

Guided Lab: Amazon Athena Data Querying and Table Creation

Description

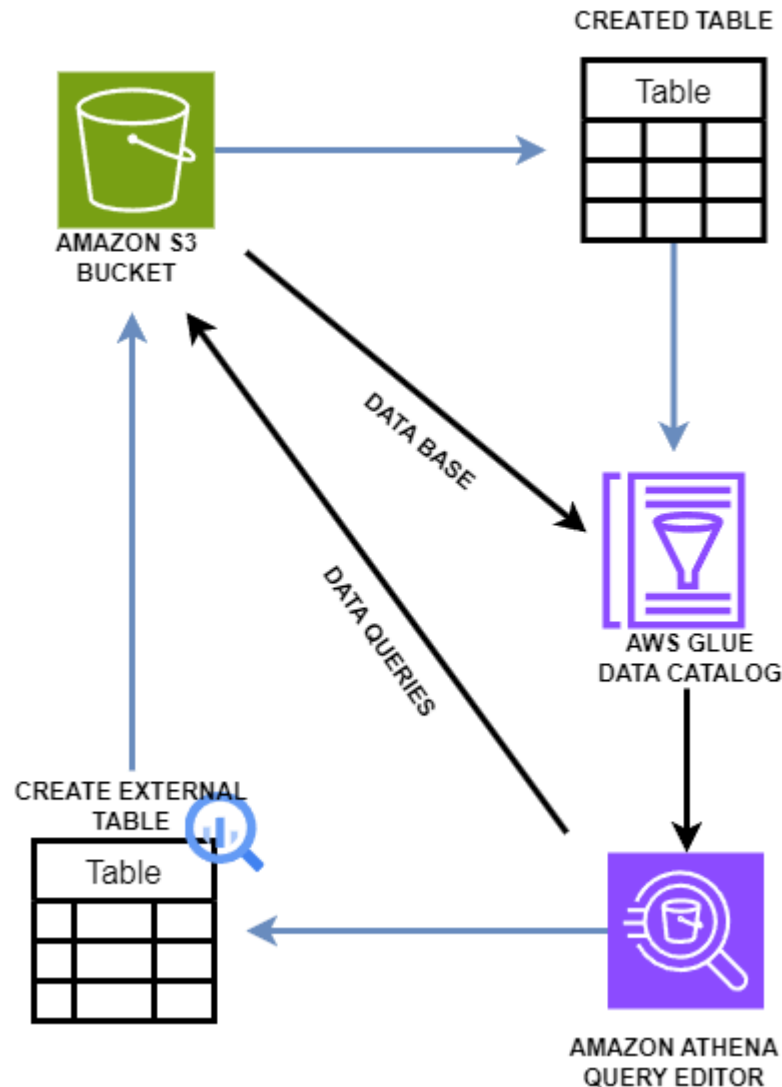
Welcome to our guided lab, where we'll dive into the world of data analytics using Amazon Athena, AWS S3, and the AWS Glue Data Catalog. This hands-on session is designed to introduce you to the power and flexibility of analyzing structured and semi-structured data stored in S3, utilizing the serverless interactive query service provided by Athena. Whether you're a data analyst, engineer, or just a data enthusiast, this lab will equip you with the knowledge to efficiently query and analyze data at scale.

Overview of Steps

In this lab, you'll go through a series of steps that will take you from setting up your environment to executing queries that will provide insights into your data. Here's a brief outline of what we'll cover:

1. **Environment Setup:** You'll start by creating an S3 bucket and uploading an activity log file, setting the stage for our analysis.
2. **AWS Glue Data Catalog Configuration:** Next, we'll create a database using the AWS Glue Data Catalog. This serves as a central metadata repository that Athena will leverage to understand the structure of your data.
3. **Table Creation in Athena:** With the metadata repository in place, you'll learn how to define and create a table in Athena that maps to the structure of your activity log data stored in S3. This is a critical step that enables Athena to execute SQL queries against your data.

