



**SWE4004**

**CLOUD COMPUTING PROJECT REPORT**

Title: End-to-End ML Model Deployment on AWS Cloud Using CI/CD Pipeline

**By:**

**21MIS7066**

**21MIS7041**

**21MIC7068**

**Under the guidance of:**

**Prof. S. Bharath Bhushan**

## Table of Contents

1. Abstract .....	3
2. Problem Statement .....	4
3. Introduction .....	5
4. Objectives .....	6
5. Problems in the Existing System .....	7
6. Methodology .....	9
- Setting Up the AWS Environment .....	9
- Configuring GitHub for CI/CD .....	10
- Implementing the CI/CD Pipeline .....	11
- Deploying the ML Model .....	14
7. Results .....	15
8. Conclusion .....	17
9. Future Scope .....	18
10. References .....	20
11. Appendices .....	21

## **ABSTRACT**

This project focuses on deploying end-to-end machine learning (ML) projects in a production environment using AWS Cloud and a CI/CD pipeline. The deployment process involves creating an AWS Ubuntu instance and setting up a private repository on Amazon Elastic Container Registry (ECR) to store Docker images. The CI/CD pipeline is implemented using GitHub Actions, with a GitHub self-hosted runner to automate the process of pushing updates to the AWS instance.

The pipeline is designed to ensure seamless integration, testing, and deployment of ML models. It includes stages for continuous integration, where code is linted and unit tests are run, and continuous delivery, where Docker images are built, tagged, and pushed to ECR. The final stage, continuous deployment, pulls the latest Docker image and runs it on the AWS instance, making the ML model available for production use.

This approach leverages the scalability and reliability of AWS services, combined with the automation capabilities of GitHub Actions, to create a robust and efficient deployment process. The project demonstrates the integration of cloud infrastructure with modern DevOps practices, enabling rapid and reliable deployment of ML models in a production setting.

## **PROBLEM STATEMENT**

Deploying machine learning (ML) models into production environments is fraught with challenges, particularly regarding efficiency, reliability, and scalability. Traditional manual deployment methods are error-prone, time-consuming, and inconsistent, leading to prolonged downtime, suboptimal model performance, and increased operational costs. Moreover, integrating continuous integration and continuous deployment (CI/CD) pipelines with ML workflows in cloud environments adds layers of complexity. Organizations struggle with setting up and managing cloud resources, automating build and deployment processes, ensuring scalability and reliability, and seamlessly integrating these processes with existing development workflows. This project aims to address these issues by developing a robust CI/CD pipeline that automates the end-to-end deployment of ML models on AWS, leveraging cloud infrastructure and modern DevOps practices to ensure efficient, reliable, and scalable production deployments.

## INTRODUCTION

In the rapidly evolving field of machine learning (ML), the deployment of models into production environments has become a critical yet challenging task. Ensuring that ML models can be deployed efficiently, reliably, and at scale is essential for organizations looking to leverage their data-driven insights in real-time applications. Traditional methods of deploying ML models often involve manual processes that are not only time-consuming but also susceptible to errors and inconsistencies. These manual approaches can lead to significant downtime, reduced performance of models, and increased operational costs.

The integration of continuous integration and continuous deployment (CI/CD) pipelines with ML workflows presents a modern solution to these challenges. CI/CD pipelines automate the processes of testing, building, and deploying software, enabling more frequent and reliable updates. However, implementing such pipelines for ML models in cloud environments, particularly on platforms like AWS, introduces its own set of complexities. These include the setup and management of cloud resources, automation of build and deployment processes, and ensuring that the deployed models can scale effectively and remain reliable under varying loads.

This project aims to address these challenges by developing an end-to-end CI/CD pipeline for deploying ML models into production on AWS Cloud. The solution involves creating an AWS t2.micro Ubuntu instance, setting up a private repository on Amazon Elastic Container Registry (ECR) for storing Docker images, and using a GitHub self-hosted runner to automate the deployment process. By leveraging the scalability and reliability of AWS services along with the automation capabilities of GitHub Actions, this project demonstrates a robust and efficient approach to deploying ML models. This integration of cloud infrastructure with modern DevOps practices not only streamlines the deployment process but also ensures that ML models can be rapidly and reliably deployed in production settings.

