# Amazon Web Services
# Data Engineering Immersion Day

Lab 2. ETL with AWS Glue

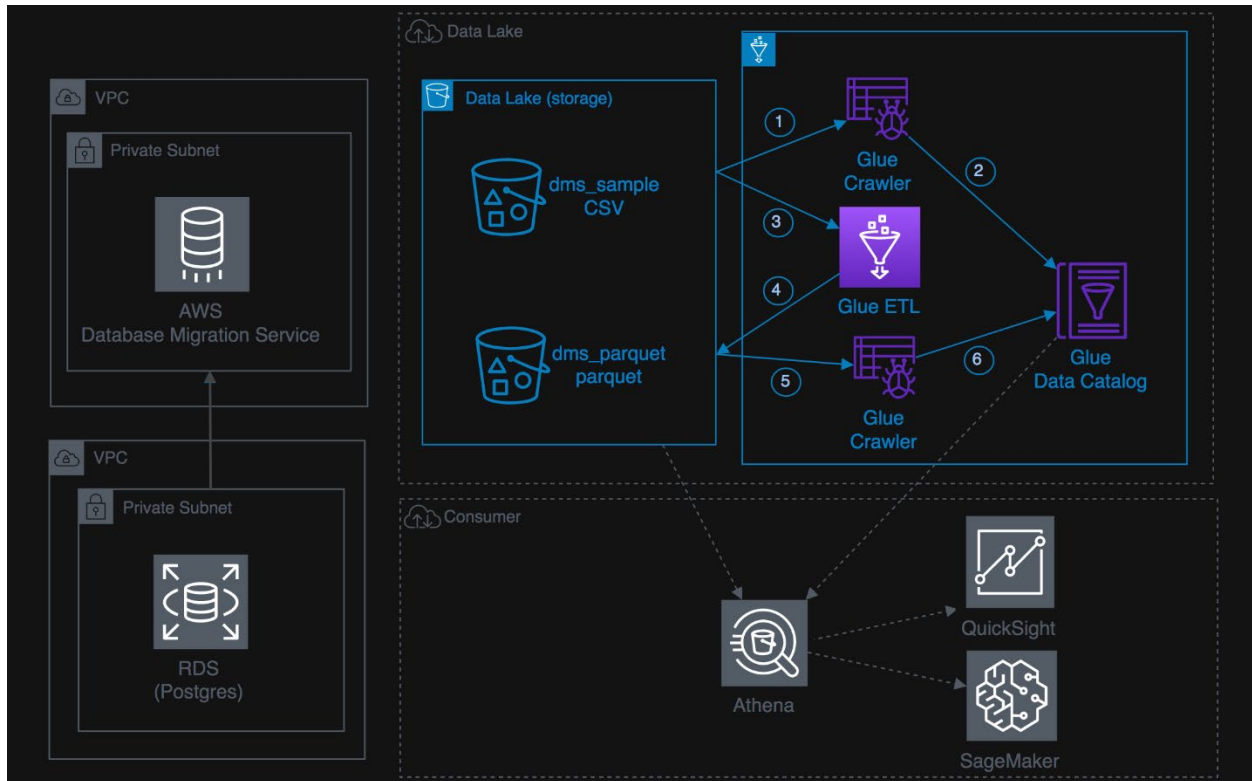*July 2021*

# Table of Contents

# Introduction

This lab will give you an understanding of the AWS Glue – a fully managed data catalog and ETL service



**Prerequisites**

1. Completed Lab 1. Hydrating the Data Lake with DMS

2. Or complete Lab1. Copy Source Data

**Tasks Completed in this Lab:**

In this lab you will be completing the following tasks. You can choose to complete only **Part-(A)** to move to next lab where tables can be queried using Amazon Athena and Visualize with Amazon QuickSight

1. PART-(A): Data Validation and ETL
2. PART- (B): Glue Job Bookmark Functionality(Optional)
3. PART- ( C ): Glue Workflows(Optional)

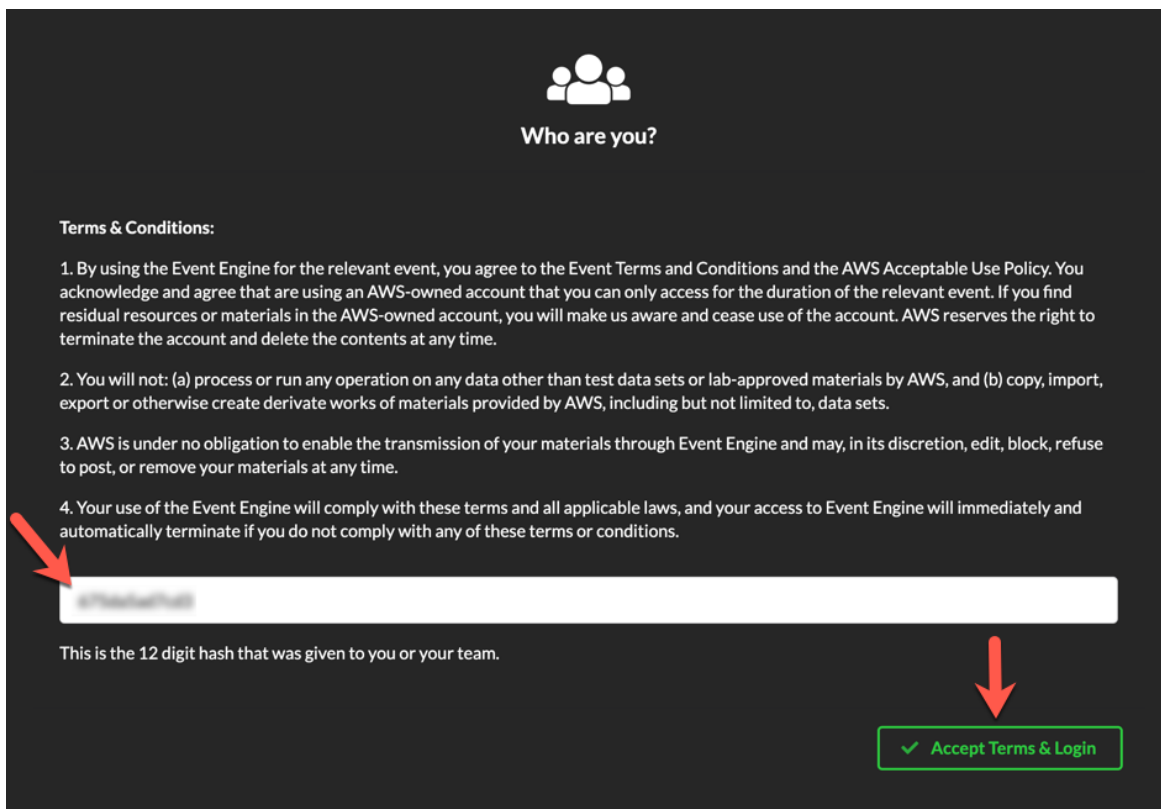The Lab is also available - https://aws-dataengineering-day.workshop.aws/

