

GHG and Methane Emissions Monitoring and Analysis Using AWS Glue and AWS Quicksight

Big Data Analytics Project Report

CIS 602: Special Topics in CIS, Summer 2024

Group 17

Members:

(3)Sindhuja Baikadi(02128756)

(38)Vaishnavi Paineni(02125170)

Project Title: GHG and Methane Emissions Monitoring and Analysis Using AWS Glue and AWS Quicksight

Abstract:

The "Dynamic GHG and Methane Monitoring and Analysis Using AWS Glue and AWS Quicksight" project aims to establish a comprehensive and scalable system for tracking real-time GHG and Methane emissions. By leveraging Amazon S3 for data ingestion, AWS Glue for data processing, and AWS Quicksight for Analysis, the project will enable continuous monitoring and in-depth analysis of GHG and Methane levels. This system is designed to provide immediate insights and critical information to support environmental management and policy-making efforts.

Data Set:

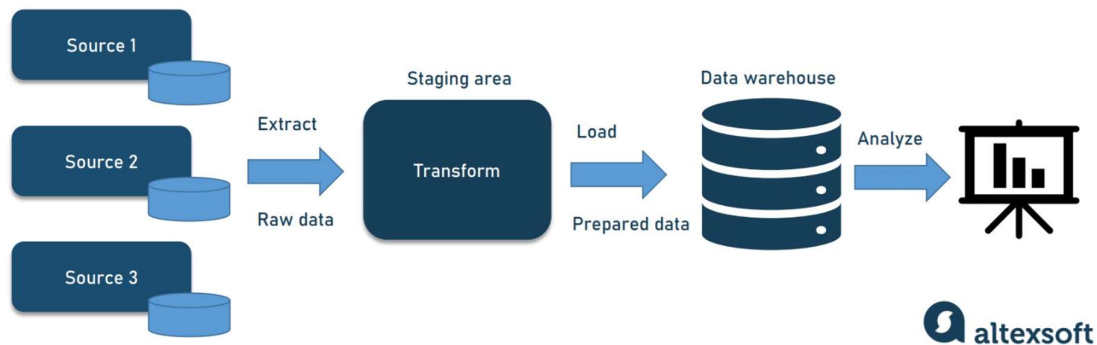
The dataset contains data from all the countries over the world. It has the following columns: Entity(Country), Code, Year, GHG, and Methane. The data is recorded over some time from 1850 to 2022.

- <https://ourworldindata.org/co2-and-greenhouse-gas-emissions>

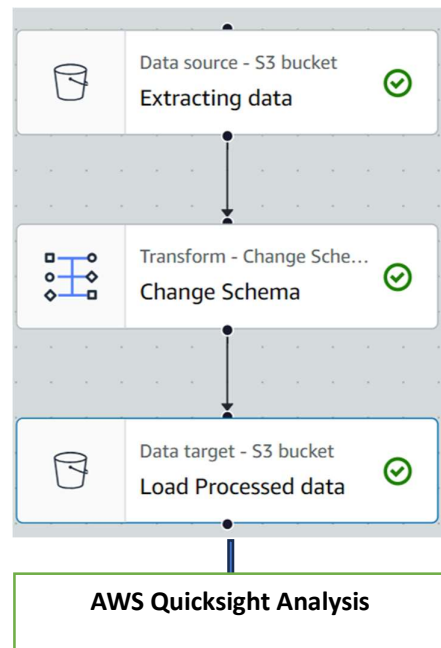
Objectives:

- 1. Data Acquisition:** Implement Amazon S3 Data Stream to consistently collect GHG and Methane emissions data from multiple sources such as industrial sensors, satellite data, and public APIs.
- 2. Data Preprocessing and Enrichment:** Use AWS Glue to clean, normalize, and enhance incoming GHG and Methane emissions data with additional contextual information like geographic and temporal details.
- 3. Instantaneous Data Insights:** Design real-time analytics to measure critical GHG and Methane metrics (e.g., concentration levels, emission trends) for specific regions, providing rapid insights into changes in GHG and Methane emissions.
- 4. Querying the data through Athena:** The data stored in the s3 bucket i.e., raw data and processed data is queried using Athena.
- 5. Structured Data Storage:** Store processed GHG and Methane emissions data in Amazon S3, systematically organized by geographic location, and indexed using AWS Glue Data Catalog for easy retrieval and analysis.
- 6. Scheduled Data Analysis:** Set up AWS Glue jobs to conduct regular batch processing, including historical GHG and Methane data analysis to detect long-term trends and anomalies in emissions.
- 7. Interactive Data Representation:** Utilize Amazon QuickSight or similar data visualization tools to create dynamic dashboards and reports that display real-time and historical GHG and Methane emissions data in an easily interpretable format.
- 8. Performance Monitoring:** Amazon CloudWatch continuously monitors the system's performance and health, ensuring reliability and efficiency.

AWS Data Pipeline:



Our Data Pipeline:

