



COMP4964 Data Pipeline CI/CD Automation

A Demonstration of Seamless Integration and Deployment for Data
Workflows

Group 5

Neriyel Reyes

Dipenvir Kaur

Automated Data Pipeline: Project Overview

This project showcases a robust CI/CD pipeline designed for data processing, ensuring reliability and efficiency from code commit to deployment.

Automated Jenkins Pipeline

Hosted securely on DigitalOcean, orchestrating the entire CI/CD workflow.

GitHub Webhook Integration

Triggers automatic builds upon every code push, maintaining continuous integration.

AWS Lambda Processing

Serverless function designed to efficiently process CSV files from S3 buckets.

Automated Unit

Tests code quality and prevents deployment of broken or faulty code.



Key Technologies Driving Our Automation

Our solution leverages industry-leading tools to create a resilient and scalable data pipeline.

CI/CD Orchestration: Jenkins

The heart of our continuous integration and continuous delivery system, Jenkins automates the build, test, and deployment phases.

Version Control: GitHub

Manages our codebase and integrates webhooks for real-time pipeline triggers.

Serverless Compute: AWS Lambda

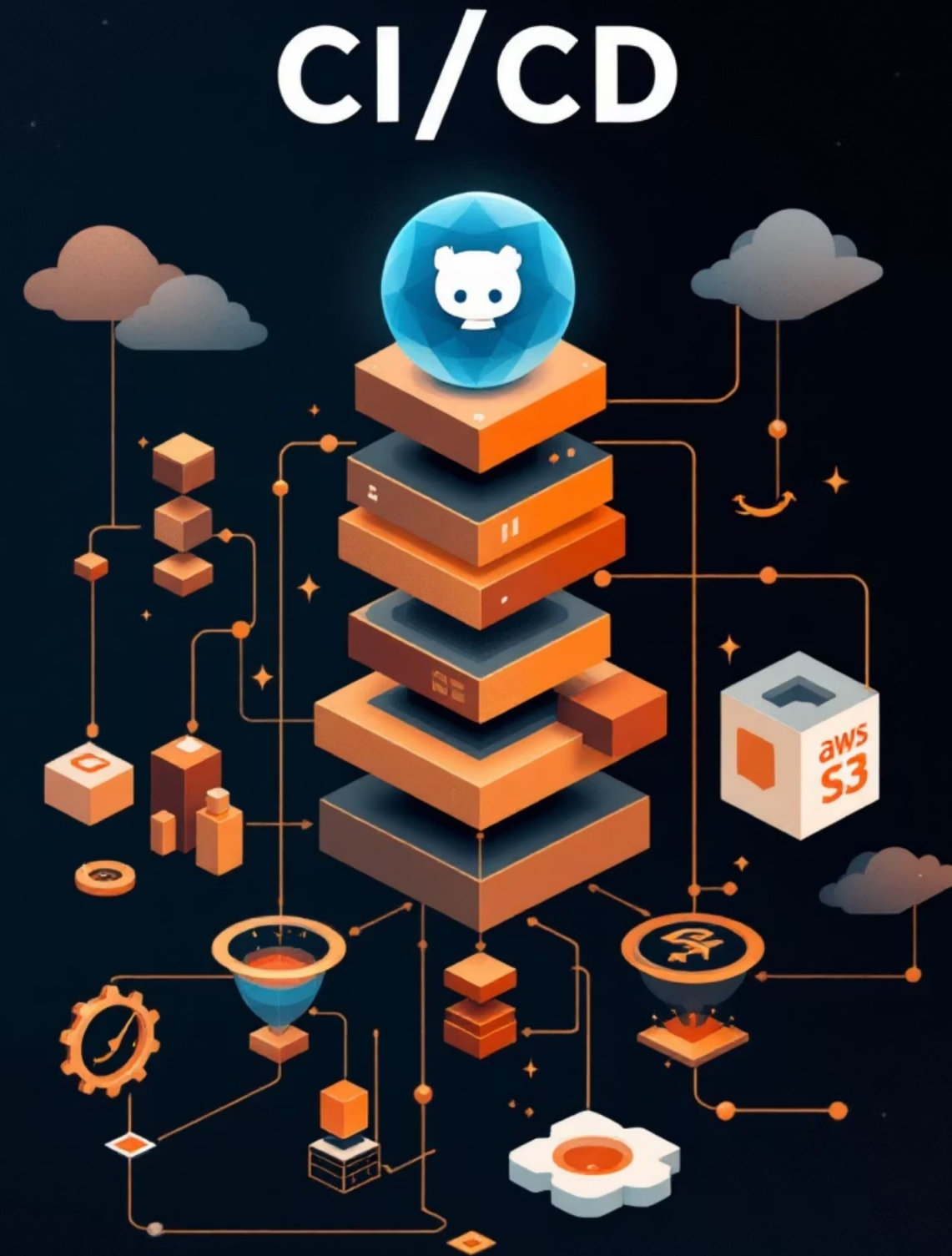
Executes our data processing logic without server provisioning or management.