# Get Started Using the Lab Environment

Please skip this section if you are running the lab on your own AWS account.

Today, you are attending a formal event and you will have been sent your access details beforehand. If in the future you might want to perform these labs in your own AWS environment by yourself, you can follow instructions on GitHub - https://github.com/aws-samples/data-engineering-for-aws-immersion-day.

A 12-character access code (or 'hash') is the access code that grants you permission to use a dedicated AWS account for the purposes of this workshop.

1. Go to https://dashboard.eventengine.run/, enter the access code and click Proceed:
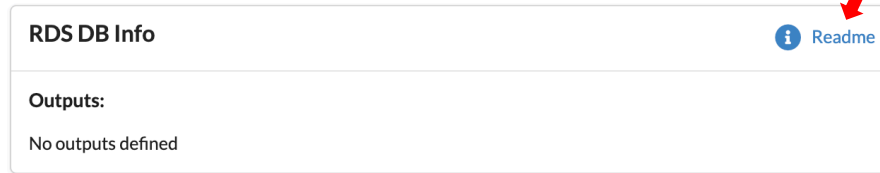


2. On the Team Dashboard web page you will see a set of parameters that you will need during the labs. Best to save them to a text file locally, alternatively you can always go to this page to review them. Replace the parameters with the corresponding values from here where indicated in subsequent labs:

Because you're at a formal event, some AWS resources have been pre-deployed for your convenience, for example:
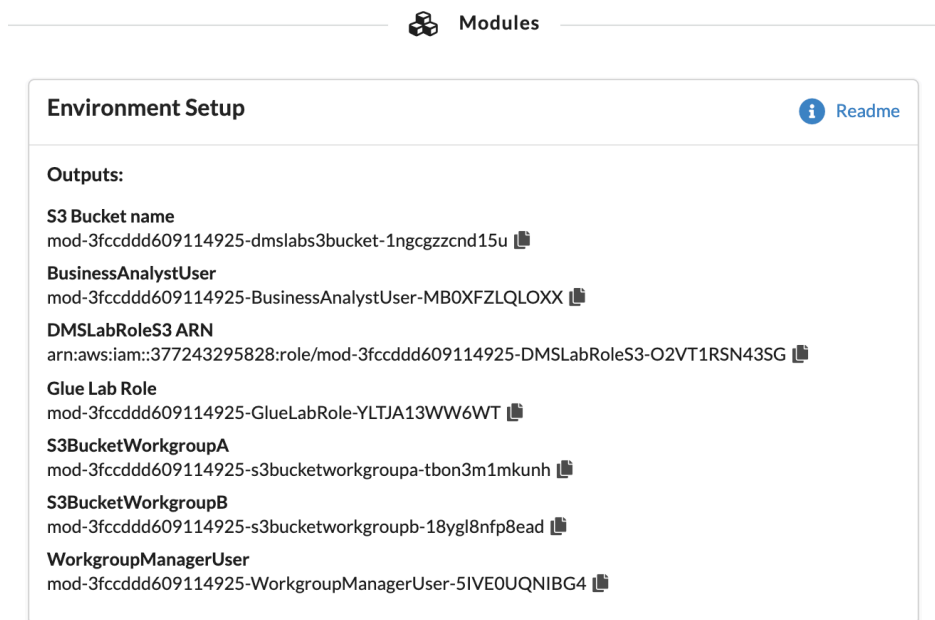
- The source database connection in RDS DB Info module

