# Guided Lab: Querying Data with Amazon Athena and AWS Glue Crawler Integration
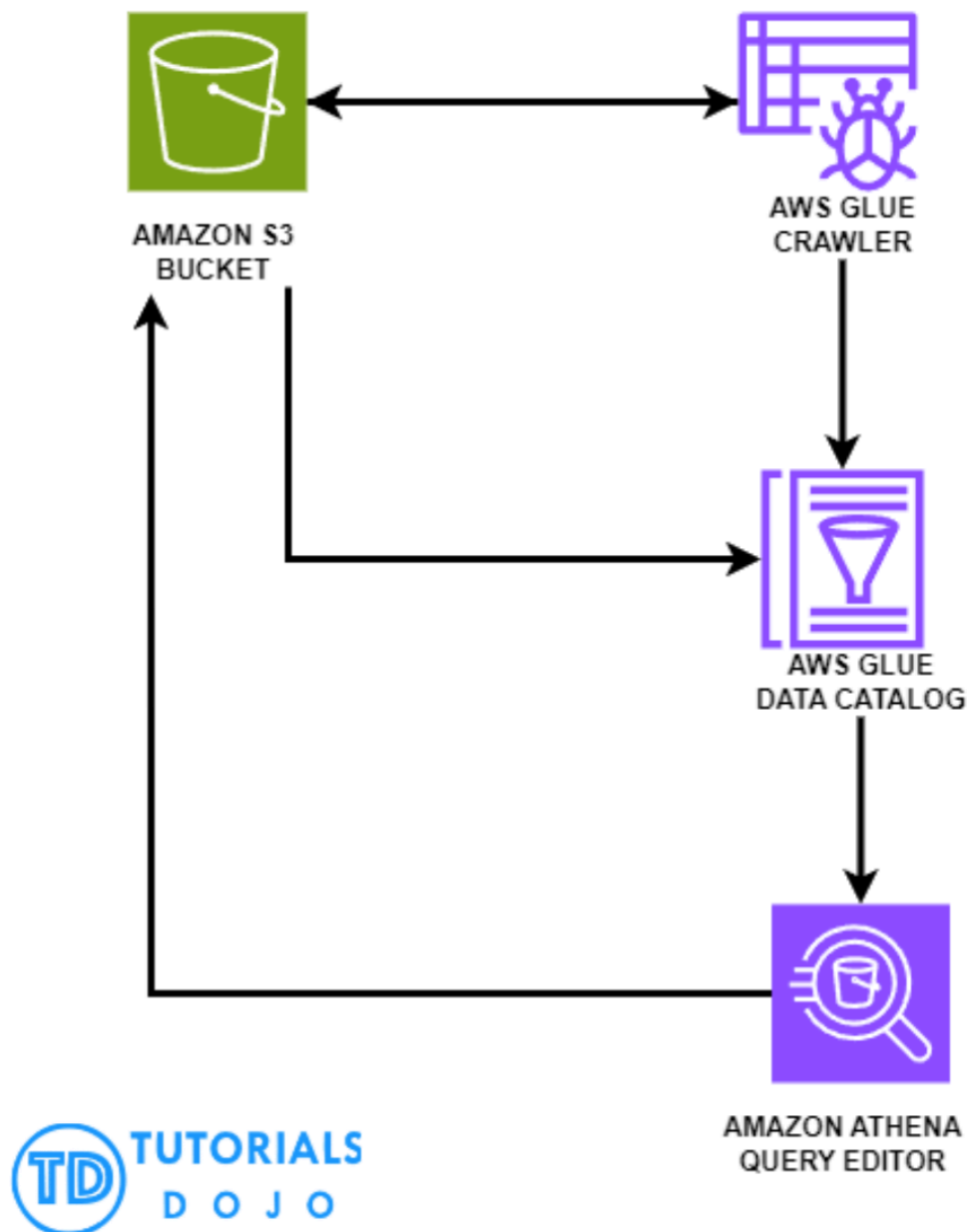
## Description

Data analytics has become an indispensable part of business strategy and decision-making. Amazon Web Services (AWS) provides a suite of scalable and flexible services designed for data analytics. Among these services, Amazon S3, Athena, and Glue (for data cataloging and data crawling) stand out for their ability to store massive datasets, query data directly in place, and organize data across various data stores efficiently.

**Overview of Steps:**

1. **Setting Up Amazon S3 Bucket**: Your data needs a place to reside. Amazon S3 serves as the foundation, providing a secure, scalable, and durable storage solution. Here, you'll store the raw data files that Athena will query.

