**Caremetrix**

# Data Pipeline Architecture Documentation

**Project Name: Caremetrix**

**Team: Synthea001**

## Introduction

This document provides a comprehensive overview of a robust, end-to-end data pipeline designed and implemented using Amazon Web Services (AWS). The pipeline is engineered to ingest, process, transform, catalogue, and analyse synthetic healthcare data generated by the **Synthea™** simulation platform.

## About Synthea™

Synthea™ is an open-source patient population simulator that produces realistic synthetic health records, including demographics, medical conditions, medications, procedures, and healthcare encounters. Unlike real patient data, Synthea's synthetic data maintains the statistical properties and distributions of real populations while ensuring privacy and compliance by avoiding the use of any personal identifiable information (PII). This makes it ideal for testing, research, and developing healthcare analytics platforms without privacy concerns.

## Pipeline Purpose and Goals

The primary purpose of this data pipeline is to transform raw synthetic healthcare data into clean, structured, and query-optimized datasets ready for business intelligence and analytics applications. The pipeline achieves this by leveraging AWS's serverless and managed services to minimize operational overhead while maximizing scalability, reliability, and performance.

## Key Architectural Components

- **Data Ingestion and Storage**: Raw synthetic data from Synthea is ingested and stored in Amazon S3, serving as a cost-effective, durable, and highly available data lake.
- **Data Pre-processing**: An AWS Glue job performs lightweight cleaning tasks such as removing extra commas, adding ingestion timestamps, and normalizing file formatting immediately after data ingestion, preparing the data for further processing.

- **Data Transformation and ETL**: AWS Glue jobs handle comprehensive cleansing, deduplication, data type standardization, and convert data to Parquet format for efficient storage and faster querying.
- **Metadata Cataloguing and Schema Enforcement using Glue**: AWS Glue Crawlers automatically infer schemas and update the Glue Data Catalog, while pre-defined schemas ensure consistency and prevent schema drift is not efficient. So, we decided to use Glue Job instead.
- **Data Querying and Analytics**: Amazon Athena enables serverless, interactive SQL queries on Parquet data.
- **Visualization and Reporting**: By using Power BI, we are facilitating the reports and stats.

## Architecture Components

| Component | Service Used | Purpose |
|---|---|---|
| Raw Storage | Amazon S3 | Store unprocessed Synthea data |
| Pre-processing | AWS Glue | Basic cleanup before deeper ETL |
| ETL & Transformation | AWS Glue | Data cleansing, deduplication, conversion |
| Metadata Catalog | Glue Crawler | Automatic schema inference |
| Schema Management | AWS Glue | Schema enforcement & validation |
| Query Layer | Amazon Athena & Glue | SQL-based exploration & data marts |
| BI/Visualization | Power BI | Dashboards, reports, insights |

## 3. Data Import Workflow

## Local Synthea Data Generation

- Clone the Synthea repository to your local environment.
- Compile and run the simulation with parameters specifying the number of patients and a random seed for reproducibility. For example, to generate data for 1,000 patients with a seed of 1749537673085, run:

```
. /run_synthea -p 1000 -s 1749537673085
```

- The simulation outputs CSV files into the `output/csv/` directory, including multiple tables such as patients, encounters, conditions, and more.