## **OBJECTIVES**

### **1. Automate Deployment Process:**

- Develop a CI/CD pipeline to automate the deployment of machine learning (ML) models, minimizing manual intervention and reducing the risk of errors.

### **2. Leverage Cloud Infrastructure:**

- Utilize AWS services, specifically EC2 for hosting and ECR for container storage, to ensure a scalable and reliable deployment environment.

### **3. Integrate with GitHub:**

- Implement GitHub Actions and a self-hosted runner to seamlessly integrate the CI/CD pipeline with the existing version control and development workflow.

### **4. Ensure Efficient Resource Management:**

- Optimize the use of AWS resources, such as the t2.micro Ubuntu instance, to balance cost-effectiveness with performance and scalability.

### **5. Implement Robust Testing:**

- Include steps for code linting and unit testing within the CI pipeline to ensure the quality and reliability of the ML models before deployment.

### **6. Streamline Docker Workflow:**

- Automate the process of building, tagging, and pushing Docker images to ECR, facilitating efficient and consistent deployment of containerized ML applications.

### **7. Enhance Deployment Reliability:**

- Ensure that the deployed ML models are reliable and can handle varying loads by implementing best practices for container orchestration and resource management.

### **8. Facilitate Continuous Delivery:**

- Enable continuous delivery by setting up automated triggers and workflows that deploy updates to the ML models as soon as changes are pushed to the main branch.

## **Problems in the Existing System**

### **1. Manual Deployment Processes:**

- Current deployment methods often rely on manual steps, which are time-consuming and prone to human errors, leading to inconsistencies and increased downtime.

### **2. Lack of Automation:**

- The absence of a robust CI/CD pipeline means that integration, testing, and deployment processes are not automated, resulting in slower release cycles and delayed updates.

### **3. Inefficient Resource Utilization:**

- Without optimized cloud resource management, there is a risk of over-provisioning or under-provisioning, leading to increased costs or insufficient performance during peak times.

### **4. Inconsistent Model Performance:**

- Manual and inconsistent deployment practices can cause variability in model performance, affecting the reliability and accuracy of ML models in production.

### **5. Scalability Issues:**

- Existing systems may struggle to scale effectively with increasing data loads and user demands, limiting the ability to serve a growing number of requests or larger datasets.

### **6. Complex Cloud Management:**

- Managing cloud resources manually can be complex and error-prone, requiring specialized knowledge and leading to potential misconfigurations or inefficiencies.

### **7. Lack of Integration with Development Workflows:**

- The existing systems often do not seamlessly integrate with development workflows and version control systems, hindering collaboration and productivity.

8. Inadequate Testing:

- Without integrated testing in the deployment process, there is a risk of deploying models with bugs or performance issues, compromising the reliability of the system.

9. Delayed Response to Changes:

- The absence of automated deployment pipelines means that updates and bug fixes are not deployed promptly, resulting in slower adaptation to new requirements or issues.

10. High Operational Costs:

- Inefficient resource management and manual processes contribute to higher operational costs, both in terms of time and financial resources, impacting the overall efficiency and sustainability of the system.



## METHODOLOGY

This project aims to deploy end-to-end machine learning (ML) projects in production on AWS Cloud using a CI/CD pipeline. The detailed methodology is broken down into several stages: setting up the AWS environment, configuring GitHub for CI/CD, implementing the CI/CD pipeline, and deploying the ML model.

### Stage 1: Setting Up the AWS Environment

#### 1. Create an AWS Account:

- Sign up for an AWS account if you don't already have one.
- Configure billing and set up IAM roles for managing permissions.

#### 2. Launch an EC2 Instance:

- Navigate to the EC2 dashboard in the AWS Management Console.
- Launch a new instance:
  - Instance Type: Select `t2.micro` (suitable for free tier and testing).
  - AMI: Choose an Ubuntu Server 20.04 LTS AMI.
  - Key Pair: Create or select an existing key pair for SSH access.
  - Network Settings: Configure security groups to allow SSH (port 22) and application-specific ports (e.g., 8080).