4. **Data Querying:** Armed with a structured view of your data, you'll run several SQL queries in Athena. These queries will range from basic data retrieval to more complex aggregation queries, designed to familiarize you with Athena's querying capabilities and help you derive meaningful insights from your dataset.



Prerequisites

This lab assumes you have experience creating Amazon S3 Bucket and are familiar with its basic components.

If you find any gaps in your knowledge, consider taking the following labs:

- Creating an Amazon S3 bucket.
- Querying Data with Amazon Athena and AWS Glue Crawler Integration.

Objectives

In this lab, you will:

- Learn how to query data directly from S3 using Amazon Athena.
- Learn how to create an external table using Amazon Athena query.
- Use AWS Glue to create a data catalog (database) for organizing data from Amazon S3.

Subscribe to access AWS
PlayCloud Labs

Lab Steps

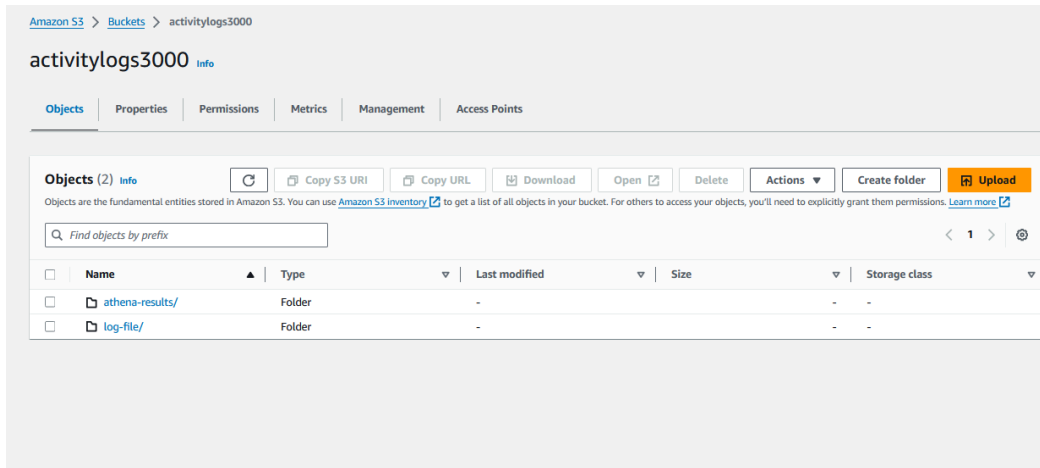
Prepare Your Environment Prepare Your Environment

1. Create an S3 Bucket:

- Give your bucket a unique name,
- Leave the rest of the settings as default and click **Create bucket**

2. Create folders in your S3 Bucket:

- Navigate to your newly created bucket.
- Create two folders
 - **athena-results**
 - **log-file**



3. Download this sample activity log file:

<https://media.tutorialsdojo.com/public/ActivityLog>

4. Upload the file in the /log-file of the s3 bucket folder you created

Set Up AWS Glue Data Catalog

1. **Navigate to AWS Glue:** In the AWS Management Console, find and select AWS Glue.

2. Create a Database:

- In the Glue Console, select **"Databases"** from the left-hand menu, then click **"Add database"**.
- Name your database uniquely, like **activitylogs3000db**, and provide a description if desired.

Create a database

Create a database in the AWS Glue Data Catalog.

Database details

Name

Database name is required, in lowercase characters, and no longer than 255 characters.

Description - *optional*

Descriptions can be up to 2048 characters long.

