# P. E. S. College of Engineering, Mandya

## Department of Information Science and Engineering

An Internship Presentation on

# Data Engineering Virtual Internship

**Submitted by:**
**Rakesh S**
**4PS21IS045**
**Dept. of ISE,**
**PESCE, Mandya**

**Under the Guidance of**
**M R Suresh**
**Associate Professor,**
**Dept. of ISE,**
**PESCE, Mandya**

# Contents

- **Introduction**
- **About the Course**
- **Key Learnings**
- **Hands-on Work**
- **Project**
- **Conclusion**

# **Introduction**

- I completed a 10-week Data Engineering Virtual Internship from January to March 2025.

- The internship was organized by AICTE (All India Council for Technical Education) in collaboration with EduSkills and AWS Academy.

- This program was delivered through the National Internship Portal and was open to engineering students across India.

- The curriculum provided in-depth exposure to AWS Cloud and Data Engineering concepts, hands-on labs, and real-world projects.

- The internship enabled me to gain practical skills in cloud computing, data pipeline design, and AWS services.

- Certificate issued by AICTE, EduSkills, and AWS Academy.

# About the course

The internship combined two in-depth AWS Academy programs:

- **AWS Academy Cloud Foundations**

- **AWS Academy Data Engineering**

## Course Objectives:

- Provided a foundational understanding of cloud computing and AWS services.

- Covered essential data engineering concepts and practical skills for designing data pipelines on AWS.

- Prepared students for AWS certification exams and industry roles.

# Cloud Service Models

1. **Infrastructure as a Service (IaaS):**
   - Provides the basic building blocks for cloud IT: networking, virtual computers, and storage.
   - Offers the highest level of flexibility and management control over IT resources.
   - Most similar to traditional IT resources familiar to many IT departments and developers.

2. **Platform as a Service (PaaS):**
   - Reduces the need to manage underlying infrastructure (hardware, OS).
   - Focuses on deploying and managing applications, not servers or storage.

3. **Software as a Service (SaaS):**
   - Delivers a complete product managed by the provider, typically end-user applications.
   - Users only need to consider how to use the software, not how it is maintained or hosted (e.g., web-based email)

# Cloud Deployment Models

1. **Cloud:**
   - Applications are fully deployed and run in the cloud, either built natively or migrated from existing infrastructure.
   - Can use low-level infrastructure or higher-level managed services.

2. **Hybrid:**
   - Connects cloud-based resources with on-premises infrastructure.
   - Enables organizations to extend and integrate their IT resources.

3. **On-premises (Private Cloud):**
   - Resources are deployed in-house, often using virtualization.
   - Offers dedicated resources but lacks many cloud benefits; sometimes called private cloud

# Key Advantages of Cloud Computing

- Trade capital expense for variable expense

- Massive economies of scale

- Stop guessing capacity

- Increase speed and agility

- Stop spending on running and maintaining data centers

- Go global in minutes

# Web Services and AWS Overview

**Web Service:**

- Software accessible over the internet using standardized formats (XML, JSON) for API interactions.

**What is AWS?**

- A secure cloud platform offering a broad set of global cloud-based products (compute, storage, database, networking, etc.).
- Provides on-demand access, flexibility, and pay-as-you-go pricing.
- Services are modular and work together like building blocks

**Accessing AWS Services**

- AWS Management Console
- AWS Command Line Interface (CLI)
- AWS Software Development Kits (SDKs)

# AWS Global Infrastructure Overview

**AWS Regions**

- An AWS Region is a distinct geographic area with multiple, isolated locations known as Availability Zones.

- AWS operates Regions worldwide, each designed to provide fault tolerance and stability by isolating them from one another.

- Data stored in one Region is not automatically replicated to another; cross-region replication must be configured by the user.