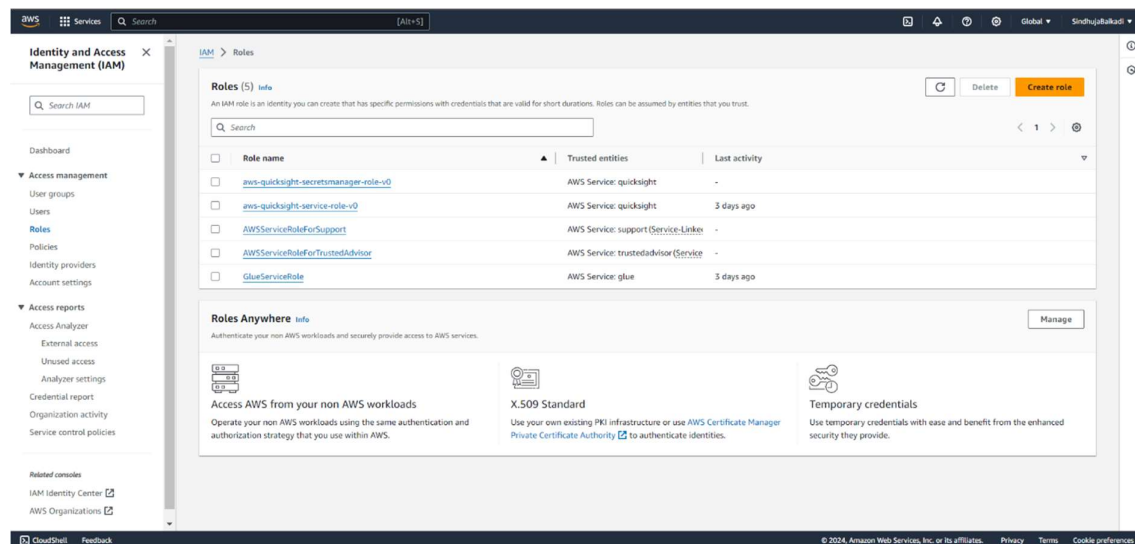


- 1. Data Ingestion:** Ingest data related to GHG and Methane into an AWS S3 bucket.
- 2. Data Transformation:** Use AWS Glue to preprocess and transform the incoming data on GHG and Methane. This step includes data cleaning, renaming columns, converting data types, and converting the data into Parquet format. Additionally, catalog the data using AWS Glue Data Catalog for streamlined access and discovery.
- 3. Data Storage:** Store the transformed data in Amazon S3.
- 4. Batch Data Processing:** Configure AWS Glue jobs to execute periodic batch processing tasks on the GHG and Methane data.
- 5. Data Visualization:** Employ a visualization tool like Amazon QuickSight to create interactive dashboards for analyzing and presenting the data.

Step-1: Setting IAM Role

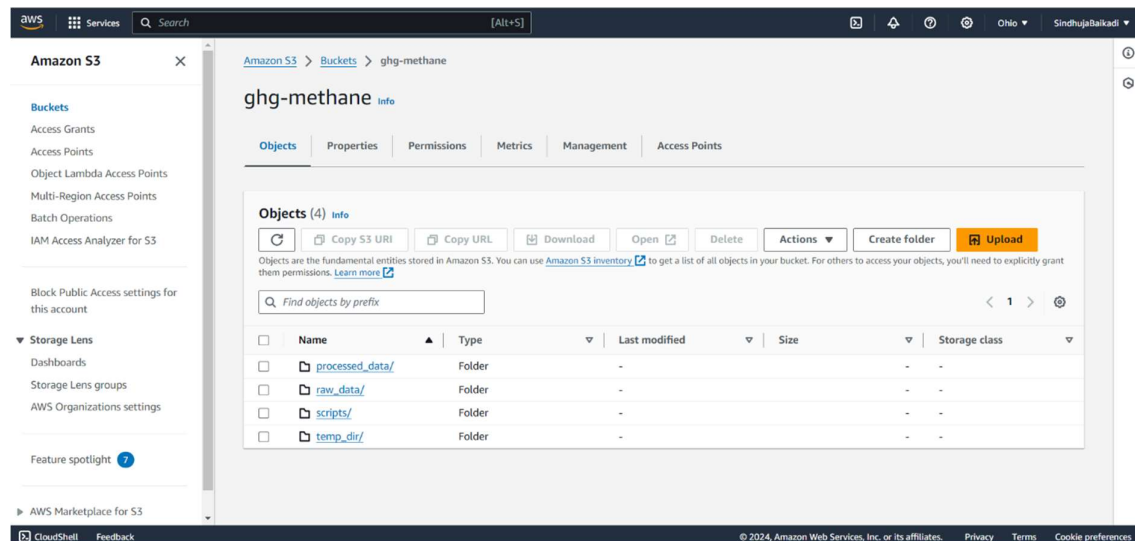
An IAM role in AWS is a set of permissions that allows users or services to perform actions on AWS resources. Unlike IAM users, roles are not tied to specific individuals and provide temporary access through assumed credentials. They are useful for granting access to resources without needing to share long-term credentials.

Created an IAM role with the name **“GlueServicerole”** with administrative access so that it can be used in all applications of AWS like AWS GLUE.



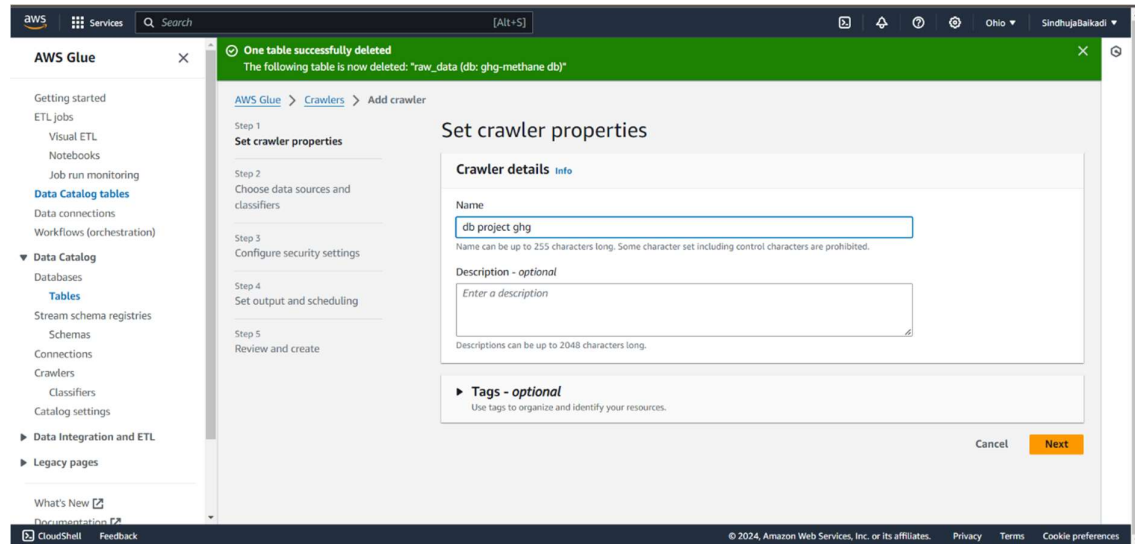
Step-2: Creating S3 buckets

An S3 bucket named **"ghg-methane"** was set up, containing subfolders for different purposes. The **`raw_data`** subfolder is used for uploading the initial raw data, and **`processed_data`** stores data after it has been transformed. A **`scripts`** folder holds the Py-Spark scripts necessary for these transformations, and a **`temp_dir`** folder is designated for temporary storage by AWS Glue during the querying process.

