industrial machines**. Instead of relying on cloud computing, AI computations occur **locally on the device**.

### 1.2 What is IoT (Internet of Things)?

**IoT (Internet of Things)** is a network of **connected devices** that collect and exchange data over the internet. Examples include **smart thermostats, industrial sensors, security cameras, and wearable devices**.

### 1.3 Why Integrate Edge AI with IoT?

- ✓ **Faster Processing** – AI inference happens on **edge devices** instead of cloud servers.
- ✓ **Reduced Latency** – Critical AI decisions (e.g., **self-driving cars, industrial robots**) are **real-time**.
- ✓ **Enhanced Security & Privacy** – Less data sent to the cloud means **better data privacy**.
- ✓ **Reduced Bandwidth Usage** – Edge computing lowers **network congestion and cloud storage costs**.
- ✓ **Better Reliability** – AI applications work **offline or with minimal internet dependency**.

## CHAPTER 2: KEY COMPONENTS OF EDGE AI & IoT

### 2.1 Edge AI Hardware

Device Type	Examples	Use Case
<b>Microcontrollers (MCUs)</b>	Raspberry Pi, Arduino	Home automation, small AI tasks
<b>Edge GPUs &amp; TPUs</b>	NVIDIA Jetson Nano, Google Coral	AI vision, robotics, real-time analytics
<b>Smart Sensors</b>	LiDAR, Temperature sensors	Autonomous vehicles, environmental monitoring
<b>Industrial Gateways</b>	Siemens IoT Gateway	Predictive maintenance, manufacturing

### 2.2 IoT Communication Protocols

- ◆ **MQTT (Message Queuing Telemetry Transport)** – Lightweight protocol for **low-bandwidth IoT communication**.
- ◆ **HTTP/HTTPS** – Standard web-based IoT communication.
- ◆ **CoAP (Constrained Application Protocol)** – Efficient for **low-power devices**.
- ◆ **LoRaWAN** – Used for **long-range communication in smart agriculture & city monitoring**.

### 2.3 AI Models on Edge Devices

1. **Computer Vision** (Face recognition, defect detection).
2. **Natural Language Processing (NLP)** (Voice assistants, chatbots).
3. **Predictive Analytics** (Smart manufacturing, healthcare AI).

#### 4. Anomaly Detection (Cybersecurity, fraud prevention).

---

### 📌 CHAPTER 3: SETTING UP AN EDGE AI IoT SYSTEM

#### 3.1 Selecting the Right Hardware

- For simple AI tasks, use Raspberry Pi or Arduino.
- For complex AI inference, use NVIDIA Jetson Nano or Google Coral.

#### 3.2 Installing Required Software

```
# Install TensorFlow Lite (for running AI models on edge devices)
```

```
pip install tflite-runtime
```

```
# Install MQTT for IoT communication
```

```
pip install paho-mqtt
```

#### 3.3 Edge AI + IoT Example: Real-Time Object Detection

##### Step 1: Install Required Libraries

```
import cv2
```

```
import tensorflow.lite as tflite
```

```
import numpy as np
```

```
import paho.mqtt.client as mqtt
```

##### Step 2: Load AI Model on Edge Device

```
# Load TensorFlow Lite Model
```

```
interpreter = tflite.Interpreter(model_path="mobilenet.tflite")
```

```
interpreter.allocate_tensors()
```

### Step 3: Capture and Process Camera Feed

```
# Open camera
cap = cv2.VideoCapture(0)

while True:
    ret, frame = cap.read()
    if not ret:
        break
    # Preprocess frame for AI model
    input_data = cv2.resize(frame, (224, 224))
    input_data = np.expand_dims(input_data,
                                axis=0).astype(np.float32)
    # Run inference
    interpreter.set_tensor(interpreter.get_input_details()[0]['index'],
                           input_data)
    interpreter.invoke()
    output =
    interpreter.get_tensor(interpreter.get_output_details()[0]['index'])

    # Display result
    cv2.imshow('Edge AI Object Detection', frame)
    if cv2.waitKey(1) & 0xFF == ord('q'):
```

```
break
```

```
cap.release()
```

```
cv2.destroyAllWindows()
```

#### Step 4: Send Data to IoT Cloud Using MQTT

```
# Connect to IoT broker
```

```
client = mqtt.Client()
```

```
client.connect("broker.hivemq.com", 1883, 60)
```

```
# Send detected object data to IoT cloud
```

```
client.publish("iot/edge-ai", "Detected Object: Car")
```