Database settings

Location - *optional*

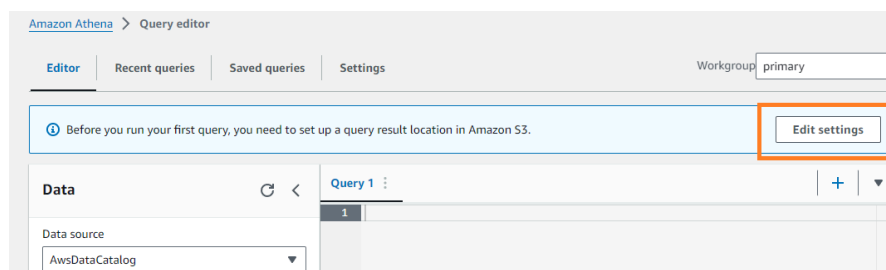
Set the URI location for use by clients of the Data Catalog.

[Cancel](#)[Create database](#)

- Click Create

Create a Table in Athena Using the Glue Data Catalog

1. **Open Athena:** Navigate to the Athena service in the AWS Console.
2. **Set Up Query Result Location in S3:**
 - In Athena, you'll see a query editor. Before executing any queries, you need to specify an S3 bucket location where query results will be saved.
 - Click on the **Edit setting**



- Click in **Manage**
- Browse and select on the S3 Bucket folder athena-results or enter in the **Location of query result**

```
s3://your-bucket-name/athena-results/
```

Amazon Athena > Query editor > Manage settings

Manage settings

Query result location and encryption

Location of query result - optional
Enter an S3 prefix in the current region where the query result will be saved as an object.

Q s3://activitylogs3000/athena-results X View Browse S3

You can create and manage lifecycle rules for this bucket
Use Amazon S3 lifecycle rules to store your query results and metadata cost effectively or to delete them after a period of time.
[Learn more](#)

Lifecycle configuration

Expected bucket owner - optional
Specify the AWS account ID that you expect to be the owner of your query results output location bucket.

Enter AWS account ID

☐ **Assign bucket owner full control over query results**
Enabling this option grants the owner of the S3 query results bucket full control over the query results. This means that if your query result location is owned by another account, you grant full control over your query results to the other account.

☐ **Encrypt query results**

Cancel Save

- Click on **Save**

3. Create Table:

- Navigate back to the **Editor's** tab, you'll create a table that references your activity log in S3. Use a DDL statement like the following, adjusting paths, table names, and column definitions as needed:

```

CREATE EXTERNAL TABLE IF NOT EXISTS
activitylogs3000db.activity_log (
    `date` date,
    `time` string,
    `location_code` string,
    `placeholder1` string,
    `ip_address` string,
    `log_type` string,
    `domain` string,
    `activity_path` string,
    `response_code` bigint,
    `placeholder2` string,
    `placeholder3` string
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\t'
STORED AS TEXTFILE
LOCATION 's3://your-bucket-name/log-file/';

```

REMEMBER TO CHANGE THE PLACEHOLDER `your-bucket-name` **IN**
THE `s3://your-bucket-name/log-file/`

- Click **Run**

The screenshot displays the Amazon Athena Query Editor interface. On the left, the 'Data' panel shows the 'Data source' as 'AwsDataCatalog', the 'Database' as 'activitylogs3000db', and a list of tables including 'activity_log'. The main editor area shows the SQL query being executed, which is the same CREATE EXTERNAL TABLE statement as shown above. Below the query editor, there are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Create'. The 'Query results' tab is active, showing a green bar indicating the query is 'Completed'. Below this, it shows 'Query successful.' and performance metrics: 'Time in queue: 90 ms', 'Run time: 428 ms', and 'Data scanned: -'. A 'Reuse query results' toggle is also visible.