**RDS DB Info**                                    ⓘ Readme

**Outputs:**

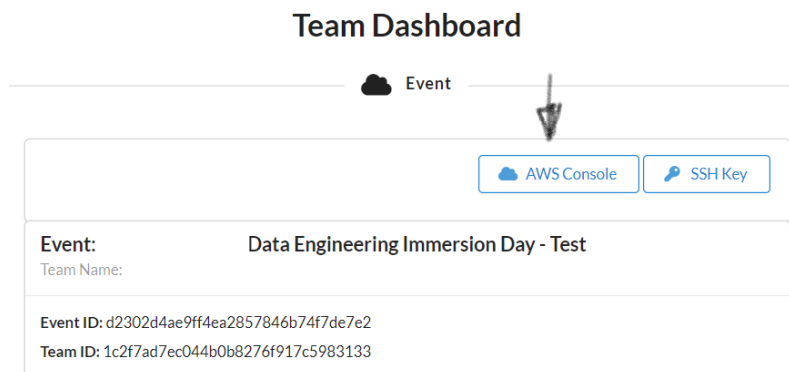No outputs defined

- S3 Bucket, IAM role for the Glue lab etc

⬡ Modules

**Environment Setup**                              ⓘ Readme

**Outputs:**

**S3 Bucket name**
mod-3fccddd609114925-dmslabs3bucket-1ngcgzzcnd15u 📋

**BusinessAnalystUser**
mod-3fccddd609114925-BusinessAnalystUser-MB0XFZLQLOXX 📋

**DMSLabRoleS3 ARN**
arn:aws:iam::377243295828:role/mod-3fccddd609114925-DMSLabRoleS3-O2VT1RSN43SG 📋

**Glue Lab Role**
mod-3fccddd609114925-GlueLabRole-YLTJA13WW6WT 📋

**S3BucketWorkgroupA**
mod-3fccddd609114925-s3bucketworkgroupa-tbon3m1mkunh 📋

**S3BucketWorkgroupB**
mod-3fccddd609114925-s3bucketworkgroupb-18ygl8nfp8ead 📋

**WorkgroupManagerUser**
mod-3fccddd609114925-WorkgroupManagerUser-5IVE0UQNIBG4 📋

3. On the Team Dashboard, please click AWS Console to log into the AWS Management Console:

**Team Dashboard**

☁ Event

☁ AWS Console      🔑 SSH Key

**Event:**         Data Engineering Immersion Day - Test
Team Name:

**Event ID:** d2302d4ae9ff4ea2857846b74f7de7e2
**Team ID:** 1c2f7ad7ec044b0b8276f917c5983133

4. Click Open Console. For the purposes of this workshop, you will not need to use command line and API access credentials:



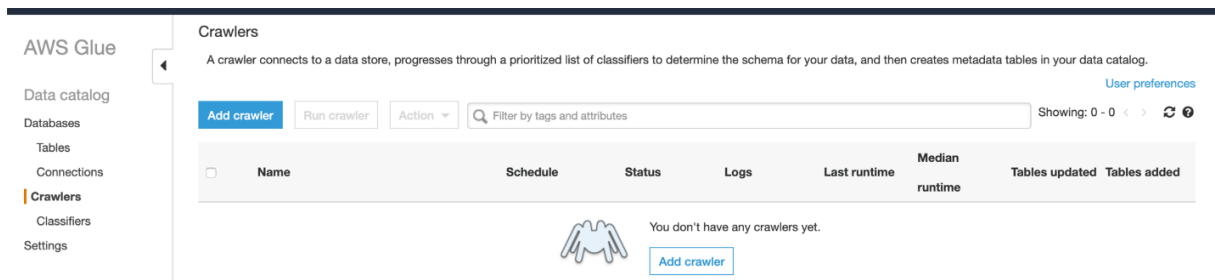Once you have completed these steps, you can continue with the rest of this lab.

# PART A: Data Validation and ETL

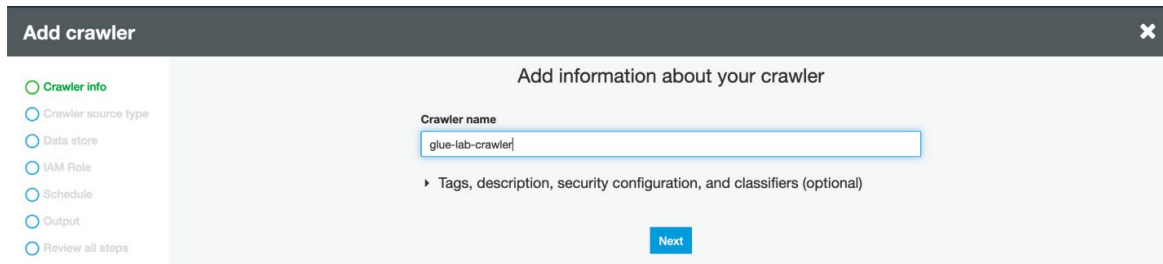## Create Glue Crawler for initial full load data

1. Navigate to the [AWS Glue service](#)



2. On the AWS Glue menu, select **Crawlers**.



3. Click **Add crawler**.
4. Enter **glue-lab-crawler** as the crawler name for initial data load.
5. Optionally, enter the description. This should also be descriptive and easily recognized and Click **Next**.



6. Choose **Data stores**, **Crawl all folders** and **Click Next**

7. On the **Add a data store** page, make the following selections:
   a. For Choose a data store, click the drop-down box and select **S3**.
   b. For Crawl data in, select **Specified path in my account**.
   c. For Include path, browse to the target folder stored CSV files, e.g., **s3://xxx-dmslabs3bucket-xxx/tickets**
8. Click **Next**.



9. On the **Add another data store page**, select **No**. and Click **Next**.



10. On the **Choose an IAM role** page, make the following selections:
    a. Select **Choose an existing IAM role**.
    b. **For IAM role**, select <stackname>-**GlueLabRole**-<RandomString> pre-created for you.
       For example "dmslab-student-GlueLabRole-ZOQDII7JTBUM"

11. Click **Next**.



12. On the Create a schedule for this crawler page, for Frequency, select **Run on demand** and Click **Next**.



13. On the Configure the crawler's output page, click **Add database** to create a new database for our Glue Catalogue.



14. Enter **ticketdata** as your database name and click **create**

15. For **Prefix added to tables (optional)**, leave the field empty.
16. For **Configuration options (optional)**, select **Add new columns only** and keep the remaining default configuration options and Click **Next**.



17. Review the summary page noting the Include path and Database output and Click **Finish**. The crawler is now ready to run.

18. Tick the crawler name, click **Run crawler** button.



Crawler will change status from starting to stopping, wait until crawler comes back to ready state (the process will take a few minutes), you can see that it has created **15 tables**.

19. In the AWS Glue navigation pane, click **Databases** > **Tables**. You can also click the **ticketdata** database to browse the tables.

## Data Validation Exercise

1. Within the Tables section of your **ticketdata** database, click the person table.

You may have noticed that some tables (such as person) have column headers such as col0,col1,col2,col3. In absence of headers or when the crawler cannot determine the header type, default column headers are specified.

This exercise uses the person table in an example of how to resolve this issue.

2. Click **Edit Schema** on the top right side.



3. In the Edit Schema section, double-click **col0** (column name) to open edit mode. Type "id" as the column name.

Repeat the preceding step to change the remaining column names to match those shown in the following figure: `full_name`, `last_name` and `first_name`

4. Click **Save**.

## Data ETL Exercise

**Pre-requisite:** To store processed data in parquet format, we need a new folder location for each table, eg. the full path for sport_team table look like this –

"s3://<s3_bucket_name>/tickets/**dms_parquet/sport_team**"

Glue will create the new folder automatically, based on your input of the full file path, such as the example above. Please refer to the user guide in terms of how to manually create a folder in S3 bucket.