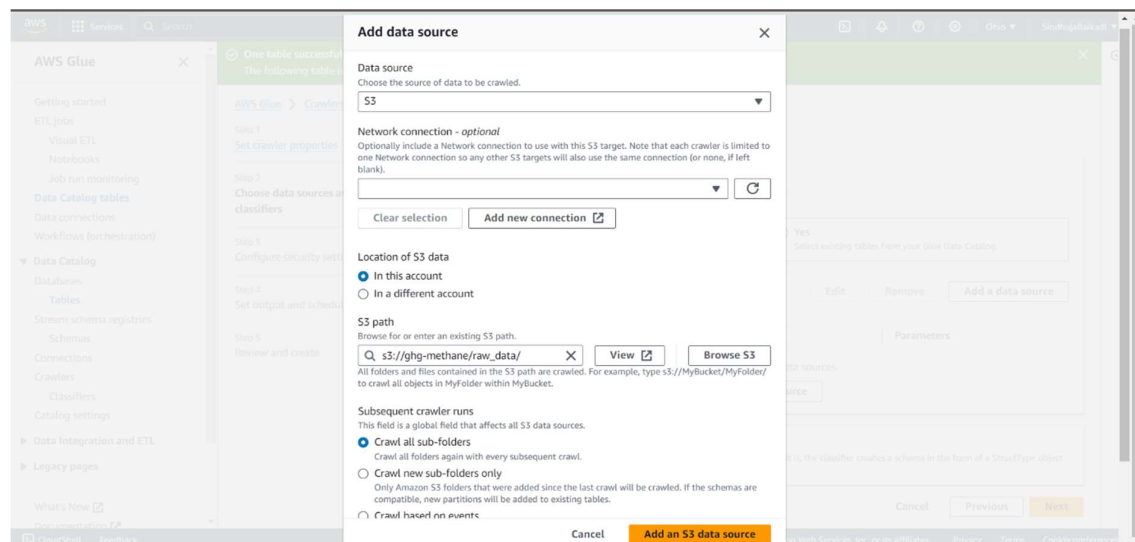


Step-3: AWS GLUE – ETL Setup and Process

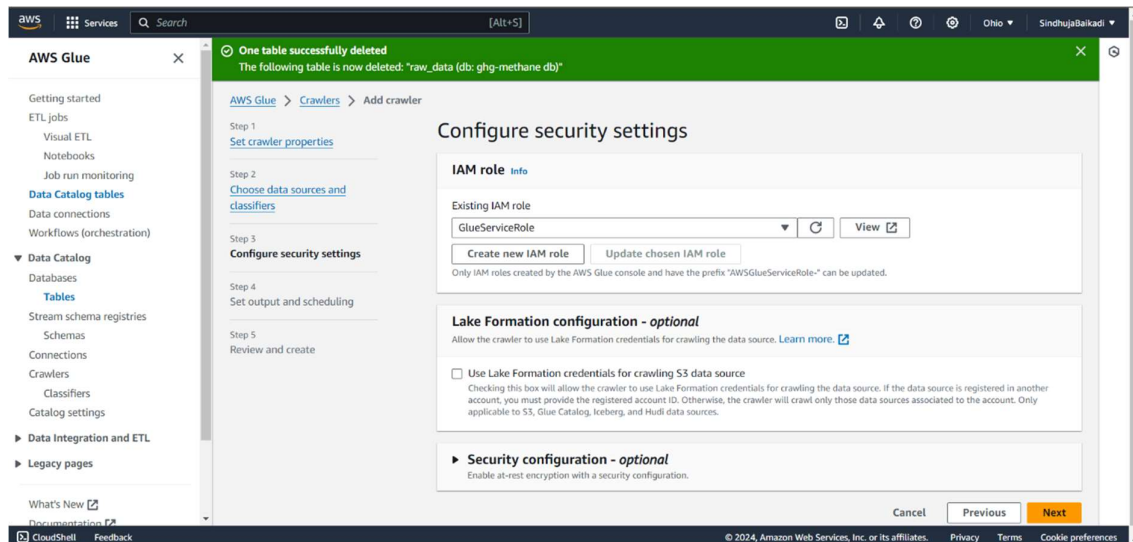
Open the AWS Management Console, search for "AWS Glue," and select it to access the AWS Glue console. In the navigation pane, choose Tables under Databases, then click "Add Tables using a crawler" and enter "db project ghg" as the name. Click "Next."



To add a data source, select S3 and choose "In this account" for the data location. Enter the S3 bucket path for raw_data, set the crawler to "Crawl all sub-folders," and add the S3 data source.

