**Architecture Design:**
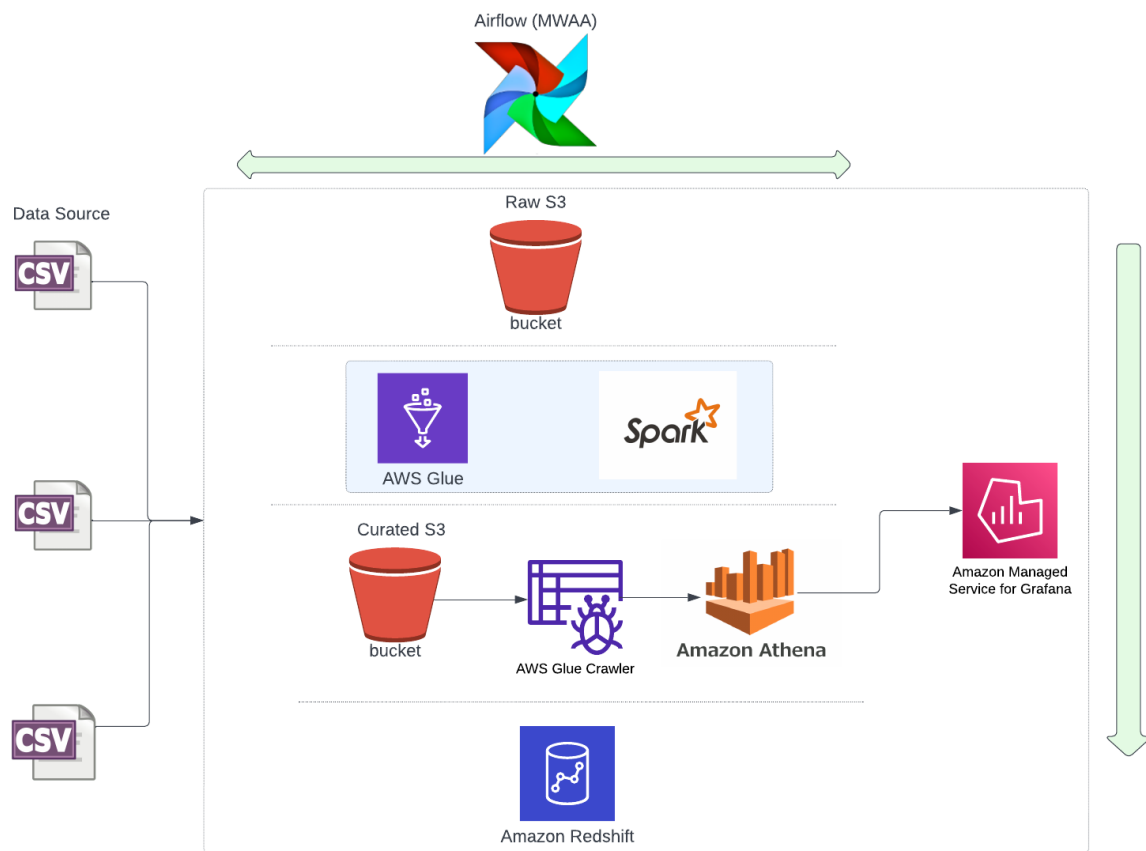


The whole architecture is based on AWS technologies primarily ETL Glue, Spark, Amazon Athena, AWS Glue Crawler, and Amazon managed Grafana. If the job needs to be executed periodically it can be achieved by either scheduling the ETL Glue jobs or using MWAA (Managed Airflow in AWS).

The process is the traditional ETL Pipeline setup where the raw data is read from the CSV files stored in either S3 raw buckets or fetched from the APIs or external services, it is then transformed using Apache Spark and loaded again to S3 curated buckets. From there the data is crawled using AWS Glue Crawler and exported to Glue Catalog tables. It is then exposed to Athena for querying using SQL. Finally, the Athena can be used a source for the Grafana Dashboards.

## Steps of this ETL Pipeline:
1. **Identify the Data Sources**
   For this POC (proof of concept) we are using multiple CSV files of weather data.
   The dataset consist of weather data from three different cities in India with attributes like temperature, pressure, humidity, and precipitation etc.

Source: https://www.kaggle.com/datasets/hiteshsoneji/historical-weather-data-for-indian-cities

2. **Create a ETL Glue, and Spark job to transform the data**
   For transforming the huge volume of data, there is a need of distributive architecture like spark which can perform the parallel computation on your data. We are using ETL Glue jobs in AWS to run the Spark script. Below are the steps to create and deploy the Glue jobs
   - Create a ETL Glue job from console (AWS Glue -> Jobs)



   Here you can select the Glue version, Spark Version, Python, and Worker Type

- Add the Spark Script (etl.py) under the script section



- Explore the job status from the Runs section

- Once the job is succeeded, create the Glue crawler to crawl the final data from curated S3 bucket. (AWS Glue -> Crawlers)



- Create a Database (weather) and Table (aggreagated_metricsoutput) in Glue Catalog.
- Run the Glue Crawler job to export data to Glue Tables (AWS

- Expose the Glue Tables to AWS Athena and Query using SQL



3. Use Athena as a Data Source for Grafana/Amazon Quicksight (dashboard attached as pdf separately)

*Note: This can also be achieved programmatically using Airflow, but it is out of scope for this project as we are dealing with one time static CSV files.*