2. **Creating a Database in AWS Glue Data Catalog**: Think of the database as a container or namespace within which you'll organize your data. It doesn't store data itself but acts as a logical grouping mechanism for your tables, which represent different datasets or aspects of your data.

3. **Adding Tables to the Database**: Tables define the schema or structure of your data (such as columns and data types) and point to the actual data stored in S3. This step is crucial because it tells Athena how to interpret the raw data during queries. You can create tables manually by defining the schema or automatically using crawlers that scan your data in S3 and infer the schema. **In this lab, we will create tables using Glue Crawler.**

4. **Querying Data with Amazon Athena**: With your data in S3, a database to organize your tables, and tables to define your data schema, you're now ready to use Athena to run SQL queries directly

against your data. Athena's serverless nature means you don't manage any infrastructure, focusing solely on analyzing your data.

## Prerequisites

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

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

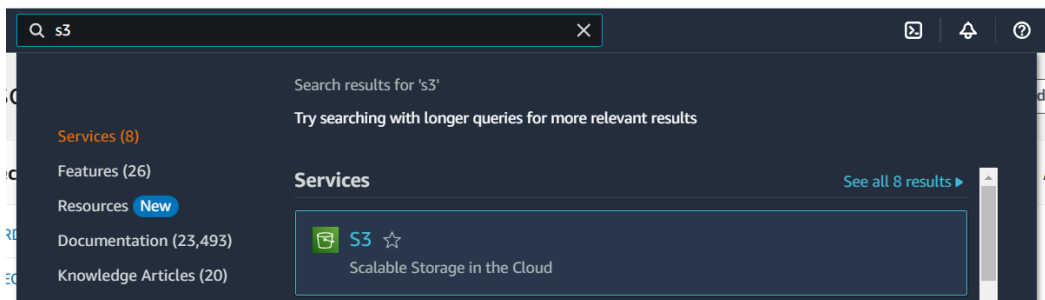- Creating an Amazon S3 bucket.

# Objectives

In this lab, you will:

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

**Subscribe to access AWS**

**PlayCloud Labs**

# Lab Steps

## Setting Up Amazon S3 Bucket

1. Log into AWS Management Console and navigate to the S3 service.



2. Create a new bucket



- Provide a unique name
- Leave everything as default settings.
- Click **Create Bucket**

3. Download this file for this lab

https://media.tutorialsdojo.com/public/Philippine_Tourist_Spots.csv

4. Create folders and Upload data files

- Create two folders and name them:
    - file
    - query



- Upload the file:
    - Open **file/** folder of your bucket.
    - Click on **Upload** and upload the file you downloaded previously.

## Setting Up AWS Glue Data Catalog

1. Navigate to the AWS Glue service in the AWS Management Console.

## 2. Creating a database

- To create a database, you need to
    - Click on **Databases** in the left corner of the window.



- Fill the **Name** with a unique database name and add a **Description** if desired.
- Click on **Create database**

## Create a database

Create a database in the AWS Glue Data Catalog.

### Database details

**Name**

playcloud3000db

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

**Description - *optional***

Enter text

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**

- You should see your newly created database afterward.

**AWS Glue** ✕

Getting started
ETL jobs
    Visual ETL
    Notebooks
    Job run monitoring
Data Catalog tables
Data connections
Workflows (orchestration)

▼ Data Catalog
  **Databases**
    Tables
  Stream schema registries
    Schemas
  Connections
  Crawlers
    Classifiers
  Catalog settings

▶ Data Integration and ETL

▶ Legacy pages

AWS Glue > Databases

## Databases (7)

Last updated (UTC) April 16, 2024 at 06:06:34

Edit     Delete     **Add database**

A database is a set of associated table definitions, organized into a logical group.

Filter databases

| | Name ▲ | Description ▽ | Location ... ▽ | Created on (UTC) ▽ |
|---|---|---|---|---|
| ☐ | playcloud3000db | - | - | April 15, 2024 at 12:20:34 |

## Setting Up AWS Glue Data Crawler

1. **Adding tables** by Glue Crawler.

- To create a Glue Crawler
  - Click on **Crawlers**

- Click on **Create crawler**
- Fill the **Name** for your crawler add description if desired.