## CHAPTER 4: APPLICATIONS OF EDGE AI & IoT

### 4.1 Smart Cities

 **Traffic Management** – AI detects congestion and adjusts traffic signals.

 **Surveillance Systems** – AI-powered CCTV for **real-time threat detection**.

 **Environmental Monitoring** – IoT sensors measure **air quality** and pollution.

### 4.2 Industrial IoT (IIoT)

 **Predictive Maintenance** – AI detects **machine failures before they happen**.

 **Energy Optimization** – AI reduces **power consumption in factories**.

- 🔍 **Quality Inspection** – AI cameras detect **manufacturing defects** in real-time.

### 4.3 Healthcare AI & IoT

- ⌚ **Wearable Health Monitoring** – Smartwatches detect **heart rate, oxygen levels, and ECG**.
- 🏥 **AI in Hospitals** – Real-time patient monitoring **alerts doctors of emergencies**.
- 🧠 **AI-Powered Diagnosis** – Edge AI scans X-rays and **detects diseases instantly**.

---

## CHAPTER 5: SECURITY CHALLENGES IN EDGE AI & IoT

### 5.1 Key Security Risks

- 🔴 **Data Privacy Issues** – AI models process sensitive **user data**.
- 🔴 **Cyberattacks** – Hackers can attack **IoT networks**.
- 🔴 **Model Poisoning** – Adversarial attacks modify AI predictions.

### 5.2 Solutions for Securing Edge AI & IoT

- ✓ **Use Secure Boot & Encryption** – Protects edge devices from unauthorized firmware.
- ✓ **Deploy AI Security Models** – Detects **anomalies & cyber threats**.
- ✓ **Access Control & Authentication** – Uses **blockchain & biometrics** for secure login.

## CHAPTER 6: SUMMARY

Feature	Edge AI	IoT
Purpose	Runs AI models on edge devices	Connects devices to the internet
Processing	On-device (local)	Cloud or local
Latency	Low (real-time)	Higher (depends on cloud speed)
Use Case	Self-driving cars, security cameras	Smart homes, industrial automation

## CHAPTER 7: CONCLUSION & NEXT STEPS

### 7.1 Key Takeaways

-  Edge AI + IoT reduces latency and improves efficiency.
-  AI models can run on IoT devices using hardware accelerators.
-  Security is crucial for Edge AI deployments.

### 7.2 Next Steps

-  Try deploying Edge AI models on Raspberry Pi.
-  Experiment with real-time IoT sensors & AI integration.
-  Explore AI-powered robotics with IoT connectivity.

---



# AI MODEL MONITORING & PERFORMANCE OPTIMIZATION

---

## 📌 CHAPTER 1: INTRODUCTION TO AI MODEL MONITORING & OPTIMIZATION

### 1.1 What is AI Model Monitoring?

AI model monitoring involves tracking the **performance, accuracy, and reliability** of machine learning (ML) models in **production**. It ensures that models **stay accurate over time** and detect issues like **data drift, bias, and concept drift**.

### 1.2 Why is Model Monitoring Important?

- Ensures **model accuracy** over time.
- Detects **drift in data and predictions**.
- Prevents **bias and fairness issues**.
- Helps in **debugging and retraining**.

### 1.3 What is Model Optimization?

AI model optimization refers to improving a model's **accuracy, efficiency, and robustness** by fine-tuning hyperparameters, reducing computational costs, and ensuring high performance.

## 📌 CHAPTER 2: KEY CHALLENGES IN AI MODEL MONITORING

### 2.1 Challenges in AI Model Deployment

✖ **Data Drift** – Changes in input data distribution affect model predictions.

✖ **Concept Drift** – The relationship between features and target

changes over time.

✖ **Model Degradation** – Accuracy drops due to outdated training data.

✖ **Performance Bottlenecks** – High latency or expensive computation costs.

## 2.2 Symptoms of a Degrading AI Model

- ◆ Increased **error rate** or **misclassifications**.
- ◆ Poor performance on **new data**.
- ◆ Increased **bias** or **fairness issues**.
- ◆ Unexpected **spikes in resource usage**.