1. In the left navigation pane, under ETL, click **AWS Glue Studio**.



2. Choose "**View Jobs**"

3. Leave the "Visual with a source and target" option selected, and press "**Create**"



4. Select the "Data source - S3 bucket" at the top of the graph.

5. In the panel on the right under "Data source properties - S3", choose the "**ticketdata**" database from the drop down.

6. For Table, select the **sport_team** table.



7. Select the "ApplyMapping" node. In the Transform panel on the right and change the data type of "id" column to double in the dropdown.

Lab 2. ETL with AWS Glue



8. Select the "Data target - S3 bucket" node at the bottom of the graph, and change the Format to **Parquet** in the dropdown. Under *Compression Type*, select **Uncompressed** from the dropdown.
9. Under "S3 Target Location", select "**Browse S3**" browse to the "mod-xxx-dmslabs3bucket-xxx" bucket, select "**tickets**" item and press "**Choose**".



10. In the textbox, append **dms_parquet/sport_team/** to the S3 url. The path should look similar to s3://mod-xxx-dmslabs3bucket-xxx/tickets**/dms_parquet/sport_team/** - don't forget the "/" at the end. The job will automatically create the folder.



11. Finally, select the **Job details** tab at the top. Enter **Glue-Lab-SportTeamParquet** under Name.
12. For "**IAM Role**", select the role named similar to mod-xxx-**GlueLabRole-**xxx.

14

13. Scroll down the page and under "**Job bookmark**", select "**Disable**" in the drop down. You can try out the bookmark functionality later in this lab.



14. Press the "**Save**" button in the top right-hand corner to create the job.

15. Once you see the "**Successfully created job**" message in the banner, click the "**Run**" button to start the job.
16. Select "**Jobs**" from the navigation panel on the left-hand side to see a list of your jobs.
17. Select "**Monitoring**" from the navigation panel on the left-hand side to view your running jobs, success/failure rates and various other statistics.



18. Scroll down to the "**Job runs**" list to verify that the ETL job has completed successfully. This should take about 1 minute to complete.



19. We need to repeat this process for an additional 4 jobs, to transform the **sport_location, sporting_event, sporting_event_ticket** and **person** tables.

    During this process, we will need to modify different column data types. We can either repeat the process above for each table, or we can clone the first job and update the details. The steps below describe how to clone the job - if creating manually each time, follow the above steps but make sure you use the updated values from the tables below.

20. Return to the "**Jobs**" menu, and select the "**Glue-Lab-SportsTeamParquet**" job by clicking the small circle next to the name.

21. Under the "**Actions**" dropdown, select "**Clone job**". Update the job as per the following tables, then "**Save**" and "**Run**".

## 1. Sport_Location:

Create a **Glue-Lab-SportLocationParquet** job with the following attributes:

| Task / Action | Attribute | Values |
|---|---|---|
| "Data source - S3 bucket" node | Database | ticketdata |
| | Table | sport_location |
| "Transform - ApplyMapping" node | Schema transformations | None |
| "Data target - S3 bucket" node | Format | Parquet |
| | Compression Type | Uncompressed |
| | S3 target path | `tickets/dms_parquet/sport_location/` |
| "Job details tab" | Job Name | `Glue-Lab-SportLocationParquet` |
| | IAM Role | xxx-GlueLabRole-xxx |
| | Job bookmark | Disable |

## 2. Sporting_Event:

Create a **Glue-Lab-SportingEventParquet** job with the following attributes:

| Task / Action | Attribute | Values |
|---|---|---|
| "Data source - S3 bucket" node | Database | ticketdata |
| | Table | sporting_event |
| "Transform - ApplyMapping" node | Schema tranformations | column "start_date_time" => TIMESTAMP |
| | | column "start_date" => DATE |
| "Data target - S3 bucket" node | Format | Parquet |
| | Compression Type | Uncompressed |
| | S3 target path | `tickets/dms_parquet/sporting_event/` |
| "Job details tab" | Job Name | `Glue-Lab-SportingEventParquet` |
| | IAM Role | xxx-GlueLabRole-xxx |
| | Job bookmark | Disable |

### 3. Sporting_Event_Ticket:

Create a **Glue-Lab-SportingEventTicketParquet** job with the following attributes:

| Task / Action | Attribute | Values |
|---|---|---|
| "Data source - S3 bucket" node | Database | ticketdata |
| | Table | sporting_event_ticket |
| "Transform - ApplyMapping" node | Schema tranformations | column "id" => DOUBLE |
| | | column "sporting_event_id" => DOUBLE |
| | | column "ticketholder_id" => DOUBLE |
| "Data target - S3 bucket" node | Format | Parquet |

| Task / Action | Attribute | Values |
|---|---|---|
| | Compression Type | Uncompressed |
| | S3 target path | `tickets/dms_parquet/sporting_event_ticket/` |
| "Job details tab" | Job Name | `Glue-Lab-SportingEventTicketParquet` |
| | IAM Role | xxx-GlueLabRole-xxx |
| | Job bookmark | Disable |

*4. Person:*

Create a **Glue-Lab-PersonParquet** job with the following attributes:

| Task / Action | Attribute | Values |
|---|---|---|
| "Data source - S3 bucket" node | Database | ticketdata |
| | Table | person |
| "Transform - ApplyMapping" node | Schema tranformations | column "id" => DOUBLE |
| "Data target - S3 bucket" node | Format | Parquet |
| | Compression Type | Uncompressed |
| | S3 target path | `tickets/dms_parquet/person/` |
| "Job details tab" | Job Name | `Glue-Lab-PersonParquet` |
| | IAM Role | xxx-GlueLabRole-xxx |
| | Job bookmark | Disable |

## Create Glue Crawler for Parquet Files

1. In the Glue Studio naviation menu, select **Crawlers** to open the Glue Crawlers page in a new tab. Click **Add crawler**.



