**aws**

**SUMMIT
ONLINE**

ASEAN

# Using IaC with Terraform to Provision Big Data Platform on Amazon EMR

**Dwi Fahni Denni**

**Lead DevOps Engineer (ZebraX)**
**AWS Community Builders - Indonesia**

# Agenda

- **Architecture Big Data Platform**

- **GitOps Flow**

- **Demo**

  - **Deploy Amazon EMR Cluster**
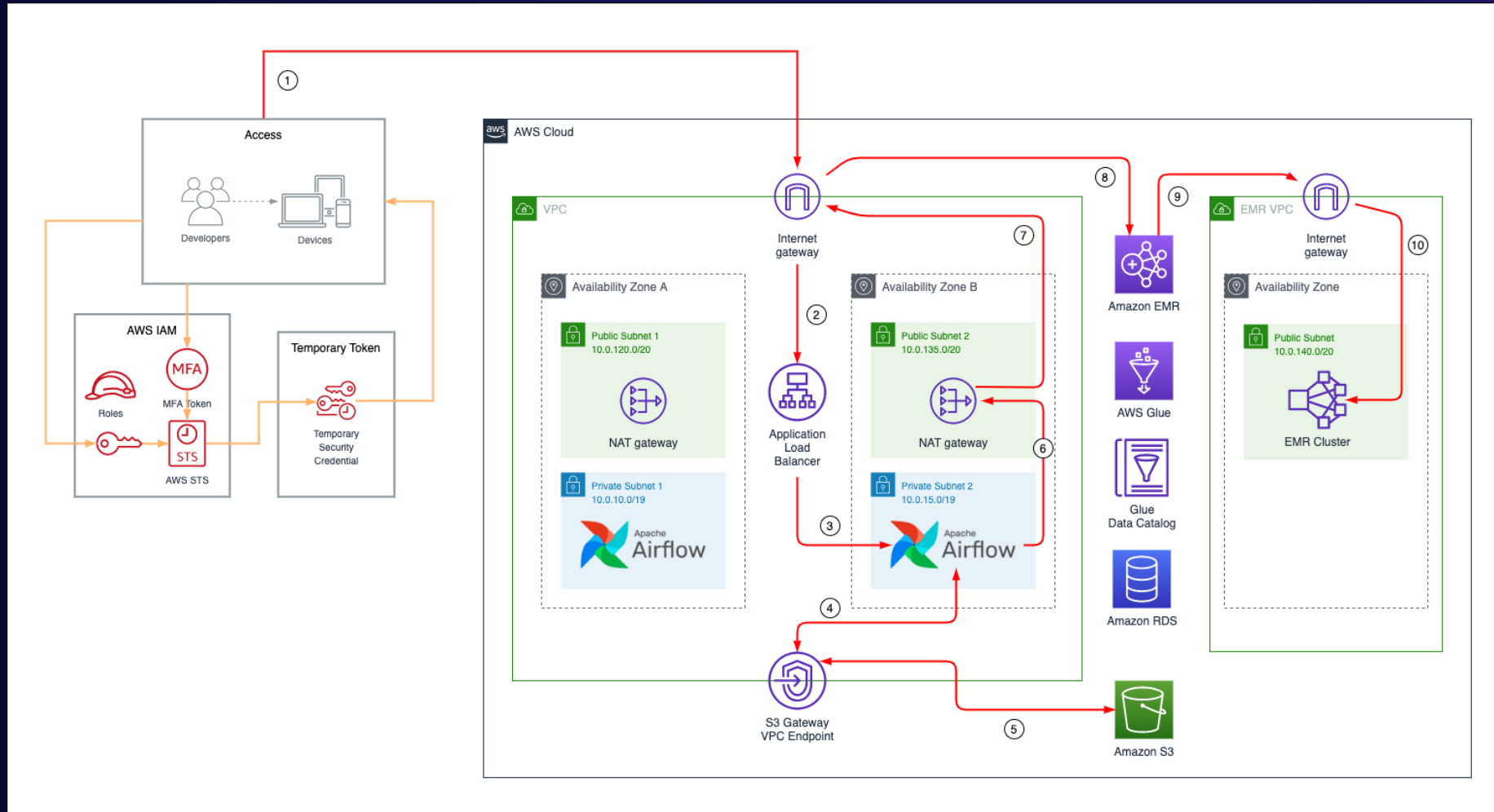
# Abstract

**USING IAC WITH TERRAFORM TO PROVISION BIG DATA PLATFORM ON AMAZON EMR**

- Train multiple ML Models in a serial mode it could be challenging, time consuming, and not effective. We need to create parallel modelling which provides tremendous benefit in building a variety of models by speeding up the process through parallelization so the model building process becomes more efficient.