File Storage: AWS S3

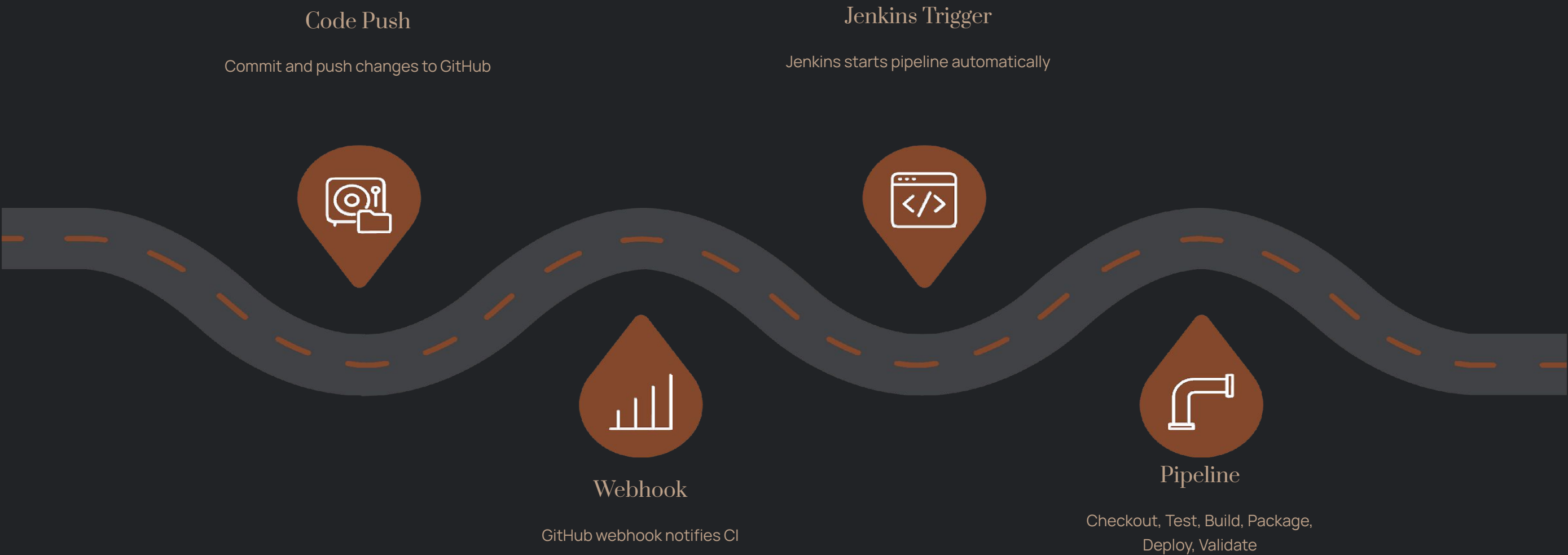
Provides secure, scalable, and highly available storage for our raw and processed data files.

Infrastructure as Code: CloudFormation

Defines and provisions all AWS infrastructure resources using a declarative template, ensuring consistency and repeatability.