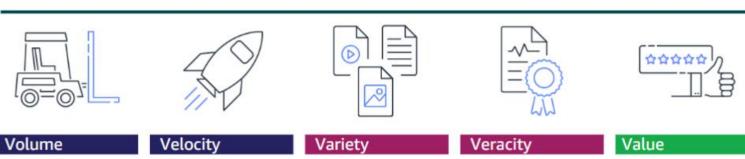- Some Regions, such as AWS GovCloud (US) and China Regions, have special access or compliance requirements

**Availability Zones (AZs)**

- Each Region consists of two or more Availability Zones, which are separate physical locations with independent power, networking, and cooling.

- An Availability Zone may contain one or more data centers, and customers deploy resources at the AZ level, not the data center level.

# The Elements of Data

The five Vs of data, which map to the questions that a data engineer must answer to design a good data infrastructure.

**Data characteristics that drive infrastructure decisions**

| Volume | Velocity | Variety | Veracity | Value |
|---|---|---|---|---|
| How big is the dataset? How much new data is generated? | How frequently is new data generated and ingested? | What types and formats? How many different sources does the data come from? | How accurate, precise, and trusted is the data? | What insights can be pulled from the data? |

# Ingesting and Preparing Data

**Data Pipeline Layers**

- Ingestion: Extract data from external sources and bring it into the pipeline.

- Temporary Storage: Hold data for initial processing or staging.

- Processing/Transformation: Convert data into usable formats, either before or after loading into permanent storage.

**Traditional ETL (Extract, Transform, Load)**

Best for structured data and traditional analytics.
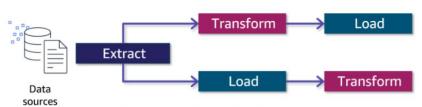
Data is cleaned and formatted before storage, ready

for immediate analysis.

**Modern ELT (Extract, Load, Transform)**

Ideal for handling large volumes of unstructured or

semi-structured data. Enables flexible, on-demand transformations for diverse analytics needs.



**ETL and ELT**

**Extract transform load (ETL)**
1. **Extract** structured data.
2. **Transform** the data into a format that matches the destination.
3. **Load** the data into structured storage for defined analytics.

**Extract load transform (ELT)**
1. **Extract** unstructured or structured data.
2. **Load** the data into the data lake in the format that is as close to the raw form as possible.
3. **Transform** the data as needed for analytic scenarios.

# Storing and Organizing Data

## Types of cloud storage

### Block storage

- Offers dedicated, low-latency storage

- Is scalable and offers high performance

- Is similar to local direct attached storage or a storage area network (SAN)

- Example: Amazon Elastic Block Storage (Amazon EBS)

### File storage

- Stores data as files

- Is highly scalable

- Is ideal for storage such as content repositories and media stores

- Example: Amazon Elastic File System (Amazon EFS)

### Object storage

- Stores unstructured, semistructured, or structured data

- Is highly scalable

- Uses a unique identifier for each object

- Has a lower cost than traditional storage

- Example: Amazon Simple Storage Service (Amazon S3)

# Processing Big Data

## Types of data processing

| Batch data processing | Streaming data processing |
| --- | --- |
| • Infrequently accessed (cold) data querying | • Frequently accessed (hot) data querying |
| • Processes input data in batches at varying intervals | • Processes data sequentially and incrementally in near real time |
| • Tolerates structured and unstructured data | • Capable of processing less predictable data on a massive scale |
| • Capable of deep analysis of big datasets | • Enables analysis of continually generated data |
| • Examples: Amazon EMR, Apache Hadoop | • Examples: Amazon Kinesis Data Streams, Apache Spark Streaming |

# Processing Data for ML

**Three Layers of AWS ML Infrastructure**

1. **Compute, Networking, and Storage Layer**
   1. Compute: EC2 P3/P4 (training), G4/Inf1 (inference), Elastic Inference (GPU acceleration), Elastic Fabric Adapter (low-latency networking)
   2. Storage:
      1. Amazon S3: Scalable object storage for ML datasets
      2. Amazon EBS: High-performance block storage
      3. Amazon EFS/FSx: File system storage for large datasets and shared code

2. **Framework Layer**
   1. Supports popular ML frameworks via Deep Learning AMIs and Deep Learning Containers.