- In this presentation, we will solving our problem in training multiple ML models using Spark Panda's UDF (a python libraries for building ML models) inside Amazon EMR cluster.

- We will also learn how to provisioning Amazon EMR cluster with Terraform as Infrastructure-as-Code tools.

**Chapter 1:**

# Architecture Big Data Platform
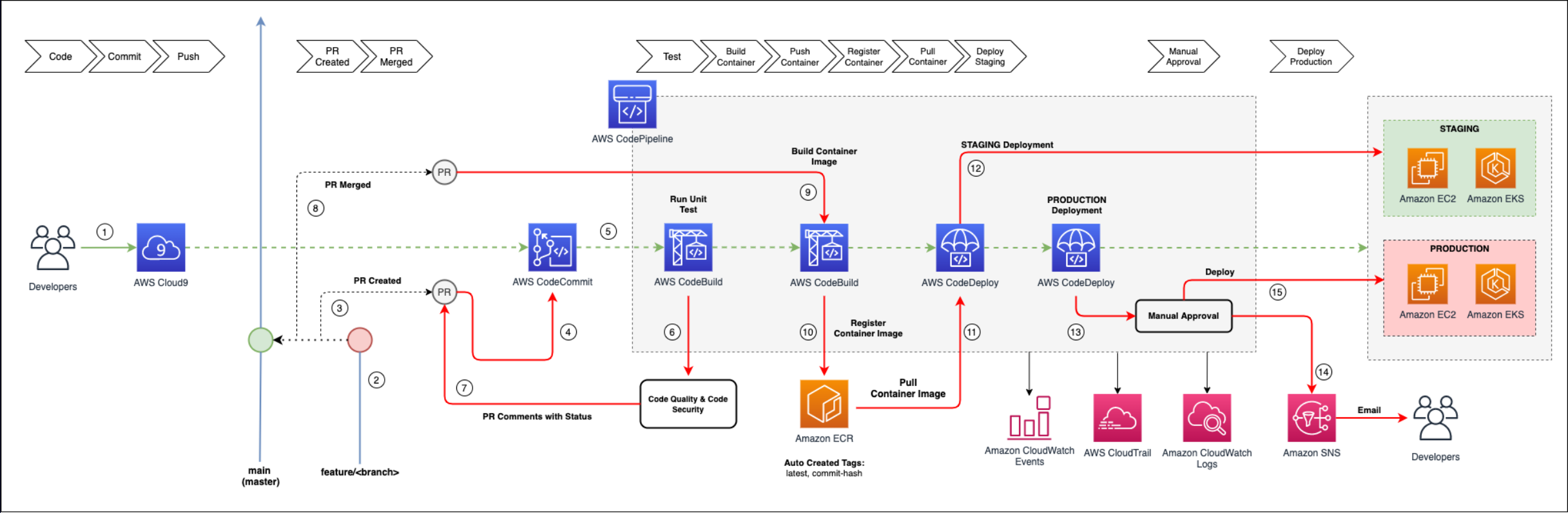
# Architecture Big Data Platform

**Chapter 2:**
# GitOps Pipeline

# GitOps Pipeline
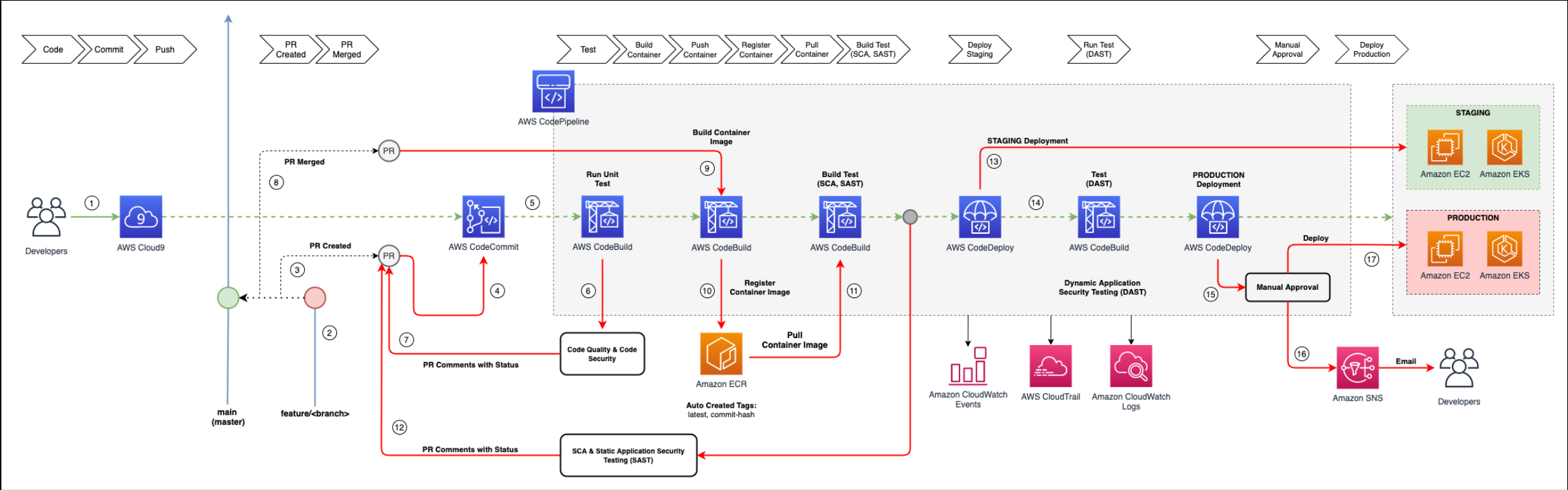
- **GitOps Flow**

- **CI/CD Pipeline Deployment**

# GitOps Pipeline

GITOPS DEVSECOPS FLOW
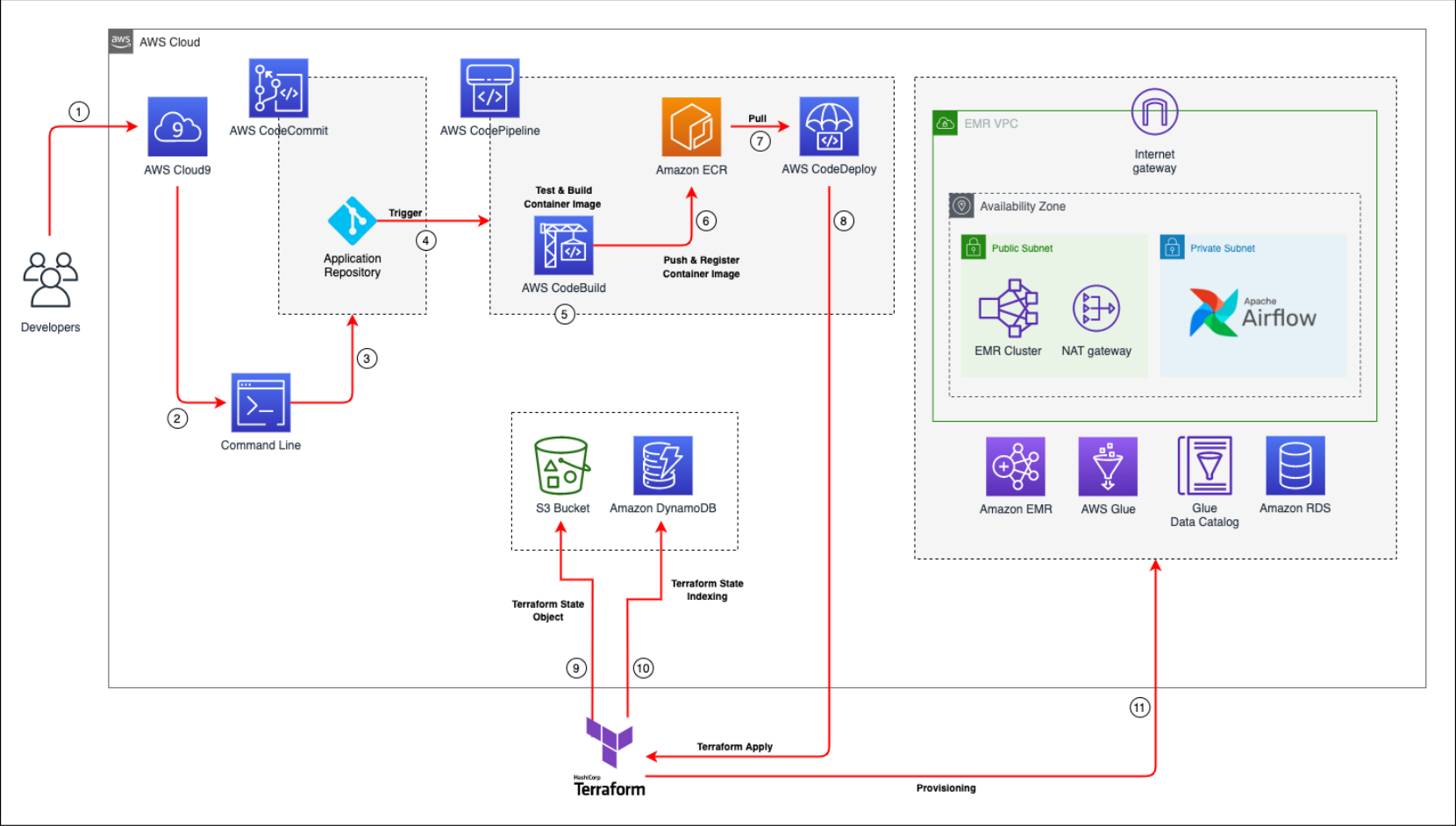
# GitOps Pipeline

Workflow of CI/CD pipeline in provisioning Amazon EMR Cluster Terraform

**Chapter 3:**

# Demo

# Provisioning Amazon EMR

- **Preparations**

- **Deploy Amazon EMR**

- **Bootstrap Script & Terraform Stack**

- **Running Spark Job**

# Provisioning Amazon EMR

- **IAM Role, Policy (AssumeRole)**

- **Infrastructure**

  - **VPC, Subnet, NAT, Internet-Gateway, DNS**

- **S3 Bucket**

  - **Bucket for Terraform State**

  - **Bucket for EMR Bootstrap & EMR Configuration**