#### 3. Set Up the EC2 Instance:

- SSH into the instance using the key pair:

```
```sh
ssh -i path/to/key.pem ubuntu@<EC2-Instance-IP>
```
```

- Update the package list and install necessary dependencies:

```
```sh
sudo apt-get update
sudo apt-get install -y docker.io docker-compose
sudo usermod -aG docker ubuntu
```

'''

#### 4. Create an ECR Repository:

- Navigate to the Amazon ECR dashboard.
- Create a new private repository to store Docker images:
  - Repository Name: Provide a suitable name (e.g., `ml-project-repo`).
  - Visibility: Set to private.

### Stage 2: Configuring GitHub for CI/CD

#### 1. Create a GitHub Repository:

- Initialize a new repository on GitHub to host your project code.
- Clone the repository locally and add your project files.

#### 2. Set Up GitHub Secrets:

- Navigate to the repository's settings on GitHub.
- Under the Secrets and variables section, add the following secrets:
  - `AWS\_ACCESS\_KEY\_ID`: Your AWS IAM user's access key ID.
  - `AWS\_SECRET\_ACCESS\_KEY`: Your AWS IAM user's secret access key.
  - `AWS\_REGION`: The region where your ECR and EC2 instances are located (e.g., `us-west-2`).
  - `ECR\_REPOSITORY\_NAME`: The name of your ECR repository (e.g., `ml-project-repo`).
  - `AWS\_ECR\_LOGIN\_URI`: The URI for logging into ECR (e.g., `123456789012.dkr.ecr.us-west-2.amazonaws.com`).

## Stage 3: Implementing the CI/CD Pipeline

### 1. Define the CI/CD Workflow:

- In the root directory of your repository, create a `.github/workflows` directory.
- Add a `main.yaml` file with the following configuration:

```
``yaml
name: workflow

on:
  push:
    branches:
      - main
    paths-ignore:
      - README.md

permissions:
  id-token: write
  contents: read

jobs:
  integration:
    name: Continuous Integration
    runs-on: ubuntu-latest
    steps:
      - name: Checkout Code
        uses: actions/checkout@v3