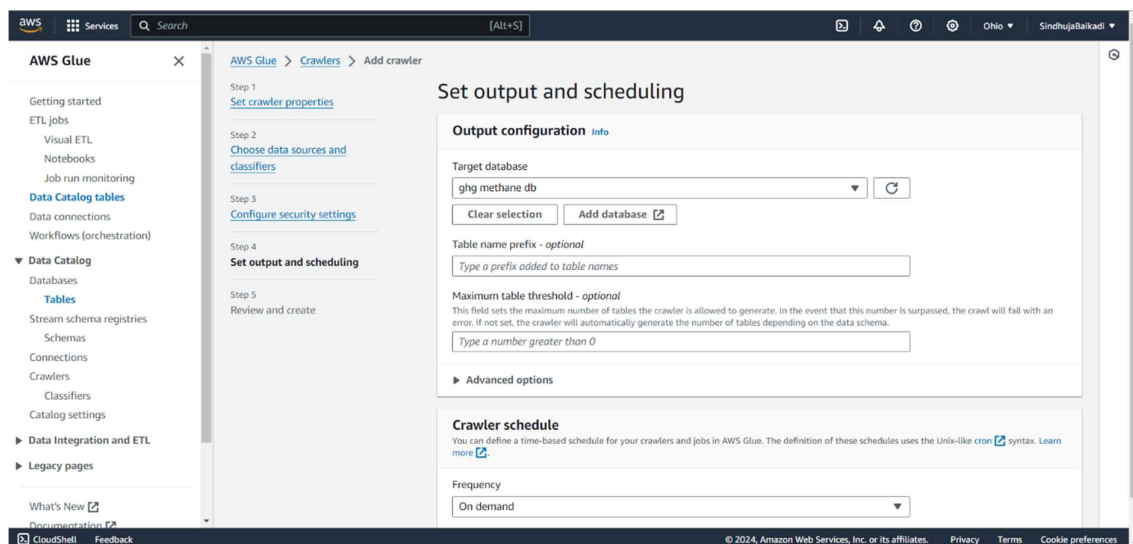


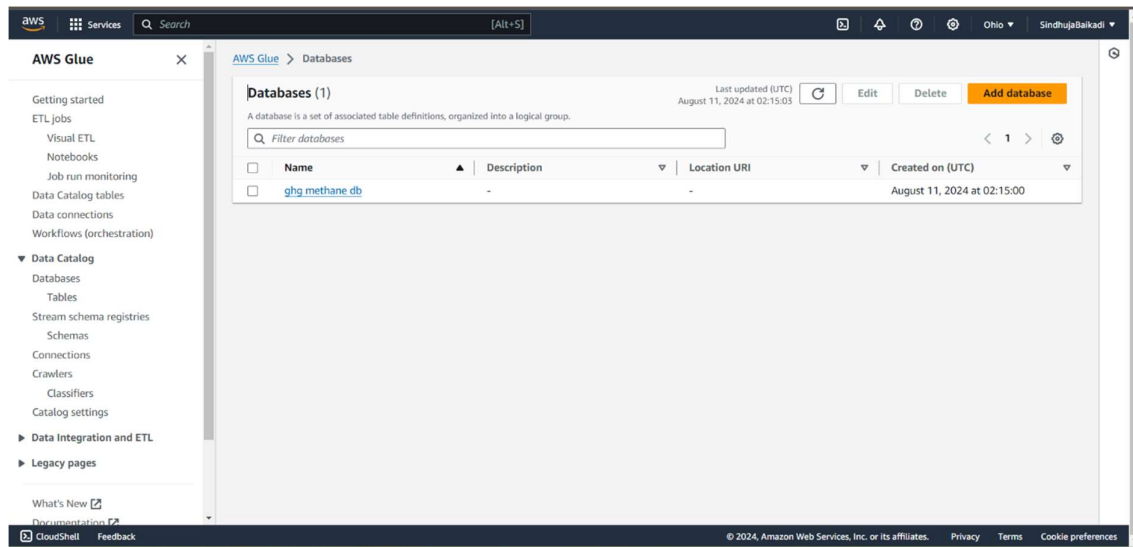
Click "Next," select "GlueServiceRole" for the IAM role, and proceed by clicking "Next" again.



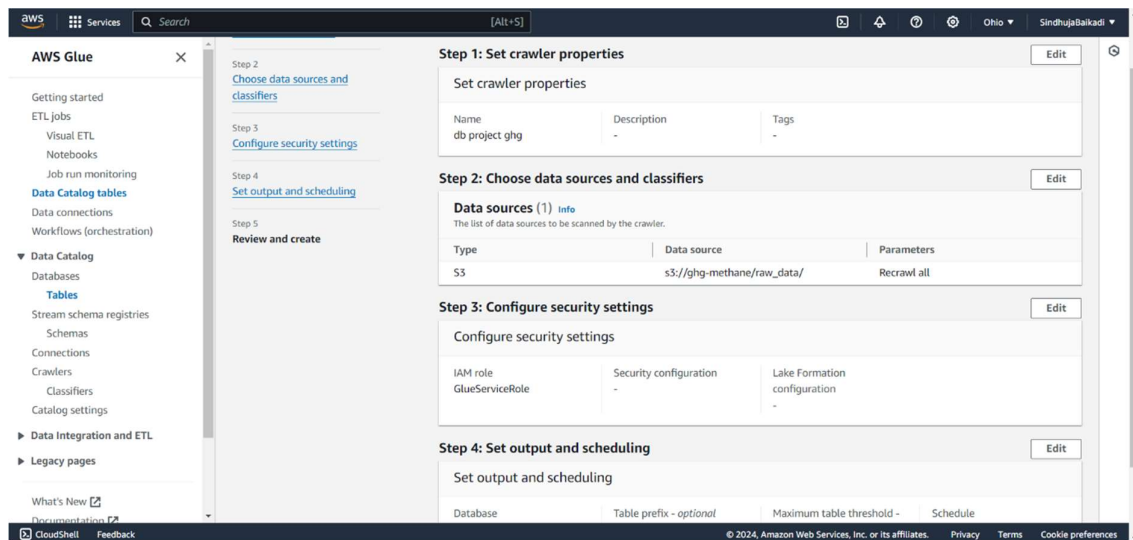
In the Output configuration section, select "Add database" to open a new tab, name the database "ghg methane db," and click "Create database."

Return to the Set output and scheduling page, and for the Target database, select "ghg methane db." In the Crawler schedule section, keep the Frequency set to "On demand" and click "Next."