- **DynamoDB**

  - **Indexing Terraform State**

- **Database (RDS) –** *optional*

# Provisioning Amazon EMR

# Provisioning Amazon EMR

# Deploy Amazon EMR

**(ETA: 30 minutes)**

# Provisioning Amazon EMR

Bootstrap script including:

- IAM (Role & Policy)

- Provisioning Master, Core & Task Instance Fleet

- Instance Type Provisioning

- Volume (EBS) Size Instance

- Hadoop Debugging

- Autoscaling

- Monitoring & Log

# Provisioning Amazon EMR

- Airflow is a useful tool to monitor workflow as Directed Acyclic Graphs (DAGs) of tasks

- The scheduler for automation pipeline and monitoring Spark job (python script) will be running under Airflow

**Appendix:**
# References

# References

| Resources | Links |
|---|---|
| Docker Container CI/CD | https://github.com/devopscorner/devopscorner-container |
| User Data Installer Scripts | https://github.com/devopscorner/scripts |
| IaC Terraform EMR | https://github.com/devopscorner/iac-terraform-emr |
| Big Data and Machine Learning | https://devopscorner.id/category/machine-learning/ |
| How to Efficiently Train Multiple ML Models on a Spark Cluster | https://medium.com/zebrax/how-to-efficiently-train-multiple-ml-models-on-a-spark-cluster-7d84512d36f0 |

# Thank you!

**Dwi Fahni Denni**

**Lead DevOps Engineer (ZebraX)**
**AWS Community Builders - Indonesia**

https://www.linkedin.com/in/dfdenni

# Learn in-demand AWS Cloud skills

## AWS Skill Builder

Access **500+ free** digital courses and Learning Plans

Explore resources with a variety of skill levels and **16+** languages to meet your learning needs

Deepen your skills with digital learning on demand

Train now

## AWS Certifications

Earn an industry-recognized credential

Receive Foundational, Associate, Professional, and Specialty certifications

Join the **AWS Certified community** and get exclusive benefits

Access **new** exam guides

# Thank you for attending **AWS Summit Online ASEAN 2022**.

Please **complete the session survey** to help us improve your Summit experience in the future.

✉ aws-asean-marketing@amazon.com

🐦 twitter.com/AWSCloudSEAsia

in linkedin.com/company/amazon-web-services

f facebook.com/AmazonWebServices

📷 instagram.com/amazonwebservices

▶ youtube.com/user/AmazonWebServices

🎮 twitch.tv/aws