# **Project Milestone 6**

## LogiStream: DataCo's Live Flow

## IE 6760 Data Warehousing & Integration

# Group 4

Student 1: Neha Patil

Student 2: Medhavi Uday Pande

Percentage of Effort Contributed by Student 1: 50%

Percentage of Effort Contributed by Student 2: 50%

Signature of Student 1: Neha Patil

Signature of Student 2: Medhavi Uday Pande

# Problem Statement

DataCo is a global e-commerce company with a complex, high-volume supply chain. While it possesses a rich historical dataset of past orders (the 180,000+ records in the CSV), its analytics are purely retrospective. Management can only analyze what went wrong (e.g., late deliveries, fraud) weeks after the fact.

The company has no real-time visibility into its operations. It cannot proactively identify a shipment that is at risk of being late, nor can it detect fraudulent order patterns as they occur. The core problem is the inability to ingest, process, and analyze high-velocity transactional and shipping data as it is generated.

This project aims to design a cloud data pipeline to solve this. The pipeline will be responsible for ingesting live data, transforming it in-stream, and combining it with historical and geospatial reference data to power real-time operational dashboards and alerts.

# Problem Definition

This project will demonstrate a modern data pipeline by ingesting data with significant volume and variety. The pipeline must handle four distinct types of data sources:

- Historical Transactional Data (Initial Batch Load): The DataCoSupplyChainDatasetRefined.csv file, serving as a one-time bulk load of ~180,000+ historical orders to provide historical context.
- Real-time Order Stream (Live Pipeline Feed): A simulated high-velocity stream (e.g., from an e-commerce API) generating JSON events for new orders as they are placed.
- Real-time Shipping Stream (Live Pipeline Feed): A simulated high-velocity feed from logistics partners (e.g., FedEx/UPS API) providing JSON-based shipping status and location updates.
- Static Reference Data (Enrichment Data): Low-volume, read-only data, including the provided GeoJSON files (routes.geojson, etc.) and dimensional tables (Customer, Product) used to enrich the live data stream.
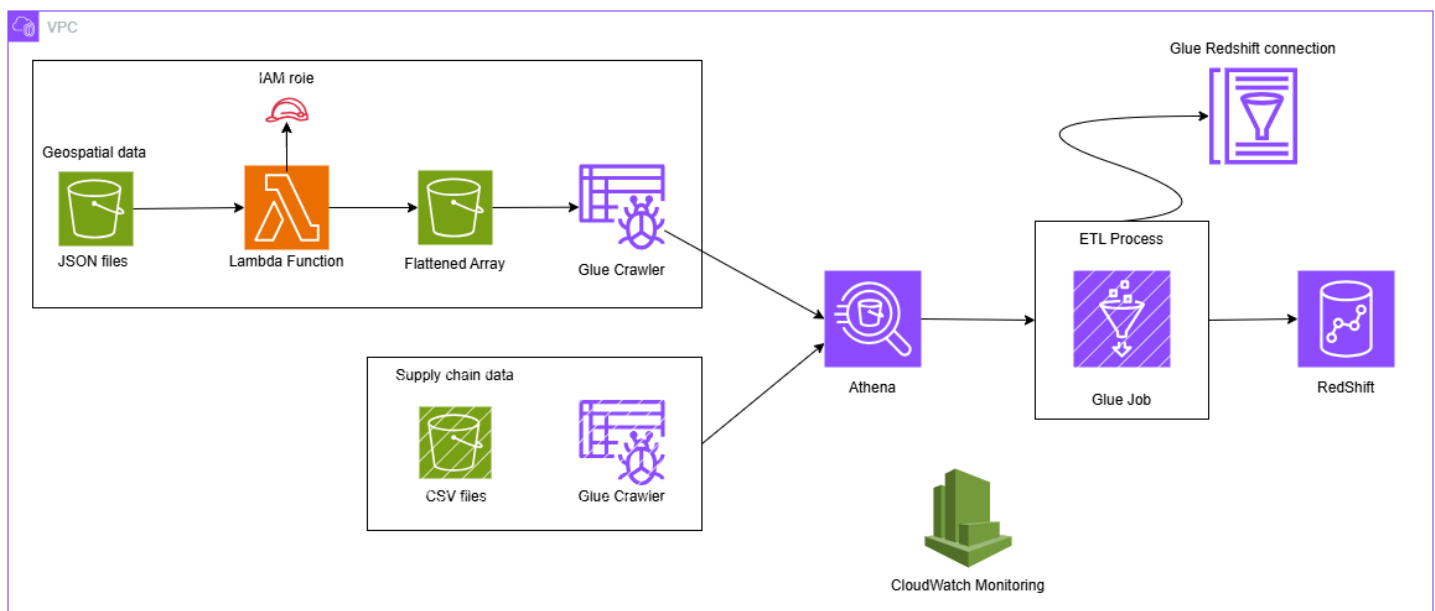
# Data Source:

https://www.kaggle.com/datasets/aaumgupta/refined-dataco-supply-chain-geospatial-dataset

# Data Pipeline

The cloud data pipeline is optimized for both structured and semi-structured data transformation. We are connecting a complex JSON source (GeoJSON) directly into the dimensional modeling process and another raw data file too. We have ingested data, transformed it using multiple AWS Services and combined it with historical and geospatial reference data to eventually enable real-time operational dashboard.

- **Ingestion Layer:** Loaded CSV (raw_data, metadata) and GeoJSON (routes.geojson) files into the Source S3 Bucket.
- **Transformation (Programming):** Used AWS Lambda Function (Python) to process the nested GeoJSON (routes.geojson), flatten the coordinates into WKT strings, and save the result as a usable CSV file back to S3.
- **Data Cataloging:** Used AWS Glue Crawlers to process the structured CSV files and the processed WKT files, loading the schemas into the Glue Data Catalog (logistream_db).
- **Data Verification:** The Data Catalog Tables were queried in Amazon Athena to validate schema integrity and data readability across all sources.
- **Pipeline Execution:** Created ETL jobs on AWS Glue ETL using a PySpark script editor, which extracted the data from the Catalog, performed dimensional modeling (Snowflake joins), and loaded the final data into the Redshift Serverless cluster.
- **Tool Configuration:** IAM Roles were set up for Lambda and Glue with necessary permissions (e.g., S3 read/write, Redshift access).
- **Documentation/Monitoring:** CloudWatch Logs were connected to monitor and troubleshoot the execution of the AWS Glue ETL Jobs.

# Architecture Diagram

## A. **Data Ingestion & Cataloging (Source Layer)**

1. **S3 Data Loading & Organization**

   All source files were loaded into a central S3 bucket and organized into 2 distinct folders: raw_data/ (main CSV), metadata/ (description CSV) in dataco-supply-chain-data and geojson/ (raw GeoJSON) in dataco-geospatial-data. Please refer to buckets 2 & 3 in the below screenshot.



**Dataco-geospatial-data :** Original source file is routes.geojson

**Dataco-supply-chain-data :** Original source files are inside folders - metadata & rawdata.

☐ metadata/DescriptionDataCoSupplyChainRefined.csv



☐ rawdata/DataCoSupplyChainDatsetRefined.csv

## 2. GeoJSON Pre-Processing (Lambda Transformation)

An AWS Lambda function (written in Python) was implemented to handle the complex, nested GeoJSON structure. This function reads the raw routes.geojson file, flattens the nested coordinate array into a Well-Known Text (WKT) string, and saves the output as a clean CSV file in a new S3 folder (processed_routes/).

### 3. Data Crawlers & Cataloging using Glue

AWS Glue Crawlers were executed to register schemas into the Glue Data Catalog (logistream). One crawler scanned the structured CSV folders (rawdata/ and metadata/). A separate crawler was executed specifically over the processed_routes/ folder. This step ensures the complex geospatial data, now simplified by Lambda, is correctly cataloged as a standard table (processed_routes) ready for joining.