Query Your Data

Now, you can **clear** the SQL Editor by clicking **Clear** bottom of the SQL editor or create a new SQL editor tab by click the **+** **(plus)** sign in the upper right of the SQL editor to add and run simple queries to analyze your activity

log. Here are a few to get you started:

1. To Select Everything from the Table

This query retrieves all the data from your activity log table. It's a good starting point to see the entire dataset you're working with.

```
SELECT * FROM "activitylogs3000db"."activity_log";
```

The screenshot shows the AWS Glue console interface. On the left, the 'Data' sidebar is visible with 'Data source' set to 'AwsDataCatalog' and 'Database' set to 'activitylogs3000db'. The 'Tables and views' section shows a table named 'activity_log'. The main panel displays the SQL query: `SELECT * FROM "activitylogs3000db"."activity_log";`. Below the query, there are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Create'. The 'Query results' tab is active, showing a 'Completed' status with 'Time in queue: 67 ms', 'Run time: 459 ms', and 'Data scanned: 4.90 KB'. The results are displayed in a table with 53 rows. The table has columns: #, date, time, location_code, placeholder 1, ip_addresses, log_type, domain, and activity_path. The data shows various activity logs from 2014-07-05.

#	date	time	location_code	placeholder 1	ip_addresses	log_type	domain	activity_path
1	2014-07-05	20:00:00	LHR3	-	10.0.0.15	ACTIVITY	example.com	/visited-cafeteria.activity
2	2014-07-05	20:00:00	MIA3	-	10.0.0.15	ACTIVITY	example.com	/attended-yoga-class.activity
3	2014-07-05	20:00:00	MIA3	-	10.0.0.15	ACTIVITY	example.com	/visited-cafeteria.activity
4	2014-07-05	20:00:00	FRA2	-	10.0.0.15	ACTIVITY	example.com	/attended-yoga-class.activity
5	2014-07-05	20:00:03	HKG1	-	10.0.0.15	ACTIVITY	example.com	/attended-yoga-class.activity
6	2014-07-05	20:00:03	HKG1	-	10.0.0.15	ACTIVITY	example.com	/visited-cafeteria.activity
7	2014-07-05	20:00:04	LAX1	-	10.0.0.15	ACTIVITY	example.com	/attended-yoga-class.activity
8	2014-07-05	20:00:04	SFO4	-	10.0.0.15	ACTIVITY	example.com	/visited-cafeteria.activity

2. To Count Total Entries

To understand the scale of your data, you can count the total number of entries in your activity log.

```
SELECT COUNT(*) AS total_entries FROM  
"activitylogs3000db"."activity_log";
```


Amazon Athena > Query editor

Editor | Recent queries | Saved queries | Settings

Workgroup: tutorials-dojo

Query 1 | Query 2 | Query 3

1 SELECT COUNT(*) AS total_entries FROM "activitylogs3000db"."activity_log";

2

Data source: AwsDataCatalog

Database: activitylogs3000db

Tables and views: Create

Filter tables and views

Tables (1): activity_log

Views (0)

SQL Ln 1, Col 75

Run again | Explain | Cancel | Clear | Create

Reuse query results up to 60 minutes ago

Query results | Query stats

Completed Time in queue: 66 ms Run time: 507 ms Data scanned: 4.90 KB

Results (1)

Search rows

#	total_entries
1	53

3. Activity Frequency

To see which activities are most common, you can count the number of occurrences for each `activity_path`.

```
SELECT activity_path, COUNT(*) AS frequency
FROM "activitylogs3000db"."activity_log"
GROUP BY activity_path
ORDER BY frequency DESC;
```

4. Activities Over Time

If you're interested in how activity changes over time, you could count the number of log entries per day.

```
SELECT date, COUNT(*) AS daily_activity_count
FROM "activitylogs3000db"."activity_log"
GROUP BY date
ORDER BY date;
```

Congratulations! You have learned how to prepare your environment, create a database and table for your data, and run queries to extract valuable insights. This knowledge will serve as a solid foundation for your future data analytics projects, empowering you to leverage AWS's scalable and serverless data analytics services efficiently. Whether for personal projects or professional tasks, the skills acquired today will enable you to tackle data analysis challenges with confidence and expertise.

One last thing! It is a good practice to clean up the resources created during this lab. Not only will it make you a better professional, but you will also become a more organized person. Happy learning!