#NOTE: For more information on the table, you can refer About Tables.xlsx

## Upload to Amazon S3

- Upload the generated CSV files to the designated **raw data** folder in your Amazon S3 bucket using AWS CLI or SDKs.
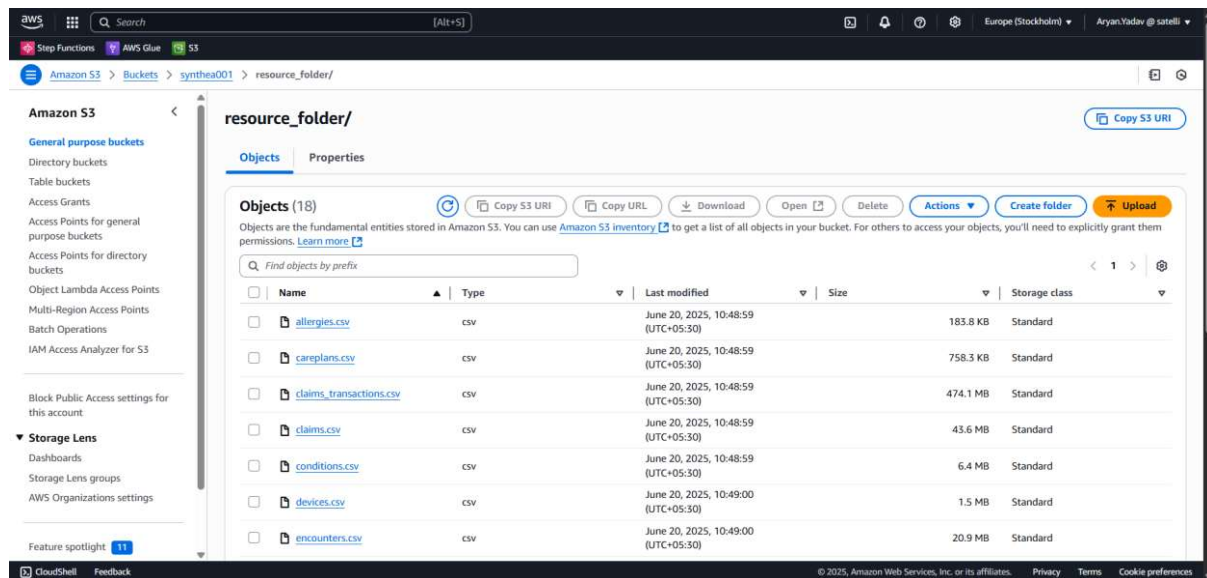- Recommended S3 folder structure:

```
s3://synthea001/resource_folder/
    └── patients.csv
    └── encounters.csv
    └── conditions.csv
    └── ...
```