2. For **Crawler name**, type **glue-lab-parquet-crawler** and Click **Next**.



3. In next screen **Specify crawler source type,** select **Data Stores** as choice **for Crawler source type** and click **Next.**
4. In Add a data store screen
   a. For **Choose a data store**, select "S3".
   b. For **Crawl data** in, select "**Specified path in my account**".
   c. For Include path, specify the S3 Path (Parent Parquet folder) that contains the nested parquet files e.g., s3://xxx-dmslabs3bucket-xxx/tickets/dms_parquet
   d. Click **Next**.

Choose S3 path

S3
  aws-glue-assets-857995645290-us-east-1
  mod-3fccddd609114925-dmslabs3bucket-1pi6kt2gfj9jm
    tickets
      dms_parquet
      dms_sample
        dms_parquet_$folder$
  mod-3fccddd609114925-s3bucketworkgroupa-2sjh4ctxinuk
  mod-3fccddd609114925-s3bucketworkgroupb-62e80elgxzwm

Select



Add a data store

**Choose a data store**

S3

**Connection**

Select a connection

Optionally include a Network connection to use with this S3 target. Note that each crawler is limited to one Network connection so any future S3 targets will also use the same connection (or none, if left blank).

**Add connection**

**Crawl data in**

○ Specified path

**Include path**

s3://mod-3fccddd609114925-dmslabs3bucket-1pi6kt2gfj9jm/tickets/dms_parquet

All folders and files contained in the include path are crawled. For example, type s3://MyBucket/MyFolder/ to crawl all objects in MyFolder within MyBucket.

▸ Exclude patterns (optional)

Back    Next

5. For Add another data store, select **No** and Click **Next**.



**Add crawler**

⊘ Crawler info
  glue-lab-parquet-crawler
⊘ Crawler source type
  Data stores
◯ Data store
  S3: s3://dmslab-stu...
◯ IAM Role

Add another data store

○ Yes
◉ No

Back    Next

Chosen data stores
S3: s3://dmslab-stud...    ✕

6.  On the Choose an IAM role page, select **Choose an existing IAM role**.
    For IAM role, select the existing role "xxx-**GlueLabRole**-xxx" and Click **Next**.



7.  For **Frequency**, select "Run On Demand" and Click **Next**.



8.  For the crawler's output database, choose your existing database which you created earlier e.g.
    "**ticketdata**"

9.  For the **Prefix added to tables** (optional), type "**parquet_**"

10. Review the summary page and click **Finish**.

11. Click **Run Crawler**. Once your crawler has finished running, you should report that tables were added from 1 to 5, depending on how many parquet ETL conversions you set up in the previous section.



Confirm you can see the tables:

1. In the left navigation pane, click **Tables**.
2. Add the filter "parquet" to return the newly created tables.

# PART B: Glue Job Bookmark (Optional):

**\*\*Pre-requisite: Completion of CDC part of DMS Lab \*\***

## Step 1: Create Glue Crawler for ongoing replication (CDC Data)

Now, let's repeat this process to load the data from change data capture.

1. On the AWS Glue menu, select Crawlers.



2. Click **Add crawler**.
3. Enter the crawler name for ongoing replication. This name should be descriptive and easily recognized (e.g., "**glue-lab-cdc-crawler**").
4. Optionally, enter the description. This should also be descriptive and easily recognized and Click



   **Next**.
5. Choose **Data Stores** as Crawler Source Type**, Crawl all folders** and Click **Next**

## Specify crawler source type

Choose Existing catalog tables to specify catalog tables as the crawler source. The selected tables specify the data stores to crawl. This option doesn't support JDBC data stores.

**Crawler source type**

◉ Data stores
◯ Existing catalog tables

**Repeat crawls of S3 data stores**

◉ Crawl all folders
◯ Crawl new folders only

Only Amazon S3 folders that were added since the last crawl will be crawled. If the schemas are compatible, new partitions will be added to existing tables.

Back    Next

6. On the Add a data store page, make the following selections:
   a. For **Choose a data store**, click the drop-down box and select **S3**.
   b. For **Crawl data in**, select **Specified path in my account**.
   c. For **Include path**, enter the **target folder** for your DMS ongoing replication, e.g., "s3://xxx-dmslabs3bucket-xxx/**cdc/dms_sample**"
7. Click **Next**.

## Add a data store

**Choose a data store**

S3

**Connection**

Select a connection

Optionally include a Network connection to use with this S3 target. Note that each crawler is limited to one Network connection so any future S3 targets will also use the same connection (or none, if left blank).

**Add connection**

**Crawl data in**

◉ Specified path

**Include path**

s3://mod-3fccddd609114925-dmslabs3bucket-1pi6kt2gfj9jm/cdc/dms_sample

All folders and files contained in the include path are crawled. For example, type s3://MyBucket/MyFolder/ to crawl all objects in MyFolder within MyBucket.

▸ Exclude patterns (optional)

Back    Next

8. On the **Add another data store page**, select **No** and Click **Next**.

**Add crawler**

⊘ Crawler info
   glue-lab-cdc-crawler
⊘ Crawler source type
   Data stores
◯ Data store
   S3: s3://mod-3fccd

Add another data store

◯ Yes
◉ No

Back    Next

9. On the **Choose an IAM role** page, make the following selections:

a. Select **Choose an existing IAM role**.
b. **For IAM role**, select **xxx-GlueLabRole-xxx**. E.g. "dmslab-student-GlueLabRole-ZOQDII7JTBUM"

10. Click **Next**.

## Choose an IAM role

The IAM role allows the crawler to run and access your Amazon S3 data stores. Learn more

○ Update a policy in an IAM role
◉ Choose an existing IAM role
○ Create an IAM role

**IAM role** ⓘ

| mod-3fccddd609114925-GlueLabRole-7OEMGU9C9TZ7 | ⌄ |

This role must provide permissions similar to the AWS managed policy, **AWSGlueServiceRole**, plus access to your data stores.

• s3://mod-3fccddd609114925-dmslabs3bucket-1pi6kt2gfj9jm/cdc/dms_sample

You can also create an IAM role on the IAM console.