End-to-End Pipeline Architecture



Jenkins File Explained

```
5 pipeline {
7
8   environment {
9     AWS_REGION = 'us-west-2'
10    STACK_NAME = 'data-pipeline-stack'
11    ENVIRONMENT = 'dev'
12    S3_DEPLOY_BUCKET = 'jenkins-sam-artifacts-195275680578-20251125123539'
13    AWS_DEFAULT_REGION = 'us-west-2'
14  }
15
16  stages {
17    stage('Checkout') {
18      steps {
19        echo 'Checking out code from GitHub...'
20        checkout scm
21      }
22    }
23
24    stage('Test') {
25      steps {
26        echo 'Running unit tests for Lambda function...'
27        sh '''
28          cd jenkins-data-pipeline
29          python3 -m pytest tests/test_lambda.py -v
30        '''
31      }
32    }
33
34    stage('Build') {
35      steps {
36        echo 'Validating SAM template...'
37        sh '''
38          echo "Checking template.yaml exists..."
39          ls -la jenkins-data-pipeline/template.yaml
40          echo "Template validated!"
41        '''
42    }
```

```
45   stage('Package') {
46     steps {
47       echo 'Preparing for deployment...'
48       sh '''
49         echo "Using pre-built template..."
50         cp jenkins-data-pipeline/template.yaml ./packaged.yaml
51         echo "Ready to deploy!"
52       '''
53     }
54   }
55
56   stage('Deploy') {
57     steps {
58       echo 'Deploying to AWS...'
59       sh '''
60         echo "=====
61         echo "CloudFormation Stack Deployment"
62         echo "=====
63         echo "Stack Name: ${STACK_NAME}"
64         echo "Region: ${AWS_REGION}"
65         echo "Environment: ${ENVIRONMENT}"
66         echo "Template: packaged.yaml"
67         echo ""
68         echo "Template size: $(wc -c < packaged.yaml) bytes"
```

```
19   post {
20     success {
21       echo 'Pipeline completed successfully!'
22       echo 'Your data pipeline is ready for use.'
23     }
24     failure {
25       echo 'Pipeline failed. Check logs above.'
26     }
27   }
28 }
29 }
```


Demo 1: Webhook Automation in Action

The Scenario:

Add a minor comment to

`lambda_handler.py`.

- Commit and push the changes to our GitHub repository.
- Observe the Jenkins dashboard: a new build automatically commences.

Verify the Jenkins logs show "**Started by GitHub push**".

Behind the Scenes:

GitHub dispatches a webhook event to our Jenkins instance at

`jenkin.neriyelreyes.org/github-webhook/`.

- Jenkins receives the notification, detects the code change, and immediately triggers the configured pipeline.
- This automation eliminates manual triggers, streamlining the development workflow.





Demo 2: CSV File Automated Cleanup



The Flow:

Upload `demo_data_messy.csv` to S3

- An S3 event notification automatically triggers the AWS Lambda function.
- The Lambda function cleans the CSV data (trims whitespace, removes duplicates).

The processed file appears in a separate S3.

- This demonstrates a second layer of robust automation, complementing the CI/CD pipeline.



Demo 3: Catching Bugs with Test Stage Validation

This demonstration highlights the critical role of automated testing, ensuring that faulty code is identified and halted before deployment, maintaining system integrity.

1

Simulating a Bug

Intentionally remove a key function (e.g., `process_csv_data`) from `lambda_handler.py` to create a critical error.

2

Automatic Trigger & Failure

Commit and push the "broken" code. Jenkins will automatically trigger, but the **Test stage will FAIL** due to `pytest` not being able to import the missing function.

3

The Safety Net

This failure prevents the deployment of non-functional code, showcasing the effectiveness of the CI/CD pipeline's testing phase as a crucial safety net.

Key Takeaways: The Power of Automation



GitHub Webhook

Enables instant pipeline triggers on every code push for continuous integration.



Jenkins Pipeline

Provides multi-stage automated testing and deployment for reliable releases.



Unit Tests

Crucial for catching bugs early in the **Test** stage, preventing broken code deployment.



S3 Triggers & Lambda

Automates data processing upon file upload, ensuring immediate and clean data availability.

salamat

ਧੰਨਵਾਦ

Any Questions?