## 📌 CHAPTER 3: AI MODEL MONITORING TECHNIQUES

### 3.1 Key Metrics to Monitor

Metric	Description
<b>Accuracy</b>	Correct predictions over total predictions
<b>Precision &amp; Recall</b>	Measure false positives and false negatives
<b>F1 Score</b>	Balance between precision and recall
<b>Drift Detection</b>	Identifies changes in data distribution
<b>Inference Time</b>	Measures latency per prediction
<b>Model Fairness</b>	Ensures unbiased predictions across demographics

### 3.2 Tools for Model Monitoring

✓ **Prometheus & Grafana** – Monitors AI model metrics in real-time.

- ✓ **MLflow** – Tracks experiments, logs models, and manages deployment.
  - ✓ **Evidently AI** – Detects **data drift** and **concept drift**.
  - ✓ **AWS SageMaker Model Monitor** – Tracks model performance on cloud platforms.
  - ✓ **TensorFlow Model Analysis (TFMA)** – Evaluates **TensorFlow models** in production.
- 

## 📌 CHAPTER 4: AI MODEL PERFORMANCE OPTIMIZATION

### 4.1 Techniques for Improving AI Model Performance

1. **Feature Engineering** – Selecting and creating **relevant features**.
2. **Hyperparameter Tuning** – Optimizing **learning rates, batch size, dropout rate**.
3. **Reducing Overfitting** – Using **regularization, dropout, and data augmentation**.
4. **Parallelization** – Using **multi-threading, GPUs, or TPUs** for efficiency.
5. **Model Quantization** – Reducing model size for **faster inference**.

### 4.2 Hyperparameter Tuning Techniques

- ✓ **Grid Search** – Exhaustively searches for the best parameters.
- ✓ **Random Search** – Randomly selects parameter combinations.
- ✓ **Bayesian Optimization** – Uses probabilistic models to find the best parameters.
- ✓ **Optuna & Hyperopt** – Automated hyperparameter tuning libraries.

## 📌 CHAPTER 5: IMPLEMENTING MODEL MONITORING IN PYTHON

### 5.1 Install Required Libraries

```
!pip install mlflow evidently pandas numpy scikit-learn
```

### 5.2 Load Sample Dataset & Train a Model

```
import pandas as pd  
  
import numpy as np  
  
from sklearn.model_selection import train_test_split  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.metrics import accuracy_score  
  
# Load dataset  
df = pd.read_csv("creditcard.csv")  
  
# Split data  
X = df.drop(columns=["Class"])  
y = df["Class"]  
  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,  
random_state=42)  
  
# Train model  
model = RandomForestClassifier(n_estimators=100,  
random_state=42)
```

```
model.fit(X_train, y_train)

# Evaluate model

y_pred = model.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))
```

---

### 5.3 Log Model Performance Using MLflow

```
import mlflow

import mlflow.sklearn

# Start MLflow tracking

mlflow.start_run()

# Log model

mlflow.sklearn.log_model(model, "random_forest_model")

mlflow.log_metric("accuracy", accuracy_score(y_test, y_pred))

mlflow.end_run()
```

---

### 5.4 Detecting Data Drift with Evidently AI

```
from evidently import ColumnDriftCalculator

from evidently.dashboard import Dashboard
```

```
from evidently.dashboard.tabs import DataDriftTab
```

```
# Create a reference dataset (historical data)
```

```
reference_data = X_train.copy()
```

```
production_data = X_test.copy()
```

```
# Check data drift
```

```
drift_calc = ColumnDriftCalculator()
```

```
drift_calc.fit(reference_data)
```

```
drift_calc.calculate(production_data)
```

```
# Generate report
```

```
report = Dashboard(tabs=[DataDriftTab()])
```

```
report.calculate(reference_data, production_data)
```

```
report.show()
```

---

## CHAPTER 6: REAL-WORLD AI MODEL MONITORING & OPTIMIZATION

### 6.1 Case Study: AI Model for Fraud Detection

A bank deploys an AI model to **detect fraudulent transactions**.

- ◆ **Challenge** – Over time, fraudsters change their tactics (**concept drift**).
- ◆ **Solution** – Use **drift detection tools** to detect changes in transaction patterns.

- ◆ **Optimization** – Fine-tune the model using **adaptive learning** strategies.

## 6.2 Case Study: AI in Healthcare

A hospital uses AI for **predicting patient readmissions**.

- ◆ **Challenge** – New health conditions emerge, leading to **data drift**.
- ◆ **Solution** – Monitor **accuracy changes** using MLflow.
- ◆ **Optimization** – Retrain the model with **new patient data** every month.