Back    Next

11. On the Create a schedule for this crawler page, for Frequency, select **Run on demand** and Click **Next**.

## Create a schedule for this crawler

**Frequency**

| Run on demand | ⌄ |

Back    Next

12. On the Configure the crawler's output page, select the existing **Database** for crawler output (e.g., "**ticketdata**").

13. For **Prefix added to tables,** specify "**cdc_**"

14. For Configuration options (optional), keep the **default** selections and click **Next**.

15. Review the summary page noting the Include path and Database target and Click **Finish**. The crawler is now ready to run.



16. Tick the crawler name "**glue-lab-cdc-crawler**", click **Run crawler** button.

17. When the crawler is completed, you can see it has "Status" as **Ready,** Crawler will change status from starting to stopping, wait until crawler comes back to ready state, you can see that it has created **2 tables**.

18. Click the database name (e.g., "**ticketdata**") to browse the tables. Specify "**cdc**" as the filter to list only newly imported tables.



## Step 2: Create a Glue Job with Bookmark Enabled

1. On the left-hand side of Glue Console, click on **Jobs** and then Click on **Add Job.**

2. On the Job properties page, make the following selections:
   a. For **Name**, type **Glue-Lab-TicketHistory-Parquet-with-bookmark**
   b. For **IAM role**, choose existing role "xxx-**GlueLabRole**-xxx"
   c. For **Type**, Select **Spark**
   d. For **Glue Version**, select **Spark 2.4, Python 3 (Glue version 2.0)** or whichever is the latest version
   e. For **This job runs**, select **A proposed script generated by AWS Glue**.
   f. For **Script file name**, use the **default**.
   g. For **S3 path where the script is stored**, provide a unique Amazon S3 path to store the scripts. (You can keep the **default** for this lab.)
   h. For **Temporary directory**, provide a unique Amazon S3 directory for a temporary directory. (You can keep the **default** for this lab.)

3. Expand the **Advanced properties** section. For Job bookmark, select **Enable** from the drop-down option.
4. Expand on the **Monitoring** options, enable **Job metrics**.
5. Click **Next**

## Configure the job properties

**Name**

Glue-Lab-TicketHistory-Parquet-with-bookmark

**IAM role** ℹ️

mod-3fccddd609114925-GlueLabRole-7OEMGU9C9TZ7

Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job. Create IAM role.

**Type**

Spark

**Glue version**

Spark 2.4, Python 3 with improved job startup times (Glue Version 2.0)

**This job runs**

🔘 A proposed script generated by AWS Glue ℹ️
⭕ An existing script that you provide
⭕ A new script to be authored by you

**Script file name**

Glue-Lab-TicketHistory-Parquet-with-bookmark

**S3 path where the script is stored**

s3://aws-glue-scripts-857995645290-us-east-1/admin

**Temporary directory** ℹ️

s3://aws-glue-temporary-857995645290-us-east-1/admin

▾ Advanced properties

**Job bookmark** ℹ️

Enable

▾ Monitoring options

☑️ Job metrics ℹ️

☐ Continuous logging

☐ Spark UI ℹ️

▸ Tags (optional)

▸ Security configuration, script libraries, and job parameters (optional)

▸ Catalog options (optional)

6. In **Choose a data source**, select **cdc_ticket_purchase_hist** as we are generating new data entries for **ticket_purchase_hist** table. Click **Next**



7. In **Choose a transform type**, select **Change Schema** and Click **Next**

Lab 2. ETL with AWS Glue

## Choose a transform type

Machine learning transforms are currently not supported for Glue 2.0.

● Change schema
  Change schema of your source data and create a new target dataset
○ Find matching records
  Use machine learning to find matching records within your source data

Back    Next

8. In Choose a data target:
   a. Create tables in your data target
   b. For **Data store**: select **Amazon S3**
   c. Format: **parquet**
   d. **Target path**: s3://xxx-dmslabs3bucket-xxx/**cdc_bookmark/ticket_purchase_history/data/**
   e. Click **Next**

## Choose a data target

● Create tables in your data target
○ Use tables in the data catalog and update your data target

Data store
  Amazon S3

Format
  Parquet

Connection
  - Select one -

  Add connection

Target path
  oi6kt2gfj9jm/cdc_bookmark/ticket_purchase_history/data/

Back    Next

9. In map the source columns to target columns window, leave everything as **default** and Click on **Save job and edit script**.

10. In the next window, review the job script and click on **Run job**, then click on **close mark** on the top right of the window to close the screen.



11. Once the job finishes its run, check the **S3 bucket** for the parquet partitioned data.



## Step 3: Create Glue crawler for Parquet data in S3

1. Once you have the data in S3 bucket, navigate to **Glue Console** and now we will crawl the parquet data in S3 to create data catalog.
2. Click on **Add crawler**

| | Name | Schedule | Status | Logs |
|---|---|---|---|---|
| | glue-lab-crawler | | Ready | Logs |
| | glue-lab-parquet-crawler | | Ready | Logs |

Tables
Connections
Crawlers
Classifiers
Schema registries
Schemas
Settings

Crawlers  A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for

Add crawler   Run crawler   Action ▾   Filter by tags and attributes

3.   In crawler configuration window, provide crawler name as **glue_lab_cdc_bookmark_crawler** and Click **Next**.

### Add information about your crawler

**Crawler name**

glue_lab_cdc_bookmark_crawler

▸ Tags, description, security configuration, and classifiers (optional)

Next

4.   In **Specify crawler source type**, select **Data stores** and **Crawl all folders**. Click **Next**

### Specify crawler source type

Choose Existing catalog tables to specify catalog tables as the crawler source. The selected tables specify the data stores to crawl. This option doesn't support JDBC data stores.

**Crawler source type**

◉ Data stores
○ Existing catalog tables

**Repeat crawls of S3 data stores**

◉ Crawl all folders
○ Crawl new folders only

Only Amazon S3 folders that were added since the last crawl will be crawled. If the schemas are compatible, new partitions will be added to existing tables.