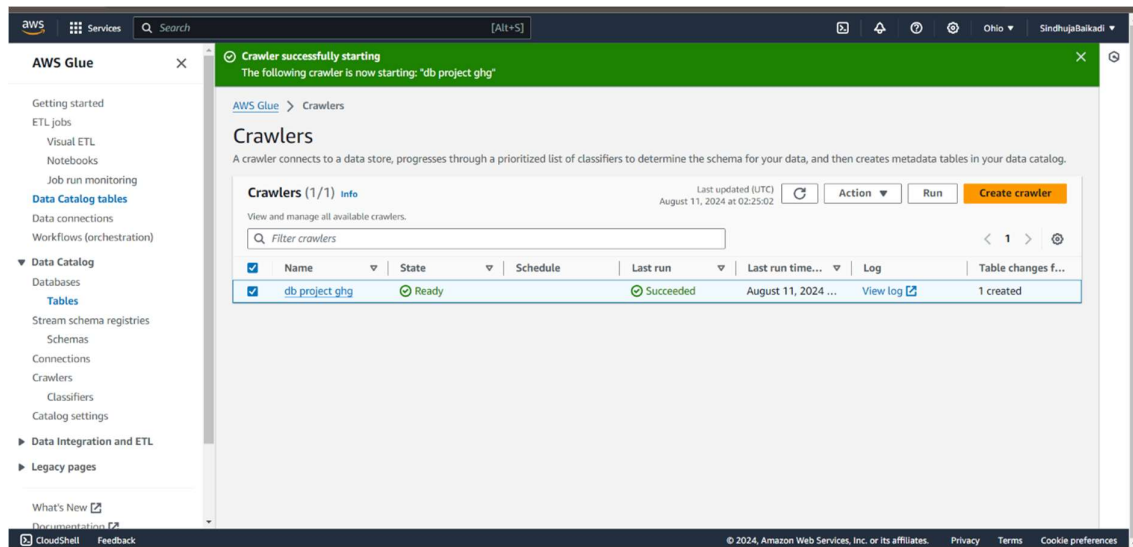




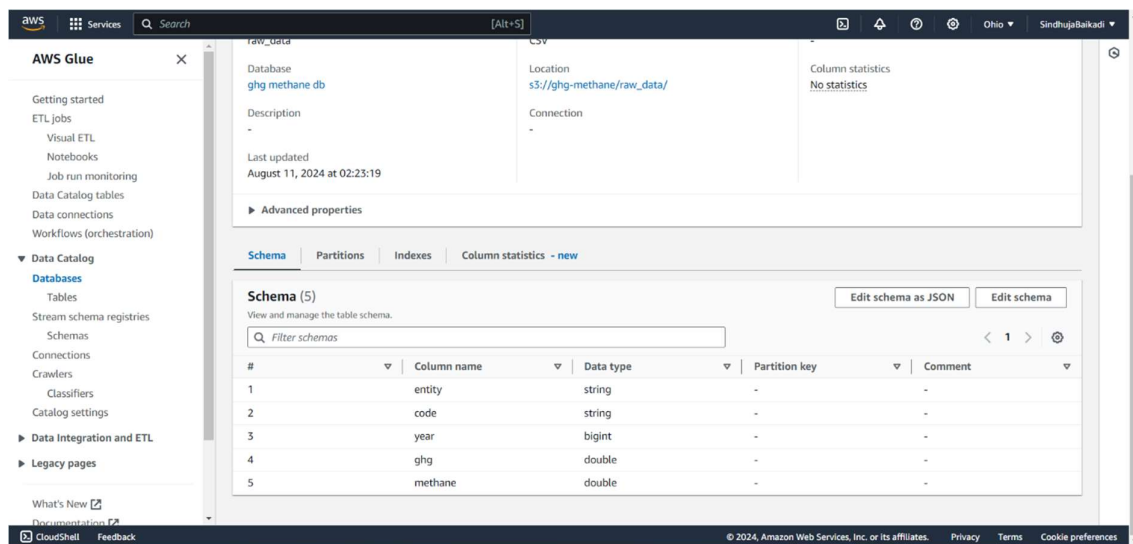
Review the configuration, then click "Create crawler."



Start the AWS Glue crawler you set up to execute the extract and load steps of your ETL process. Go to the Crawlers page, select the "db project ghg" crawler, and click Run. The status will update to Running as it creates the database and builds the metadata.



To check the AWS Glue metadata, go to the Databases section in the console, and select the "ghg methane db" database. This will show the schema with the columns identified by the crawler.



Under ETL Jobs, click on Visual ETL, **Extract** – The data was imported from the CSV file and uploaded to the designated location within the S3 bucket.

Untitled job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Data source properties - S3

Name: Extracting data

S3 source type: ☒ Data Catalog table

Database: ghg methane db

Table: raw_data

Data preview (200) Info READY

Filter sample dataset

entity	code	year	ghg
Afghanistan	AFG	1850	1.9554577
Afghanistan	AFG	1851	1.9649864
Afghanistan	AFG	1852	1.9725887
Afghanistan	AFG	1853	1.9790903

Transform – Various transformations were applied, such as checking for null values, renaming columns, and removing unnecessary columns that won't be needed in the subsequent stages of the ETL process. Dropped Code column and renamed entity column name as the country.

Untitled job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Transform

Name: Change Schema

Node parents: Extracting data

Change Schema (Apply mapping)

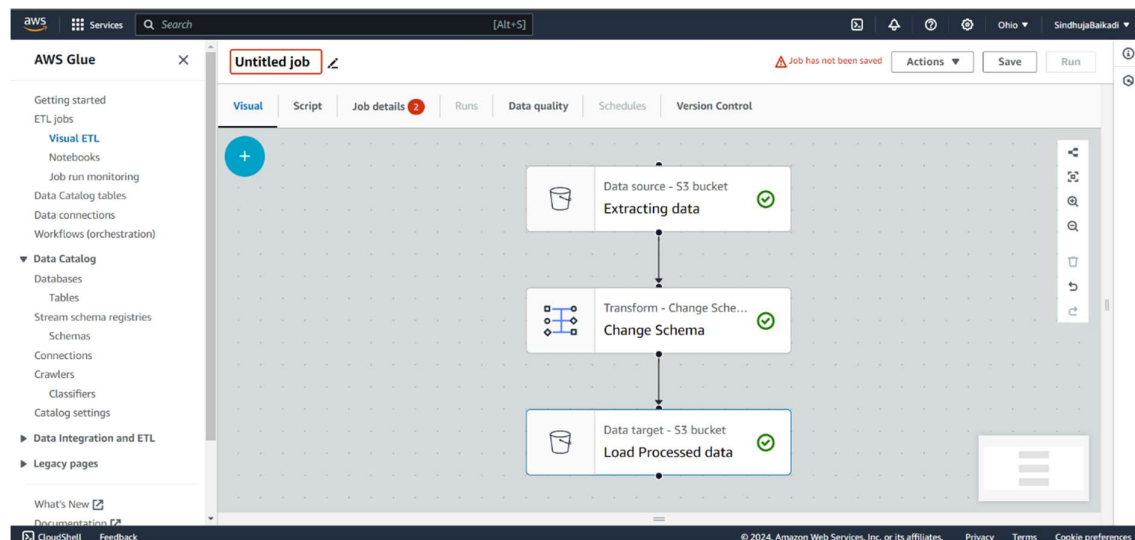
Source key	Target key	Data type	Drop
entity	entity	string	<input type="checkbox"/>
code			<input checked="" type="checkbox"/>
year	year	long	<input type="checkbox"/>
ghg	ghg	double	<input type="checkbox"/>
methane	methane	double	<input type="checkbox"/>