## CHAPTER 7: CHALLENGES & FUTURE OF AI MODEL MONITORING

### 7.1 Challenges in Model Monitoring

- ✖ **Model decay** – Performance drops over time.
- ✖ **Lack of real-time monitoring** – Requires automated tracking tools.
- ✖ **Regulatory compliance** – AI models must be auditable.

### 7.2 Future Trends in AI Optimization

- 🚀 **AutoML for real-time tuning** – Self-optimizing models.
- 🚀 **AI-powered explainability tools** – Better model interpretability.
- 🚀 **Federated Learning** – Training AI models across decentralized data sources.

## CHAPTER 8: SUMMARY

- ✓ AI model monitoring ensures that models remain **accurate and reliable**.
- ✓ Detecting **data drift & concept drift** helps maintain model

performance.

- ✓ Optimization techniques like **hyperparameter tuning** improve model efficiency.
  - ✓ Tools like **MLflow**, **Evidently AI**, and **Prometheus** assist in real-time monitoring.
- 

📌 CHAPTER 9: NEXT STEPS

- 🚀 Deploy monitoring dashboards for real-time AI tracking.
- 🚀 Use Bayesian Optimization for hyperparameter tuning.
- 🚀 Apply real-time drift detection in production AI models.

ISDM-N

---

📌 ⚡ ASSIGNMENT 1:  
🎯 DEPLOY A MACHINE LEARNING MODEL  
AS AN API USING FLASK.

ISDM-NxT



# ASSIGNMENT SOLUTION 1: DEPLOY A MACHINE LEARNING MODEL AS AN API USING FLASK

## 🎯 Objective

The goal of this assignment is to **deploy a trained machine learning model as a REST API using Flask**. This API will allow users to send data and receive predictions in real-time.

We will:

- Train and save a **machine learning model**
- Create a **Flask API** to serve the model
- Test the API using **Postman or Python requests**
- Deploy the API on a **cloud server or Docker**

## 🛠 Step 1: Install and Import Required Libraries

We will use **Flask** to create the API and **scikit-learn** to train the model.

### 1.1 Install Dependencies (if not installed)

```
pip install flask pandas numpy scikit-learn requests
```

### 1.2 Import Required Libraries

```
import pickle
```

```
import numpy as np
```

```
import pandas as pd
```

```
from flask import Flask, request, jsonify  
from sklearn.model_selection import train_test_split  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.datasets import load_iris
```

---



## Step 2: Train and Save a Machine Learning Model

We will use the **Iris dataset** and train a **Random Forest Classifier**.

### 2.1 Load and Prepare the Data

```
# Load dataset  
data = load_iris()  
X, y = data.data, data.target
```

```
# Split data into training and testing sets
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,  
random_state=42)
```

---

### 2.2 Train the Model

```
# Train a Random Forest Classifier
```

```
model = RandomForestClassifier(n_estimators=100,  
random_state=42)  
  
model.fit(X_train, y_train)
```

```
# Evaluate accuracy
```

---

```
accuracy = model.score(X_test, y_test)
print(f"Model Accuracy: {accuracy:.2f}")
```

---

## 2.3 Save the Trained Model

We will **save the model** using pickle so we can load it later in our Flask API.

```
# Save model as a pickle file
with open("iris_model.pkl", "wb") as file:
    pickle.dump(model, file)
```

---

### ➡ Step 3: Create a Flask API

Now, we will create an API using Flask that serves the model and responds to requests.

#### 3.1 Create a Flask Application

```
# Initialize Flask app
app = Flask(__name__)

# Load the trained model
with open("iris_model.pkl", "rb") as file:
    model = pickle.load(file)
```

```
# Define API endpoint for prediction
@app.route("/predict", methods=["POST"])
```

```
def predict():

    try:

        # Get JSON data from request

        data = request.get_json()

        features = np.array(data["features"]).reshape(1, -1)

        # Make prediction

        prediction = model.predict(features)[0]

        # Return prediction as JSON

        return jsonify({"prediction": int(prediction)})

    except Exception as e:

        return jsonify({"error": str(e)})

# Run the Flask app

if __name__ == "__main__":

    app.run(debug=True, host="0.0.0.0", port=5000)
```

## ➡ Step 4: Run and Test the API Locally

### 4.1 Start the Flask API

Run the Flask application:

python app.py

✓ The API will start at <http://127.0.0.1:5000>.