Back   Next

5.   In **Add a data store**:
   a.   For **Choose a data store**, select **S3**
   a.   For the **Include path**, click the folder icon and choose your target S3 bucket, then append **/cdc_bookmark/ticket_purchase_history** , e.g., "s3://xxx-dmslabs3bucket-xxx/cdc_bookmark/ticket_purchase_history"
6.   Click on **Next**

7. For **Add another data** store, select **No** and click **Next**.



8. In **Choose an IAM role**, select an existing IAM role contains **GlueLabRole** text. Something looks like this: xxx-**GlueLabRole**-xxx



9. For setting the **frequency** in create a schedule for this crawler, select "**Run on demand**". Click **Next**

10. For the crawler's output:
   a. For Database, select "**ticketdata**" database.

        b.   Optionally, add prefix to the newly created tables for easy identification. Provide the prefix as **bookmark_parquet_**

        c.   Click **Next**



11. Review all the details and click on **Finish**. Then **Run crawler**.



12. After the crawler finishes running, click on Databases, select "**ticketdata**" and view tables in this database. You will find the newly created table as "**bookmark_parquet_ticket_purchase_history**"
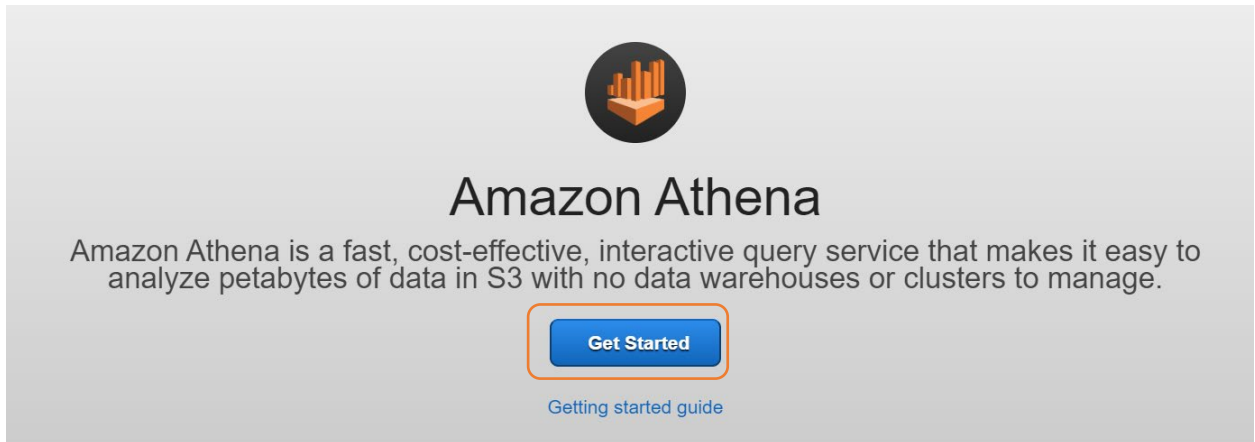


**13.** Once the table is created, click on **Action** and from **dropdown** select **View Data.**

If it's the first time you are using Athena in your AWS Account, click **Get Started**

Then click **set up a query result location in Amazon S3** at the top



In the pop-up window in the **Query result location** field, enter your s3 bucket location followed by **/**, so that it looks like **s3://xxx-dmslabs3bucket-xxx/** and click **Save**



To select some rows from the table, try running:

SELECT * FROM "ticketdata"."bookmark_parquet_ticket_purchase_history" limit 10;

To get a row count, run:

Before moving on to next step, note the rowcount.

## Step 4: Generate CDC data and to observe bookmark functionality

Ask your instructor generate more CDC data at source database, if you ran the instructor setup on your own, then make sure to follow "**Generate the CDC Data**" section from instructor prelab.

1. To make sure the new data has been successfully generated, check the S3 bucket for cdc data, you will see new files generated. Note the time when the files were generated.



2. Rerun the Glue job **Glue-Lab-TicketHistory-Parquet-with-bookmark** you created in Step 2

3. Go to the Athena Console, and rerun the following query to notice the increase in row count:

SELECT count(*) as recordcount FROM "ticketdata"."bookmark_parquet_ticket_purchase_history";

To review the latest transactions, run:

SELECT * FROM "ticketdata"."bookmark_parquet_ticket_purchase_history" order by transaction_date_time desc limit 100;

# PART C: Glue Workflows (Optional, self-paced)

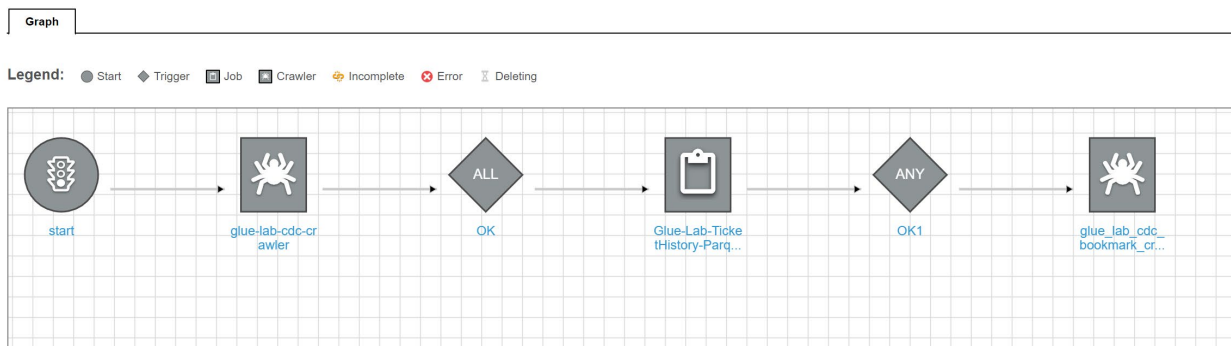**\*\*Pre-requisite before creating workflow\*\* - completed Part B**

## Overview:

In AWS Glue, you can use workflows to create and visualize complex extract, transform, and load (ETL) activities involving multiple crawlers, jobs, and triggers. Each workflow manages the execution and monitoring of all its components. As a workflow runs each component, it records execution progress and status, providing you with an overview of the larger task and the details of each step. The AWS Glue console provides a visual representation of a workflow as a graph.