Data preview (200) Info READY

Filter sample dataset

entity	year	ghg	methane
Afghanistan	1850	1.9554577	0.9578825
Afghanistan	1851	1.9649864	0.9594418
Afghanistan	1852	1.9725887	0.96073055
Afghanistan	1853	1.9790903	0.96174896

Load – After completing the transformations, the processed data was uploaded to a specific S3 location for storing processed data. Additionally, a data catalog table was created to store this processed data. The code column was removed and the entity column renamed as country.



Output of ETL Jobs:
Jobs that ran successfully after ETL.

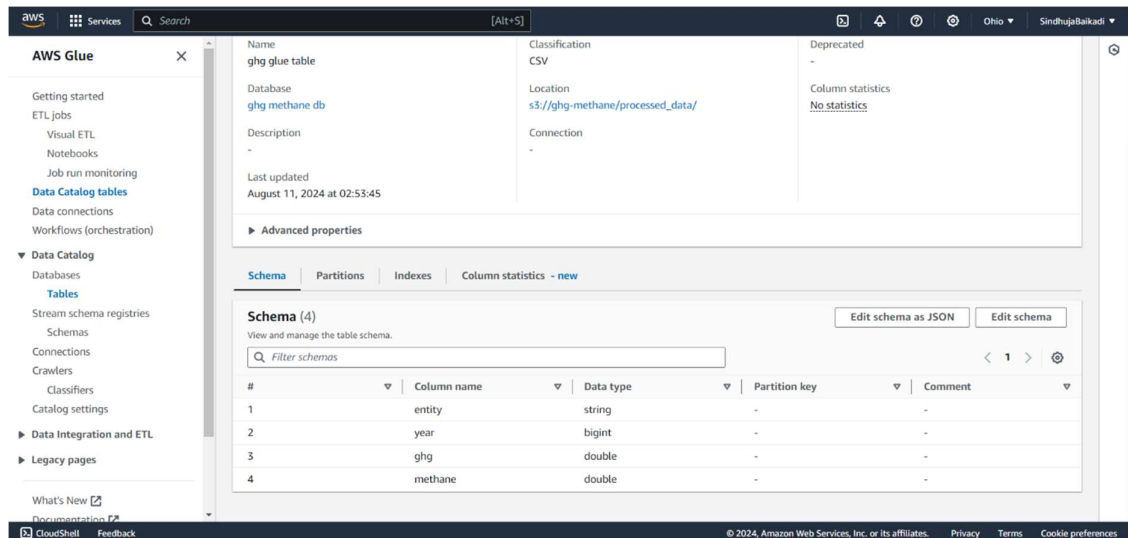
The screenshot displays the AWS Glue console interface for a job named 'Projectrole'. The job is in the 'Runs' tab, showing a table of job runs. The first run is 'Succeeded' with 0 retries, starting at 08/10/2024 22:52:59 and ending at 08/10/2024 22:53:56, with a duration of 47 seconds. Below the table, the 'Run details' section provides more information about the job run, including the job name, ID, status, and various configuration parameters.

Run status	Retries	Start time (Local)	End time (Local)	Duration	Capacity	Worker type	Glue version
Succeeded	0	08/10/2024 22:52:59	08/10/2024 22:53:56	47 s	10 DPU's	G.1X	4.0

Run details

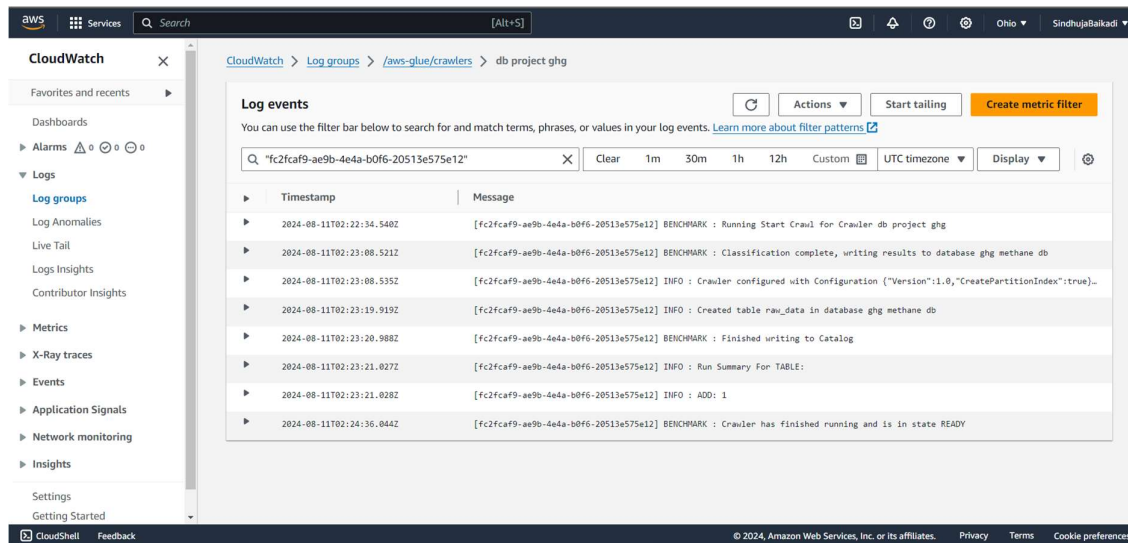
Job name	Start time (Local)	Glue version	Last modified on (Local)
Projectrole	08/10/2024 22:52:59	4.0	08/10/2024 22:53:56
Id	End time (Local)	Worker type	Log group name
jr_10e102322bec0684eaa08f4d964744e6b	08/10/2024 22:53:56	G.1X	/aws-glue/jobs
718e790aee9aa5e4eb1791dfbce1c			
Run status	Start-up time	Max capacity	Number of workers
Succeeded	10 seconds	10 DPU's	10
Retry attempt number	Execution time	Execution class	Timeout
Initial run	47 seconds	Standard	2880 minutes
Trigger name	Security configuration	Cloudwatch logs	Usage profile

Processed data in data schema:
We can see that the code column is removed.