## Pipeline Steps in Detail

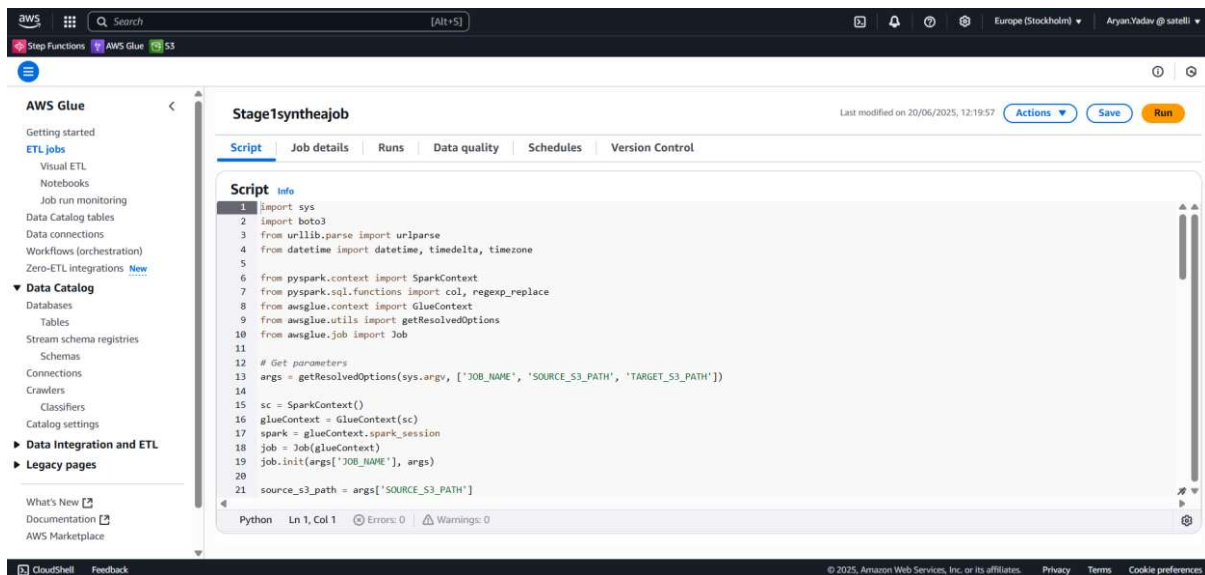## Step 1: Store Raw Data in Amazon S3

- **Service**: Amazon S3
- **Bucket Structure**:

```
s3://synthea001/resource_folder/
s3://synthea001/cleaning_stage_1/
s3://synthea001/cleaning_stage_2/
```
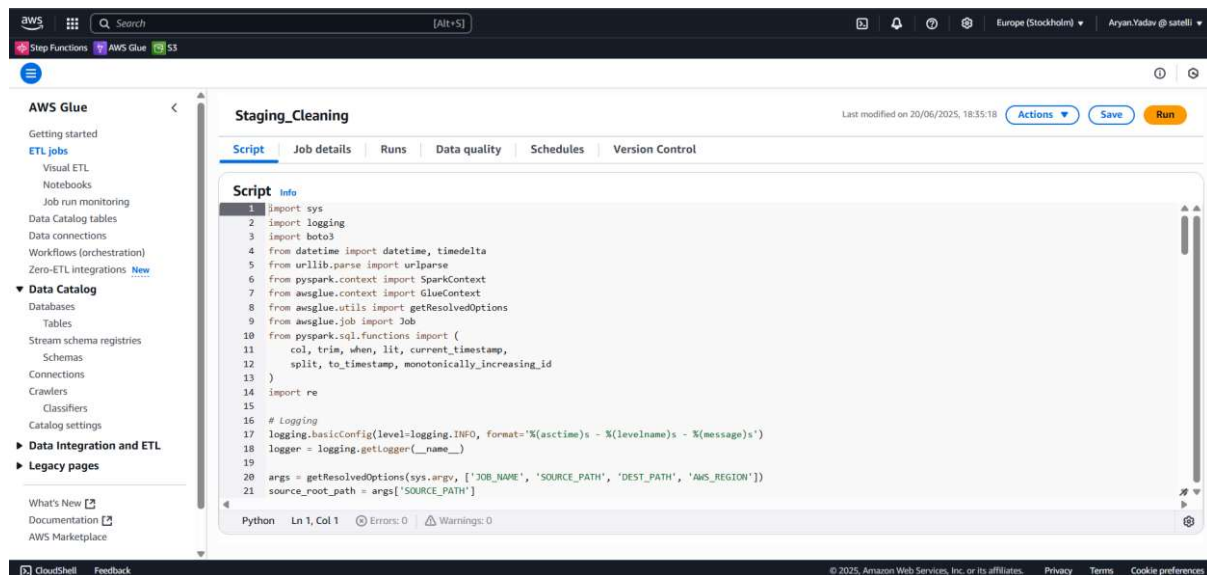
## Step 2: Pre-Processing via AWS Glue PySpark Job

- **Glue Job Role:**
  Initial lightweight preprocessing using PySpark to prepare raw CSV files for downstream transformation
- **Actions Performed:**
  - Remove trailing commas from each line to fix malformed rows.
  - Removing the redundant columns.
  - Write cleaned data to
    `s3://synthea001/cleaning_stage_1/#timestamp/`
- **Output Path:**
  Writes cleaned data to `s3://synthea001/cleaning_stage_1/#timestamp/`
- **Execution Time:**
  Typically, under 15 seconds per file, depending on file size and job parallelism
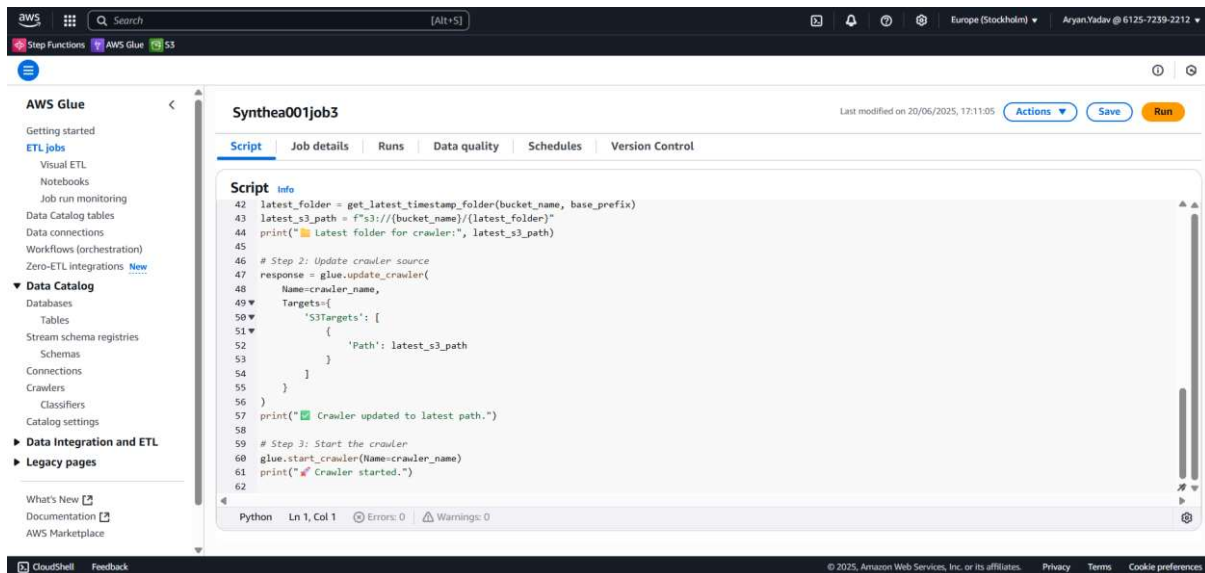
# Step 3: Data Cleansing & Parquet Conversion with AWS Glue