---

## 4.2 Test the API Using Python Requests

We can send a POST request with sample data to test the API.

```
import requests
```

```
url = "http://127.0.0.1:5000/predict"  
  
data = {"features": [5.1, 3.5, 1.4, 0.2]}  
  
response = requests.post(url, json=data)  
  
print(response.json()) # Expected output: {"prediction": 0, 1, or 2}
```

---

## 4.3 Test the API Using Postman

1. Open **Postman** and create a **POST request**.
2. Enter <http://127.0.0.1:5000/predict> in the request URL.
3. Select **Body → raw → JSON** and enter:

```
{  
  
    "features": [5.1, 3.5, 1.4, 0.2]  
  
}
```

4. Click **Send** to get the prediction response.
-

## ➡ Step 5: Deploy the Flask API on a Cloud Server

To make the API accessible over the internet, we deploy it on **AWS EC2**, **Google Cloud**, or **DigitalOcean**.

### 5.1 Deploy on an AWS EC2 Instance

1. **Launch an EC2 Instance (Ubuntu)**
2. **Connect to the server using SSH**
3. `ssh -i "your-key.pem" ubuntu@your-server-ip`
4. **Install Python and Dependencies**
5. `sudo apt update`
6. `sudo apt install python3-pip`
7. `pip3 install flask pandas numpy scikit-learn`
8. **Upload the Model and Flask App**
9. `scp -i "your-key.pem" app.py iris_model.pkl ubuntu@your-server-ip:/home/ubuntu/`
10. **Run the Flask App**
11. `python3 app.py`
12. **Access the API using Public IP**
13. `http://your-server-ip:5000/predict`

### 5.2 Deploy Using Docker

1. **Install Docker**
2. `sudo apt update`
3. `sudo apt install docker.io`

#### 4. Create a Dockerfile

5. FROM python:3.8
  6. WORKDIR /app
  7. COPY . /app
  8. RUN pip install flask pandas numpy scikit-learn
  9. CMD ["python", "app.py"]
  10. **Build and Run the Docker Container**
  11. docker build -t flask-api .
  12. docker run -p 5000:5000 flask-api
- 

#### 📌 Step 6: Monitor and Scale the API

To ensure **reliability and performance**, we can:

- ✓ Use **NGINX** as a reverse proxy for handling requests.

- ✓ Monitor API usage using **Prometheus and Grafana**.
  - ✓ Scale using **Kubernetes** for load balancing.
- 

#### ✓ FINAL SUMMARY

Step	Description
Step 1	Install dependencies and import required libraries
Step 2	Train and save a <b>Random Forest model</b> using scikit-learn
Step 3	Create a <b>Flask API</b> to serve the model

<b>Step 4</b>	Run and test the API locally using <b>Python requests &amp; Postman</b>
<b>Step 5</b>	Deploy the API on <b>AWS EC2 or Docker</b>
<b>Step 6</b>	Monitor and scale the API

## ★ CONCLUSION

- **Flask** makes it easy to deploy machine learning models as APIs.
- **APIs allow real-time predictions** for web and mobile apps.
- Deployment on **AWS or Docker** ensures **scalability and availability**.

ISDM - INSTITUTE OF SKILL DEVELOPMENT AND MANAGEMENT

---

📌 ⚡ ASSIGNMENT 2:  
🎯 INTEGRATE AI MODELS WITH AWS  
LAMBDA & GOOGLE VERTEX AI.