Step-4: Cloud Watch logs

CloudWatch Logs is an AWS service that enables you to monitor, store, and access log files from your applications, systems, and AWS services. It allows you to collect, monitor, and analyze log data in real time, set alarms based on log metrics, and automatically trigger actions when specific conditions are met. CloudWatch Logs also provides long-term storage and retention of log data, making it easier to troubleshoot and audit activities in your AWS environment.



Step-5: Querying the Data Through Athena

Query the raw and processed data stored in the S3 bucket using Amazon Athena for ad-hoc analysis and reporting.

1. How many times does each country appear in the entity column? (no of entries for each country)

Query 7 | X | **Query 15** | X | **Query 8** | X | **Query 6** | X | **Query 9** | X | **Query 10** | X | **Query 12** | X | **Query 13** | X | **Query 14** | X | **Query 11** | X

```

1 --How many times each country appears in the entity column. (no of entries for each country)
2 SELECT s.'entity', COUNT(*) as EntityCount
3 FROM 'ghg methane db'.'raw_data' s
4 GROUP BY s.'entity'

```

Query results | Query stats

Completed Time in queue: 65 ms Run time: 490 ms Data scanned: 1.34 MB

Results (208)

#	entity	EntityCount
1	Africa	175
2	Asia	175
3	Austria	175
4	Bahrain	175

2. Showing countries in the database

Query 7 | X | **Query 15** | X | **Query 8** | X | **Query 6** | X | **Query 9** | X | **Query 10** | X | **Query 12** | X | **Query 13** | X | **Query 14** | X | **Query 11** | X

```

1 --Showing countries in the database
2 SELECT DISTINCT s.'entity' as Country
3 FROM 'ghg methane db'.'raw_data' s
4 SELECT Country
5 FROM distinctCountries
6 ORDER BY Country;

```

Query results | Query stats

Completed Time in queue: 56 ms Run time: 475 ms Data scanned: 1.34 MB

Results (208)

#	Country
1	Afghanistan
2	Africa
3	Albania
4	Algeria
5	Andorra
6	Angola
7	Antigua and Barbuda
8	Argentina

3. Finding the Max and minimum GHG Emissions by Entity.

Query 7 | X | **Query 15** | X | **Query 8** | X | **Query 6** | X | **Query 9** | X | **Query 10** | X | **Query 12** | X | **Query 13** | X | **Query 14** | X | **Query 11** | X

```

1 --Finding the Max and min GHG Emissions by Entity
2 SELECT s.'entity', MIN(s.'GHG') as MIN_GHG, MAX(s.'GHG') as MAX_GHG
3 FROM 'ghg methane db'.'raw_data' s
4 GROUP BY s.'entity'
5 ORDER BY s.'entity'

```

Query results | Query stats

Completed Time in queue: 57 ms Run time: 634 ms Data scanned: 1.34 MB

Results (208)

#	entity	MIN_GHG	MAX_GHG
1	Afghanistan	0.8014482	3.1628845
2	Africa	1.6962919	9.1568937
3	Albania	1.7155292	4.7964767
4	Algeria	0.56020725	6.9836493
5	Andorra	2.607438	15.876335
6	Angola	2.8515048	36.734413
7	Antigua and Barbuda	-0.6155252	21.197968
8	Argentina	8.119172	36.59449

4. Top 5 Entities by Total Methane Emissions

The screenshot shows the AWS Glue console interface. On the left, there's a sidebar with 'Data' and 'Tables and views' sections. The 'Data' section shows the 'ghg_methane_db' database. The 'Tables and views' section shows a table named 'ghg_methane'. The main area displays a SQL query that finds the top 5 entities with the highest total methane emissions. The query is as follows:

```
1 --Find 5 entities by total methane emissions
2 --This query finds the top 5 entities with the highest total methane emissions:
3 SELECT a."entity", SUM(a."methane") as TotalMethane
4 FROM "ghg_methane_db"."raw_data" a
5 GROUP BY a."entity"
6 ORDER BY TotalMethane DESC
7 LIMIT 5
8
```

The query results are displayed in a table with 5 rows and 2 columns: #, entity, and TotalMethane. The results are as follows:

#	entity	TotalMethane
1	Qatar	3264.0834502000004
2	Bahrain	2947.682657109999
3	Kuwait	2252.28799222
4	New Zealand	2175.853894599999
5	Australia	1479.972649000002

Results for the processed data:

The screenshot shows the AWS Glue console interface. The top bar indicates the query is 'Completed' with a time in queue of 52 ms, a run time of 739 ms, and data scanned of 539.00 KB. The 'Results (10)' section shows a table with 10 rows and 5 columns: #, entity, year, ghg, and methane. The results are as follows:

#	entity	year	ghg	methane
1	entity,year,ghg,methane			
2	Afghanistan,1850,1.9554577,0.9578825			
3	Afghanistan,1851,1.9649864,0.9594418			
4	Afghanistan,1852,1.9725887,0.96073055			
5	Afghanistan,1853,1.9790903,0.96174896			
6	Afghanistan,1854,1.9851059,0.96274185			
7	Afghanistan,1855,1.9907014,0.96370953			
8	Afghanistan,1856,1.9959223,0.9646519			
9	Afghanistan,1857,2.0007322,0.96556914			
10	Afghanistan,1858,2.0051491,0.966462			