- **Service**: AWS Glue Jobs (PySpark)
- **Input**: Files in `s3://synthea001/cleaning_stage_1/#timestamp/`
- **Output**: Files in `s3://synthea001/cleaning_stage_2/#timestamp/` in **Parquet** format
- **Transformations**:
  - Remove duplicate rows
  - Convert date strings to timestamps
  - Apply consistent types across tables (e.g., `age` as integer)
  - Replace null values with None to ensure consistency
  - Standardize column names to lowercase
  - Handle Multivalued columns
  - Applied Surrogate keys and did the schema conversion handling
  - Maintained the SCD 2 in whole database
  - Convert data from CSV to optimized Parquet format for faster querying

- **Performance**:
  - Parallelized execution using AWS Glue Dynamic Frame and Spark jobs
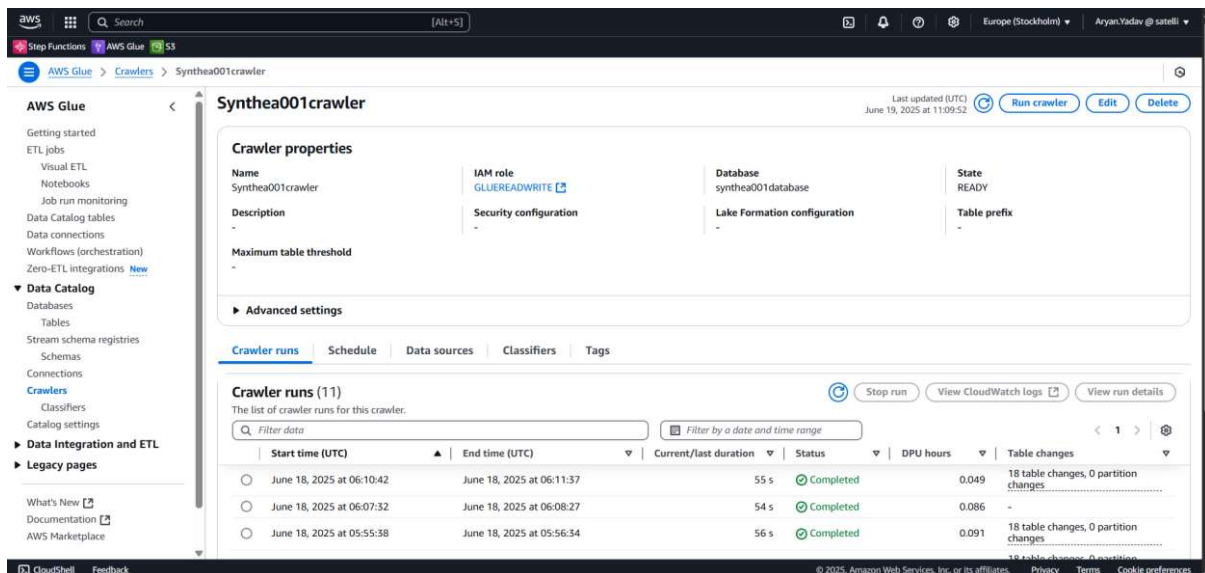  - Optimized for cost with job bookmarks (incremental processing)