ISDM-NXT

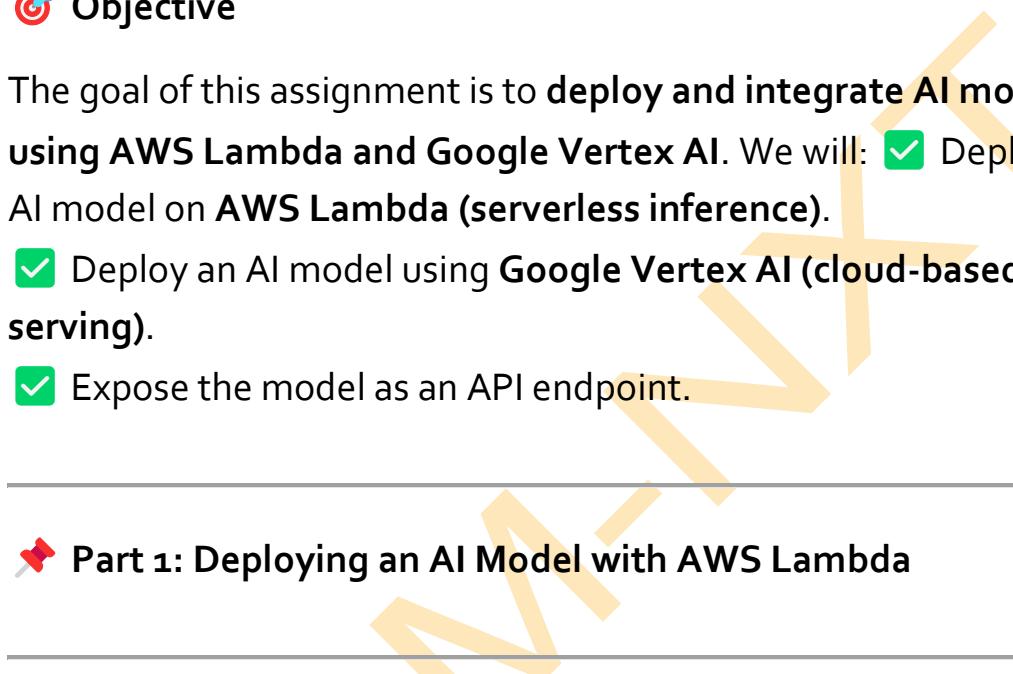
---

 **ASSIGNMENT SOLUTION 2:**  
**INTEGRATE AI MODELS WITH AWS LAMBDA**  
**& GOOGLE VERTEX AI**

---

 **Objective**

The goal of this assignment is to **deploy and integrate AI models using AWS Lambda and Google Vertex AI**. We will:  Deploy an AI model on **AWS Lambda (serverless inference)**.

- 
-  Deploy an AI model using **Google Vertex AI (cloud-based ML serving)**.
  -  Expose the model as an API endpoint.
- 

 **Part 1: Deploying an AI Model with AWS Lambda** **Step 1: Train and Save the ML Model**

We will train a **Random Forest classifier** using the **Iris dataset** and save it as a .pkl file.

 **Install Dependencies**

```
pip install scikit-learn pandas pickle5
```

 **Train and Save the Model**

```
import pickle
```

```
from sklearn.datasets import load_iris
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
from sklearn.model_selection import train_test_split

# Load dataset

iris = load_iris()

X_train, X_test, y_train, y_test = train_test_split(iris.data, iris.target,
test_size=0.2, random_state=42)

# Train model

model = RandomForestClassifier(n_estimators=100)

model.fit(X_train, y_train)

# Save the trained model

pickle.dump(model, open("model.pkl", "wb"))
```

## ➡ Step 2: Create an AWS Lambda Function

AWS Lambda allows us to deploy serverless applications.

- ◆ **Install AWS CLI & Set Up Credentials**

1. Install AWS CLI:
2. pip install awscli
3. Configure AWS credentials:
4. aws configure

Provide:

- **AWS Access Key**

- **AWS Secret Access Key**
  - **Region** (e.g., us-east-1)
- 

### 📌 Step 3: Prepare a Lambda Function

AWS Lambda requires a **handler function** to process requests.

- ◆ **Create lambda\_function.py**

```
import json
import pickle
import numpy as np

# Load trained model
model = pickle.load(open("/tmp/model.pkl", "rb"))

def lambda_handler(event, context):
    # Parse input data
    input_data = json.loads(event['body'])['features']
    input_array = np.array([input_data])

    # Make prediction
    prediction = model.predict(input_array)

    # Return response
```

```
return {  
  
    'statusCode': 200,  
  
    'body': json.dumps({"prediction": int(prediction[0])})  
  
}
```

---

#### ➡ Step 4: Deploy Lambda Function

##### ◆ Create a Deployment Package

1. Install dependencies locally:
2. pip install --target ./package numpy scikit-learn pickle5
3. Add Lambda function and model:
4. cp lambda\_function.py model.pkl package/
5. cd package
6. zip -r ..//lambda\_model.zip .
7. cd ..