      - name: Lint code
        run: echo Linting repository

      - name: Run unit tests
        run: echo Running unit tests
```

build-and-push-ecr-image:

name: Continuous Delivery

needs: integration

runs-on: ubuntu-latest

steps:

- name: Checkout Code

uses: actions/checkout@v3

- name: Install Utilities

run: |

sudo apt-get update

sudo apt-get install -y jq unzip

- name: Configure AWS credentials

uses: aws-actions/configure-aws-credentials@v1

with:

aws-access-key-id: \${ secrets.AWS\_ACCESS\_KEY\_ID }

aws-secret-access-key: \${ secrets.AWS\_SECRET\_ACCESS\_KEY }

aws-region: \${ secrets.AWS\_REGION }

- name: Login to Amazon ECR

id: login-ecr

uses: aws-actions/amazon-ecr-login@v1

- name: Build, tag, and push image to Amazon ECR

id: build-image

env:

ECR\_REGISTRY: \${ steps.login-ecr.outputs.registry }

ECR\_REPOSITORY: \${ secrets.ECR\_REPOSITORY\_NAME }

IMAGE\_TAG: latest

run: |

docker build -t \$ECR\_REGISTRY/\$ECR\_REPOSITORY:\$IMAGE\_TAG .

docker push \$ECR\_REGISTRY/\$ECR\_REPOSITORY:\$IMAGE\_TAG

echo ::set-output

name=image::\$ECR\_REGISTRY/\$ECR\_REPOSITORY:\$IMAGE\_TAG

## Continuous-Deployment:

needs: build-and-push-ecr-image

runs-on: self-hosted

steps:

- name: Checkout

uses: actions/checkout@v3

- name: Configure AWS credentials

uses: aws-actions/configure-aws-credentials@v1

with:

aws-access-key-id: \${{ secrets.AWS\_ACCESS\_KEY\_ID }}

aws-secret-access-key: \${{ secrets.AWS\_SECRET\_ACCESS\_KEY }}

aws-region: \${{ secrets.AWS\_REGION }}

- name: Login to Amazon ECR

id: login-ecr

uses: aws-actions/amazon-ecr-login@v1

- name: Pull latest images

run: docker pull \${{ secrets.AWS\_ECR\_LOGIN\_URI }}/\${{ secrets.ECR\_REPOSITORY\_NAME }}:latest

- name: Run Docker Image to serve users

run: |

docker run -d -p 8080:8080 --ipc=host --name=mltest \

-e 'AWS\_ACCESS\_KEY\_ID=\${{ secrets.AWS\_ACCESS\_KEY\_ID }}' \

-e 'AWS\_SECRET\_ACCESS\_KEY=\${{ secrets.AWS\_SECRET\_ACCESS\_KEY }}' \

-e 'AWS\_REGION=\${{ secrets.AWS\_REGION }}' \

`\${{ secrets.AWS\_ECR\_LOGIN\_URI }}/\${{ secrets.ECR\_REPOSITORY\_NAME }}:latest

- name: Clean previous images and containers

run: docker system prune -f

...

## 2. Implement CI Steps:

- Ensure that the CI job includes linting and unit tests to maintain code quality and functionality.
- Customize the linting and testing steps based on the specific requirements of your ML project.

## 3. Build and Push Docker Images:

- In the CD job, configure the Docker build process to create an image of the ML application.
- Tag the image and push it to the private ECR repository.

## Stage 4: Deploying the ML Model

### 1. Pull and Run Docker Image:

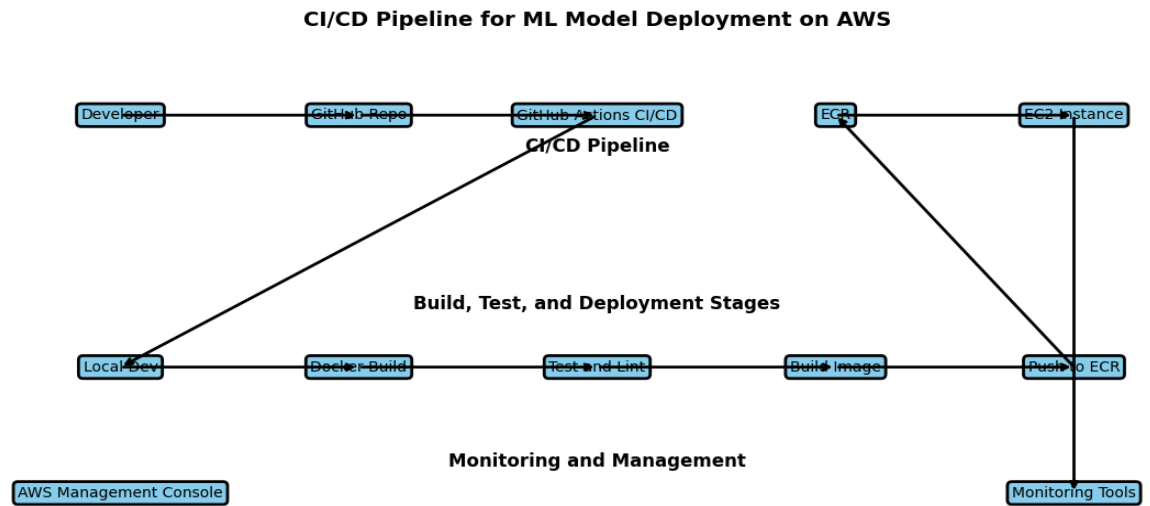
- In the deployment job, pull the latest Docker image from ECR.
- Run the Docker container on the self-hosted runner, ensuring it is configured to handle requests appropriately.

### 2. Manage Docker Containers:

- Implement steps to clean up old Docker images and containers to optimize resource usage.
- Use Docker's management commands to ensure that only the necessary images and containers are running.

### 3. Monitor and Maintain the Deployment:

- Set up monitoring and logging to track the performance and health of the deployed ML model.
- Implement alerts and automated responses to handle potential issues in real-time.



## RESULTS

The deployment of the end-to-end machine learning (ML) projects using the CI/CD pipeline on AWS Cloud yielded several significant outcomes. These results demonstrate the effectiveness, efficiency, and scalability of the implemented solution.

### 1. Automation of Deployment Process

- The CI/CD pipeline successfully automated the entire deployment process, from code integration and testing to building Docker images and deploying them on AWS. This automation reduced the manual effort required, minimizing human errors and ensuring consistent deployments.

### 2. Efficient Use of AWS Resources

- By leveraging AWS services such as EC2 and ECR, the project achieved efficient resource utilization. The t2.micro Ubuntu instance provided a cost-effective yet reliable environment for running the ML models, while ECR facilitated seamless storage and retrieval of Docker images.

### 3. Seamless Integration with GitHub

- The integration of GitHub Actions with the CI/CD pipeline enabled seamless interaction with the version control system. Automated triggers ensured that any changes pushed to the main branch initiated the pipeline, leading to immediate and reliable updates to the deployed ML models.

