# Serverless Batch Big Data Processing Pipeline using AWS

## Introduction

This document presents a serverless big data processing pipeline implemented using AWS Free Tier services. The project demonstrates how raw transactional data can be stored, processed, and analyzed efficiently without managing servers.
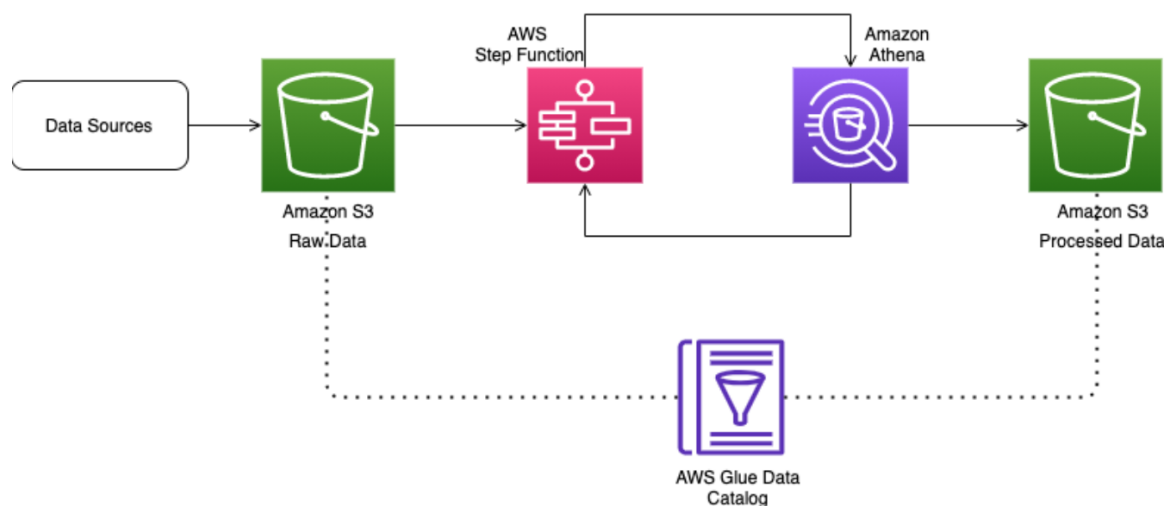
## Problem Statement

Organizations generate large volumes of raw data that is difficult to analyze directly. Traditional solutions are costly and complex. This project addresses the need for a low-cost, scalable analytics pipeline.

## Objectives

- Design a serverless big data pipeline using AWS Free Tier
- Store raw data in Amazon S3
- Process data using AWS Glue
- Convert JSON data to Parquet format
- Analyze data using Amazon Athena

## System Architecture

Local Python Script → Amazon S3 (Raw Data) → AWS Glue Crawler → AWS Glue ETL Job → Amazon S3 (Processed Data) → Amazon Athena

## AWS Services Used

Amazon S3 – Data lake storage

AWS Glue – Schema detection and ETL processing

AWS Glue Data Catalog – Metadata management

Amazon Athena – SQL-based analytics

IAM – Security and permissions

## Implementation Details

1. Data generated using Python in JSON format
2. Uploaded to S3 raw data bucket
3. Glue Crawler detects schema
4. Glue ETL converts data to Parquet
5. Athena queries processed data

Data Stores

Amazon S3  Amazon RDS  Amazon Redshift  Amazon DynamoDB

JDBC  Glue Data Catalog

Infer Schemas

④

Connects to data store ③

Crawler ①

Runs

Custom Classifiers

Connection

②

Built-in Classifiers

Writes Metadata ⑤

Data Catalog

Database  Database

Tables (Metadata)  Tables (Metadata)

## Sample Queries

SELECT SUM(price) FROM sales;
SELECT country, SUM(price) FROM sales GROUP BY country;

## Conclusion

This project successfully demonstrates a complete serverless big data pipeline using AWS services.