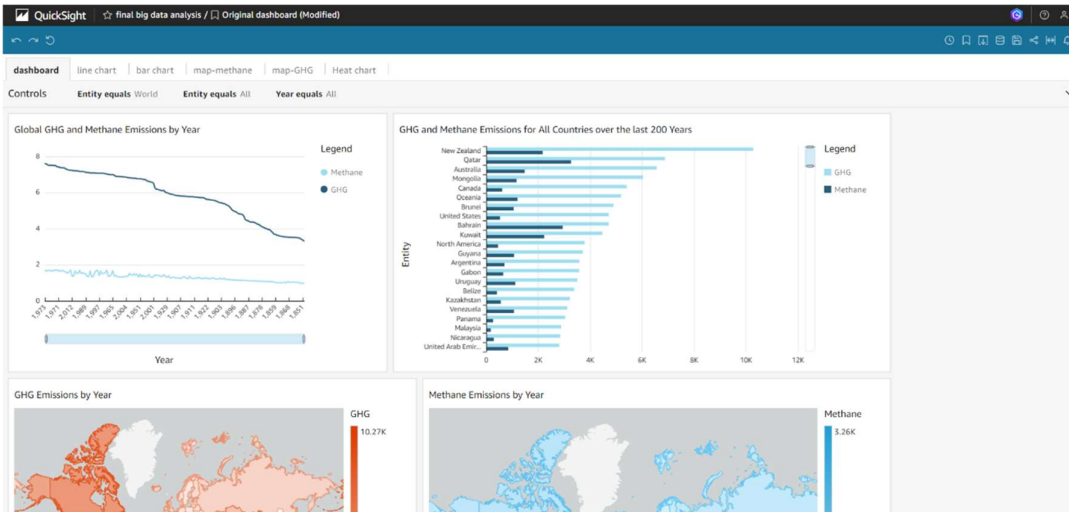
Step-6: AWS Quicksight Visualizations

- The dashboard provides a comprehensive view of both GHG and methane emissions globally, tracking their trends over time and across different countries.
- The comparative analysis highlights the countries that have historically been the largest contributors to these emissions. The heat maps offer a geographical perspective, making it easy to identify regions with the highest and lowest emissions.
- This dashboard can serve as a powerful tool for analyzing historical emission patterns and could be used to inform future environmental policies and strategies.
- **Global GHG and Methane Emissions by Year (Top Left)**

Line Chart: This chart shows the trend of global GHG and methane emissions over time. It indicates how both types of emissions have evolved, with the GHG emissions showing a more pronounced decline compared to methane emissions. In this, we have incorporated the entity filter.

- **GHG and Methane Emissions for All Countries over the Last 200 Years (Top Right)**

Bar Chart: This chart compares the cumulative GHG and methane emissions of various countries over the past 200 years. The countries are listed on the y-axis, and the x-axis represents the emission quantities. We have implemented year and entity filters in this.

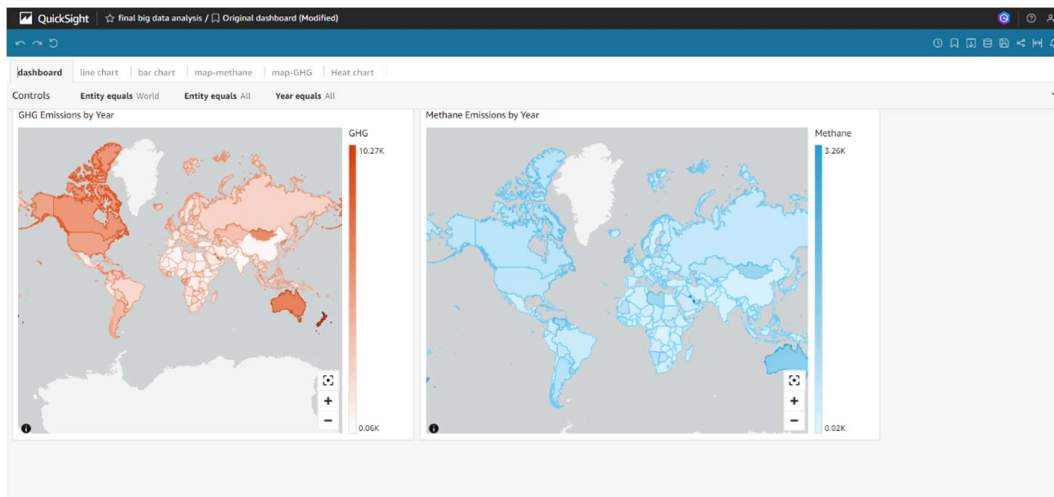


- **GHG Emissions by Year (Bottom Left)**

Heat Map: This map visually represents GHG emissions by country, with darker shades of red indicating higher levels of emissions. The color gradient on the right shows the range of emissions from 0.06K (lowest) to 10.27K (highest). The year filter is incorporated to change the year.

- **Methane Emissions by Year (Bottom Right)**

Heat Map: This map is similar to the GHG heat map but focuses specifically on methane emissions. The color gradient here ranges from 0.02K to 3.26K. Again, darker shades of blue represent higher emissions. The year filter is incorporated to change the year.



Conclusion:

The "GHG and Methane Emissions Monitoring and Analysis Using AWS Glue and AWS Quicksight" project aims to provide actionable insights and real-time GHG and Methane emissions data for various applications. Leveraging AWS services, this project will enhance decision-making and safety in scenarios influenced by GHG and Methane levels.

References:

- <https://ourworldindata.org/co2-and-greenhouse-gas-emissions>
- [Towards Data Science: Free and Reliable Weather Data Sources](<https://towardsdatascience.com/five-free-and-reliable-weather-data-sources-20b9ea6afac9>)
- [OpenWeatherMap API](<https://openweathermap.org/api>)
- [NOAA Real-time Data](<https://www.noaa.gov/education/resource-collections/data/real-time#:~:text=NOAA%20collects%20real%2Dtime%20data,%2C%20citizen%20scientists%2C%20and%20more>)
- [Weather.gov API](<https://www.weather.gov/documentation/services-web-api>)
- [OpenWeatherMap API](<https://openweathermap.org/api>)