##### ◆ Deploy to AWS Lambda

1. **Create Lambda function:**
2. aws lambda create-function \  
3. --function-name AI\_Model\_Lambda \  
4. --runtime python3.8 \  
5. --role arn:aws:iam::your-account-id:role/your-role \  
6. --handler lambda\_function.lambda\_handler \  
7. --zip-file fileb://lambda\_model.zip

## 8. Invoke Lambda function:

9. aws lambda invoke --function-name AI\_Model\_Lambda  
output.json

## 10. Deploy API Gateway for REST API:

- Open AWS Console → API Gateway.
- Create API → Connect with Lambda function.
- Deploy and get an endpoint.

---

## Part 2: Deploying an AI Model with Google Vertex AI

---

### ❖ Step 1: Train and Save Model

We assume the model is trained using the same **Iris dataset** (as done in Part 1).

#### ◆ Save Model Using Google Cloud Storage

1. Install google-cloud-storage:
2. pip install google-cloud-storage
3. Upload model to **Google Cloud Storage (GCS)**:
4. from google.cloud import storage
- 5.
6. # Set up GCS client
7. client = storage.Client()
8. bucket = client.bucket("your-bucket-name")
- 9.

```
10.      # Upload model file  
11.blob = bucket.blob("model.pkl")  
12.      blob.upload_from_filename("model.pkl")  
13.  
14.      print("Model uploaded successfully!")
```

---

➡ **Step 2: Deploy Model to Google Vertex AI**

◆ **Install Google Cloud SDK**

```
pip install google-cloud-aiplatform
```

◆ **Upload Model to Vertex AI**

```
from google.cloud import aiplatform
```

```
# Initialize Vertex AI client
```

```
aiplatform.init(project="your-project-id", location="us-central1")
```

```
# Upload model to Vertex AI
```

```
model = aiplatform.Model.upload(
```

```
    display_name="iris-classifier",
```

```
    artifact_uri="gs://your-bucket-name/model.pkl",
```

```
    serving_container_image_uri="us-docker.pkg.dev/vertex-  
    ai/prediction/sklearn-cpu.0-24:latest"
```

```
)
```

```
print(f"Model Deployed: {model.resource_name}")
```

---

### 📌 Step 3: Deploy Model as an API Endpoint

```
endpoint = model.deploy(  
    machine_type="n1-standard-4"  
)
```

```
print(f"Model Endpoint: {endpoint.resource_name}")
```

---

### 📌 Step 4: Test the Deployed Model

```
instances = [[5.1, 3.5, 1.4, 0.2]] # Test input  
  
prediction = endpoint.predict(instances)  
print("Prediction:", prediction)
```

---

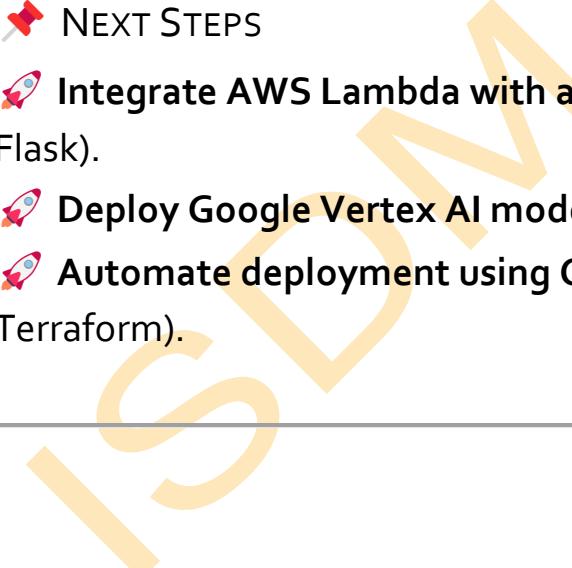
### 📌 SUMMARY

Step	AWS Lambda	Google Vertex AI
<b>Train Model</b>	Train with Scikit-learn	Train with Scikit-learn
<b>Save Model</b>	Save as model.pkl	Save to Google Cloud Storage

<b>Deploy API</b>	AWS Lambda + API Gateway	Vertex AI Endpoint
<b>Test Prediction</b>	aws lambda invoke	endpoint.predict()

---

 CONCLUSION

-  AWS Lambda provides **serverless model deployment** with low cost.
  -  Google Vertex AI offers **cloud-based ML model serving** with **scalability**.
  -  Both services expose an **API endpoint for inference**.
- 

 NEXT STEPS

-  Integrate AWS Lambda with a front-end web app (React, Flask).
  -  Deploy Google Vertex AI model on Kubernetes (GKE).
  -  Automate deployment using CI/CD pipelines (GitHub Actions, Terraform).
-