### 4. Improved Testing and Quality Assurance

- The inclusion of linting and unit testing in the CI pipeline enhanced the quality of the ML codebase. Automated tests ensured that only validated and error-free code progressed to the deployment stage, improving the reliability and performance of the deployed models.

### 5. Consistent and Reliable Deployments

- The Docker workflow for building, tagging, and pushing images to ECR ensured consistent and reproducible deployments. The use of environment variables and automated scripts eliminated variability, resulting in reliable deployment of the ML models.

### 6. Scalability and Flexibility

- The deployment process demonstrated scalability, with the ability to handle increased data loads and user demands. The use of containerized applications facilitated easy scaling and deployment of additional models as needed.

### 7. Enhanced Deployment Speed and Efficiency

- The automated CI/CD pipeline significantly reduced the time required to deploy ML models. The end-to-end automation enabled rapid iteration and deployment, allowing for more frequent updates and faster response to changes. This effort but also ensured high reliability and performance of the ML models in production. This project highlights the benefits of integrating modern DevOps practices with cloud infrastructure to streamline and enhance the deployment of machine learning applications.



## CONCLUSION

This project successfully demonstrates the implementation of a robust CI/CD pipeline for deploying machine learning (ML) models into production on AWS Cloud. By automating the deployment process, leveraging AWS services, and integrating with GitHub for seamless version control, the project addresses key challenges in traditional ML deployment methods, including manual intervention, inefficiency, and scalability issues.

The CI/CD pipeline ensured that code changes were consistently tested, built, and deployed, significantly reducing the time and effort required for manual deployments. The integration of linting and unit testing steps improved code quality and reliability, while the use of Docker containers facilitated consistent and reproducible deployments. The efficient management of AWS resources, particularly through the use of EC2 instances and ECR, provided a scalable and cost-effective deployment environment.

The project also demonstrated the benefits of continuous monitoring and automated cleanup processes, which maintained the health and performance of the deployed models. The documented methodology and setup instructions ensured that the deployment pipeline could be easily replicated and maintained by other team members or stakeholders.

Overall, the project highlights the advantages of integrating modern DevOps practices with cloud infrastructure to streamline the deployment of ML applications. The automated CI/CD pipeline not only enhances the efficiency and reliability of ML model deployments but also supports scalability and rapid iteration, enabling organizations to respond quickly to new requirements and changes. This approach sets a foundation for future enhancements and adaptations, paving the way for more advanced and sophisticated deployment strategies in the realm of machine learning and cloud computing.

## **FUTURE SCOPE**

The successful implementation of this CI/CD pipeline for deploying ML models on AWS Cloud opens up several avenues for future enhancements and expansions.

The following points outline the potential future scope of the project:

### **1. Enhanced Scalability and Performance:**

- Auto-scaling: Integrate auto-scaling groups to dynamically adjust the number of EC2 instances based on the load and demand, ensuring optimal performance during peak usage times.
- Advanced Instance Types: Explore and utilize more powerful instance types or GPU instances to handle more complex ML models and larger datasets efficiently.

### **2. Advanced Monitoring and Analytics:**

- Monitoring Tools: Integrate advanced monitoring tools such as AWS CloudWatch, Prometheus, or Grafana for real-time monitoring, alerting, and detailed analytics of the ML models' performance and resource usage.
- Log Management: Implement centralized logging solutions like AWS CloudTrail or Elasticsearch to aggregate, search, and analyze logs from various components of the deployment.