## Step 4: Crawler Trigger and Catalog Update via AWS Glue Job

- **Action:**
  Trigger the AWS Glue Crawler to scan and update the catalog with new data from timestamped Parquet folders
- **Target Location:**
  `s3://synthea001/cleaning_stage_2/#timestamp`
- **Purpose:**
  - Updates the Glue Data Catalog with latest schema from cleaned and partitioned data
  - Enables Athena to discover new partitions and tables without manual intervention
  - Ensures metadata reflects structure of transformed production-ready datasets
- **Automation:**
  Executed as part of the orchestrated pipeline after second-level cleaning is complete

## Step 5: Data Cataloguing with AWS Glue Crawlers

- **Service**: AWS Glue Crawler
- **Source:** `s3://synthea001/Cleaning_Stage2/#timestamp`
- **Purpose**:
  - Detect new tables/partitions in Parquet files
  - Create or update tables in Glue Data Catalog
- **Output**: Structured metadata in Glue Catalog Database
- **Schedule**: Triggered after Glue job completion



## Step 6: Creation of dimensions and fact table using Glue script

- Initialize Spark & Glue Context and fetch job arguments (like job name and paths).

- Create Dimension Tables (dim_*): Reads source tables (patients, encounters, etc.), adds SCD2 columns (effective_date, end_date, is_current), and saves them as Parquet tables in S3.
- Build dim_providers by joining providers and organizations.
- Create Fact Table (fact_patient_e): Joins cleaned dimension tables and adds a surrogate key (fact_patient_id) along with SCD2 columns. Write Fact Table to S3 in Parquet format and commit the Glue job.



## Step 7: Pipeline Orchestration using AWS Step Functions

**Service:** AWS Step Functions

**Input:** Triggers for AWS Glue jobs and AWS Glue Crawler executions.

**Operations:**

- Define and execute a state machine that sequentially orchestrates:
  - AWS Glue job for initial cleaning (e.g., removing commas, adding timestamps)
  - AWS Glue jobs for cleansing and Parquet conversion
  - AWS Glue job to trigger the crawler and update the metadata catalog
- Incorporate error handling and retry logic for each step

**Benefits:**

- Ensures reliable, sequential execution of all pipeline stages
- Automates retries and error recovery to minimize manual intervention
- Provides detailed execution visibility and audit trail
- Scales automatically to handle multiple or large pipeline runs



## Step 8: Reporting and Visualization with Power BI

- **Action:**
  Connect Power BI to Amazon Athena using the JDBC driver to visualize query able healthcare datasets
- **Connection Method:**
  Athena JDBC connector (configured with AWS credentials and region-specific endpoint)
- **Purpose:**
  - Build interactive dashboards and reports based on curated data marts
  - Enable non-technical stakeholders to explore insights through a familiar BI tool

- o Support ad hoc analysis using direct SQL queries executed through Athena
- **Benefits:**
  - o Leverages Power BI's advanced visualization capabilities
  - o Serverless querying via Athena ensures scalability and cost-efficiency
  - o No data movement required—connects directly to data in Amazon S3 through Athena

## Final flow of the Pipeline:

Generate the desired number of patients data from Synthea database and import the tables to AWS S3.

**Step Function**

■ S3 (Raw Layer)
Raw synthetic healthcare data from Synthea is stored here.
s3://synthea001/resource_folder

■ Stage1syntheajob - Pre Processing of data
Basic cleaning including removing extra commas from the end and further adding timestamps.

■ S3 (Staging Layer)
Holds cleaned, structured data for further processing.
s3://synthea001/cleaning_stage_1/#timestamp/

■ Staging_Cleaning - Main Cleaning
cleans the data thoroughly, maintains SCD 2, and convert the cleaned CSV to Parquet format

■ S3 (Processed Layer)
Stores data in optimized Parquet format.
s3://synthea001/cleaning_stage_2/#timestamp/

+ ■ Synthea001job3 - AWS Glue Job with automatic crawler running in the end of the script

■ AWS Glue Catalog
Stores metadata.
Database Name: Synthea001database

■ patientdatamart_001job - Make the dimensions and fact tables

■ Athena
Serverless SQL queries on S3 data using Glue Catalog.