### 4. Data Catalog Verification (Athena Querying)

The resulting tables (rawdata, metadata, processed_routes) were explicitly queried using Amazon Athena. Querying the data directly in Athena provided immediate schema validation and confirmed that all sources were correctly linked and readable before the heavy ETL phase began.

# B. Transformation & Storage (ETL Process)

**5.   Redshift Schema Creation**

The destination warehouse tables were created in the Amazon Redshift Serverless workgroup. The data model utilizes a Snowflake Schema structure to reduce redundancy and optimize analytical performance. The specialized dim_route_shapes table was created to store the WKT geometry for geospatial analysis in Tableau later.

### 6.  Master ETL Job Execution (PySpark Script)

An AWS Glue ETL Job was created using a PySpark script (Master ETL Job). This job performs the core transformation: Extracts all three sources (rawdata, metadata, processed_routes) from the Data Catalog.Transforms the data by generating dimensional hierarchies and surrogate keys. Joins the main fact dataset (rawdata) with all dimension tables, including a complex join based on Lat/Long coordinates to link to the dim_route_shapes table (GeoJSON data). Loads the final, linked dimension and fact tables into the Redshift Serverless cluster.

## Data Warehouse:

## 7. IAM Roles and Permissions

Dedicated IAM Roles were set up for the Lambda function, Glue Job, and Redshift cluster to grant necessary cross-service permissions.



## 8. Monitoring and Logging

CloudWatch Logs were connected and used throughout the process to monitor job status, execution time, and troubleshoot connectivity errors.

# Business Intelligence

### 1. Profitability Hierarchy

This Treemap visualizes the contribution of Total Profit across the product dimensional hierarchy.



**Actionable Insight:**

This visual proves the value of the Snowflake Schema design, as it requires joining the fact table to multiple dimension tables (dim_department and dim_category).

Product/Marketing Strategy: Executives can immediately identify the most profitable product segments (the largest squares) and use this insight to prioritize marketing spend and production capacity on high-value items.

### 2. Total Sales Trend Over Time:

A time-series Line Chart tracking [Total Sales] over time (Month/Year).

**Actionable Insight:**

It validates the integrity of the dim_date dimension and the accuracy of the large volume batch-loaded into the fact table.

Demand Forecasting: Managers can quickly spot seasonal peaks and troughs in demand. This informs purchasing decisions and staffing levels in warehouses, ensuring optimal inventory levels to meet future customer orders.

3. **Schedule Reliability by Shipping Mode**

A Bar Chart auditing the [Schedule Adherence %] across different shipping_mode values (carriers).

**Actionable Insight:**

This visual requires the pre-calculated Schedule Adherence measure, which compares two original raw columns.

Carrier Optimization: This directly shows which carriers or shipping methods are reliable (high adherence) and which are causing delays. Logistics managers can use this data to negotiate better rates or switch underperforming carriers.

4. **Late Order Volume by Country**

A Map visual showing the geographical distribution of [Total Orders at Late Risk] using bubble size by country.



**Actionable Insight:**

It utilizes the final calculated risk flag (late_risk = 1) and the geographic dimension, linking a key operational metric to a physical location for action.

Proactive Risk Mitigation: Large bubbles immediately signal regional logistics hotspots where intervention is required. This allows management to preemptively contact customers or deploy regional solutions, turning a retrospective problem into a proactive service opportunity.

# Conclusion

The LogiStream project successfully built a unified cloud data pipeline on AWS, transitioning DataCo from purely retrospective analytics to a proactive BI platform. By leveraging AWS Glue and PySpark to transform raw data into a clean, analytical Snowflake Schema in Amazon Redshift, the pipeline established a single source of truth. The resulting Power BI dashboard operationalizes the data, delivering actionable insights such as proactive late-delivery alerts and data-driven guidance for route optimization and inventory placement, directly fulfilling the original business problem definition.