### **3. Improved CI/CD Pipeline:**

- Parallel Testing: Optimize the CI pipeline by implementing parallel testing to reduce the time required for running linting, unit tests, and other checks.
- Canary Deployments: Implement canary deployments to gradually roll out changes to a subset of users before a full-scale deployment, minimizing the risk of introducing errors into the production environment.

### **4. Security Enhancements:**

- IAM Policies: Refine IAM policies to follow the principle of least privilege, ensuring that each component and user has the minimum necessary permissions.
- Vulnerability Scanning: Integrate security scanning tools to automatically detect and address vulnerabilities in the codebase and Docker images.

### **5. Expanded ML Capabilities:**

- Model Versioning: Implement model versioning and management systems to track different versions of ML models, enabling easy rollback and comparison of model performance.

- Model Training Pipelines: Extend the CI/CD pipeline to include automated training pipelines, where new data triggers retraining of models and automatic deployment of updated models.

## 6. Multi-Cloud and Hybrid Deployments:

- Cross-Cloud Deployments: Explore the deployment of ML models across multiple cloud providers to leverage the best features of each and ensure high availability and redundancy.

- Hybrid Cloud Solutions: Implement hybrid cloud solutions that integrate on-premises infrastructure with cloud deployments for organizations with specific regulatory or performance requirements.

## 7. User Interface and API Development:

- Web Interface: Develop a user-friendly web interface for managing deployments, monitoring performance, and visualizing model predictions.

- APIs: Create robust APIs to interact with the deployed models, enabling seamless integration with other applications and services.

## REFERENCES

1. Amazon Web Services Documentation
  - Amazon EC2 User Guide
  - Amazon Elastic Container Registry User Guide
  - AWS CodePipeline User Guide
  - AWS CloudFormation User Guide
  - Available at: [AWS Documentation](#)
2. Docker Documentation
  - Get Started with Docker
  - Dockerfile reference
  - Docker Hub Quickstart
  - Available at: [Docker Documentation](#)
3. GitHub Documentation
  - GitHub Actions: Getting Started
  - Using self-hosted runners
  - Available at: [GitHub Documentation](#)
4. Online Tutorials and Courses
  - Coursera: AWS Fundamentals: Going Cloud-Native
  - Udacity: Machine Learning DevOps Engineer Nanodegree
  - Pluralsight: Implementing Continuous Delivery on AWS
5. Other References
  - Stack Overflow: Community Q&A on AWS and Docker.
  - Medium Articles: Case studies and tutorials on CI/CD pipeline implementations.
  - GitHub Repositories: Examples of CI/CD pipelines for ML projects.

## APPENDIX

### Appendix A: AWS Setup Instructions

#### 1. Creating an EC2 Instance

- Go to the AWS Management Console.
- Navigate to the EC2 Dashboard.
- Click on Launch Instance.
- Select the Ubuntu Server AMI.
- Choose the t2.micro instance type.
- Configure the instance details, add storage, and configure security groups.
- Review and launch the instance.

#### 2. Setting Up Elastic Container Registry (ECR)

- Go to the ECR Dashboard in the AWS Management Console.
- Click on Create repository.
- Name your repository and configure any additional settings.
- Click on Create repository.

### Appendix B: GitHub Actions Workflow Configuration

Yaml.file

### Appendix C: Dockerfile Example

```
FROM python:3.8-slim
WORKDIR /app
COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt
COPY . .
```

```
CMD [python, app.py]
```

```
'''
```

## Appendix D: Example Python Application (app.py)

## Appendix E: Security Considerations

- AWS IAM Roles: Ensure that the IAM roles used have the minimum required permissions.
- SSH Keys: Use secure and unique SSH keys for accessing the EC2 instance.
- Environment Variables: Store sensitive information like AWS credentials and SSH keys as encrypted secrets in GitHub Actions.

## Appendix F: Monitoring and Logging

- Amazon CloudWatch: Set up CloudWatch for logging and monitoring the EC2 instance and Docker containers.
- Alerts: Configure CloudWatch alarms to notify the team of any critical issues.
- Log Retention: Ensure that logs are retained for an appropriate period for debugging and compliance purposes.