## Creating and Running Workflows:

Above mentioned Part A (ETL with Glue) and Part B (Glue Job Bookmarks) can be created and executed using workflows. Complex ETL jobs involving multiple crawlers and jobs can also be created and executed using workflows in an automated fashion. Below is a simple example to demonstrate how to create and run workflows.

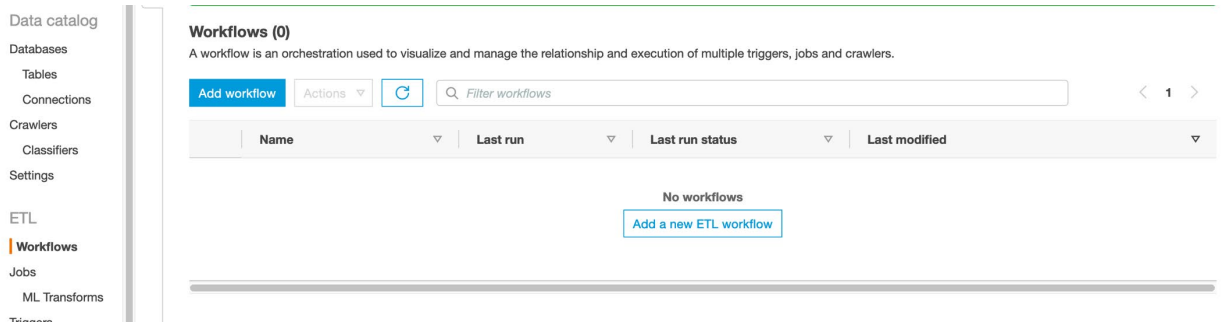

Try creating a new Glue Workflow to string together the two Crawlers and one Job from part B as follows:

On-demand trigger -> glue-lab-cdc-crawler -> Glue-Lab-TicketHistory-Parquet-with-bookmark -> glue_lab_cdc_bookmark_crawler

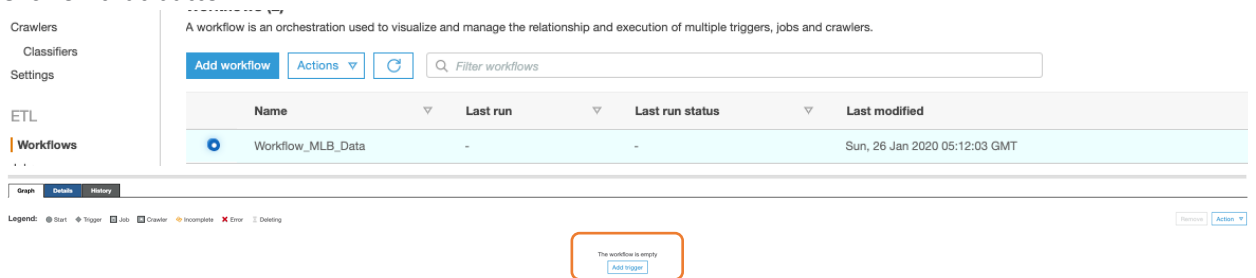


**To create a workflow:**

1. Navigate to **AWS Glue Console** and under **ETL**, click on **Workflows**. Then Click on **Add Workflow**.

2. Give the workflow name as "**Workflow_tickethistory**". Provide a description (optional) and click on **Add Workflow** to create it.

3. Click on the **workflow** and scroll to the bottom of the page. You will see an option **Add Trigger**. Click on that button.



4. In **Add Trigger** window, From Clone Existing and Add New options, click on **Add New**.
   a. Provide **Name** as "**trigger1**"
   b. Provide a **description**: Trigger to start workflow
   c. **Trigger type**: **On-demand**.
   d. Click on **Add**

   Triggers are used to initiate the workflow and there are multiple ways to invoke the trigger. Any scheduled operation or any event can activate the trigger which in turn starts the workflow



5. Click on **trigger1** to add a **new node**. New Node can be a crawler or job, depending upon the workflow you want to build.
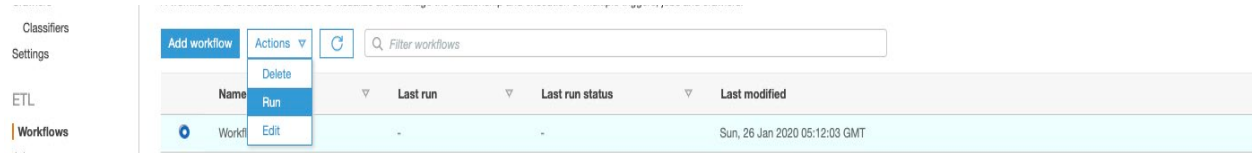
6. Click on **Add node,** a new window to add jobs or crawlers will open. Select the Crawler **glue-lab-cdc-crawler**, then **Add.**
7. Click on the crawler and **Add Trigger** provide the following:
    a. **Name**: **trigger2**
    b. **Description**: Trigger to execute job
    c. **Trigger type**: **Event**
    d. **Trigger logic**: **Start after ALL watched event.** This will make sure that job starts once Glue Crawler finishes.
    e. Click **Add**



8. After **trigger2** is added to workflow, Click on **Add node,** select job **Glue-Lab-TicketHistory-Parquet-with-bookmark,** click **Add.**
9. Click on the job and **Add Trigger** provide the following:
    a. **Name**: **trigger3**
    b. **Description**: Trigger to execute crawler
    c. **Trigger type**: **Event**
    d. **Trigger logic**: **Start after ANY watched event.** This will make sure that crawler starts once Glue job finishes processing of ALL data.
    e. Click **Add**

10. Click on **Add node,** Select the Crawler **glue_lab_cdc_bookmark_crawler**, then **Add.**
11. Select your workflow, click on **Actions->Run** and this will start the first trigger "trigger1"



12. Once the workflow is completed, you will observe that glue job and crawlers have been successfully executed.

Congratulations!! You have successfully completed this lab