3. **Workflow Services Layer**
   1. Amazon SageMaker: Integrated ML development, training, and deployment
   2. Amazon EMR: Big data processing for ML
   3. AWS Batch: Orchestrates scheduled ML training jobs
   4. Amazon ECS/EKS: Container orchestration for scalable ML environments
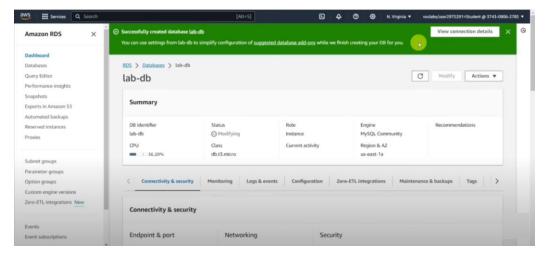   5. AWS ParallelCluster: Simplifies HPC cluster deployment

# **Automating the Pipeline**

- Benefits of Automation
  - Improves stability, consistency, and efficiency by reducing manual intervention
  - Enables rapid response to changes and failures (e.g., automatic server replacement via CloudWatch)
  - Built-in AWS monitoring and automation tools support proactive management

- Infrastructure as Code (IaC)
  - Use AWS CloudFormation or AWS CDK to define and deploy infrastructure in code
  - Repeatability: Deploy identical environments (dev, test, prod) reliably
  - Supports consistent, scalable, and testable infrastructure for analytics and ML workloads

- Orchestration with AWS Step Functions
  - Automate and visualize complex workflows (e.g., ETL pipelines)
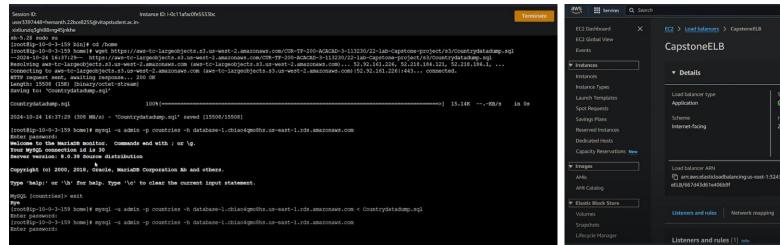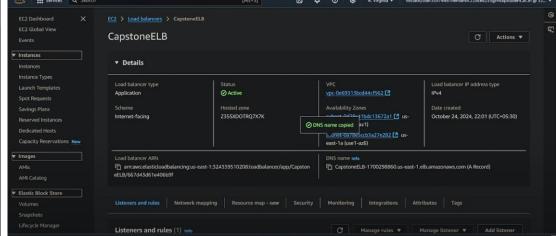  - Integrate with services like Amazon Athena for hands-on lab scenarios

# Hands-on work
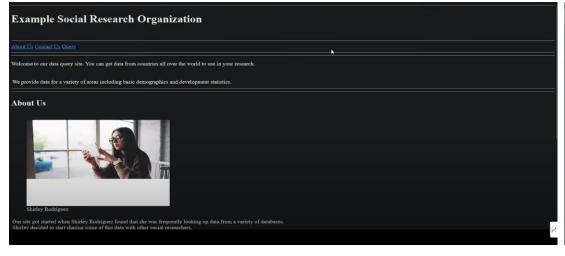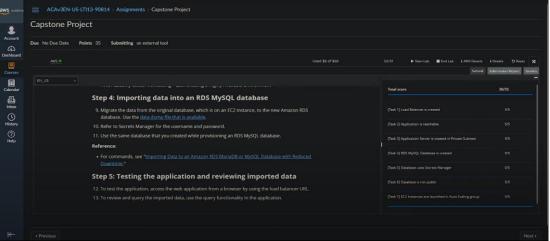


**Lab - 2 Build your VPC and Launch a Web Server**



**Lab - 5 Build a Database Server**

# Project: Capstone project

# Certification

# Thank you..