- Click on **Next**
- Click on **Add a data source** under the Data source.
  - Add these details:
    - Data source: **S3**
    - S3 path:

```
s3://<name-of-your-s3-bucket>/file/
```

*Remember to change the placeholder <name-of-your-s3-bucket> with the name of your S3 bucket*

**Add data source**  ✕

Data source
Choose the source of data to be crawled.

| S3 | ▼ |
|---|---|

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

| | ▼ | | ⟳ |
|---|---|---|---|

**Clear selection**    **Add new connection** ⬈

⊗ Error fetching connections

Location of S3 data

● In this account
○ In a different account

S3 path
Browse for or enter an existing S3 path.

| 🔍 s3://playcloud3000/file/ | ✕ | **View** ⬈ | **Browse S3** |
|---|---|---|---|

All folders and files contained in the S3 path are crawled. For example, type s3://MyBucket/MyFolder/ to crawl all objects in MyFolder within MyBucket.

Subsequent crawler runs
This field is a global field that affects all S3 data sources.

● Crawl all sub-folders
   Crawl all folders again with every subsequent crawl.

○ Crawl new sub-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.

○ Crawl based on events
   Rely on Amazon S3 events to control what folders to crawl.

☐ Sample only a subset of files

☐ Exclude files matching pattern

Cancel    **Add an S3 data source**

- Click in **Add an S3 data source**
- Click on **Next**
- Under IAM role, Select **PlayCLoud-Sandbox**

- Click on **Next**
- Under Target Database, Select the database you created.



- Click on **Next**
- Review all the details under the **Review and create**
- Click on **Create crawler**
- You should be seeing a successful window and redirected to a window similar to the image below

## 2. Run the Crawler

- After creating the crawler, Click on **Run Crawler**

  **NOTE: I**t will take a few minutes for the crawler to crawl the AWS S3 Bucket



- When the crawler finishes crawling, you should be able to see similar image below



- In the **Data Catalog Tables**, you can see that a new table should be added.
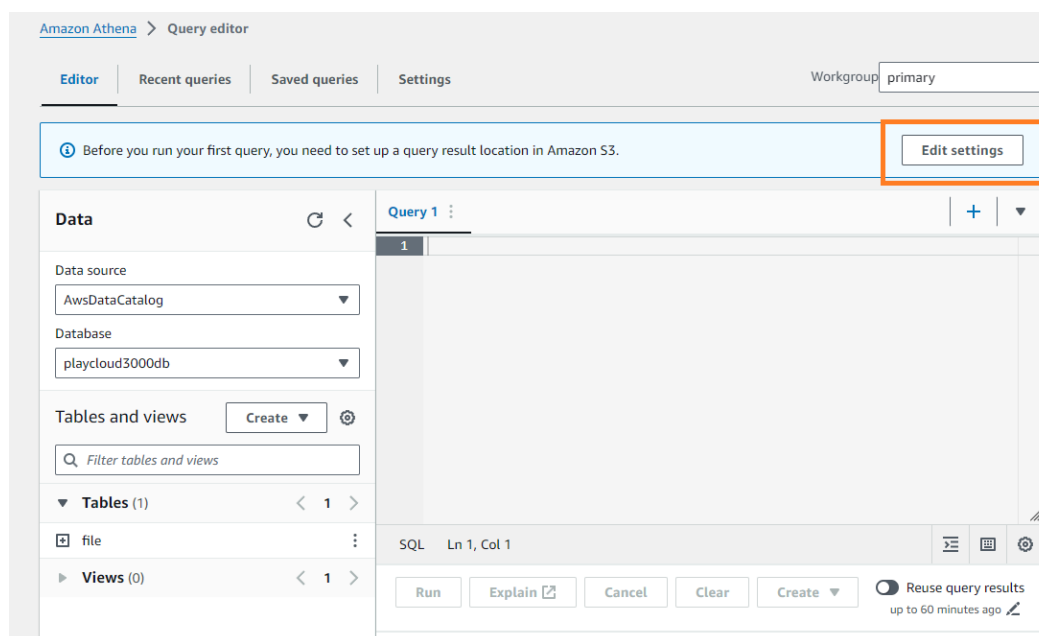
## Querying Data with Amazon Athena

1. Navigate to Amazon Athena in the AWS Management Console.



2. **Set up a query location** in Athena settings to specify an S3 bucket for storing query results.

- Click on **Edit settings**



- In the Query result location and encryption
  - Add

```
s3://<name-of-your-s3-bucket>/query/
```

Remember to change the placeholder <name-of-your-s3-bucket> with the name of your S3 bucket

- click on **Save**
- Navigate back to the Amazon Athena **Editor** tab



3. **Select the database** created in the AWS Glue Data Catalog.

- Follow the configuration below
    - Data source: **AwsDataCatalog**
    - Database: **(select the name of the database you created)**
    - The Tables should be automatically with the tables you created a while ago.

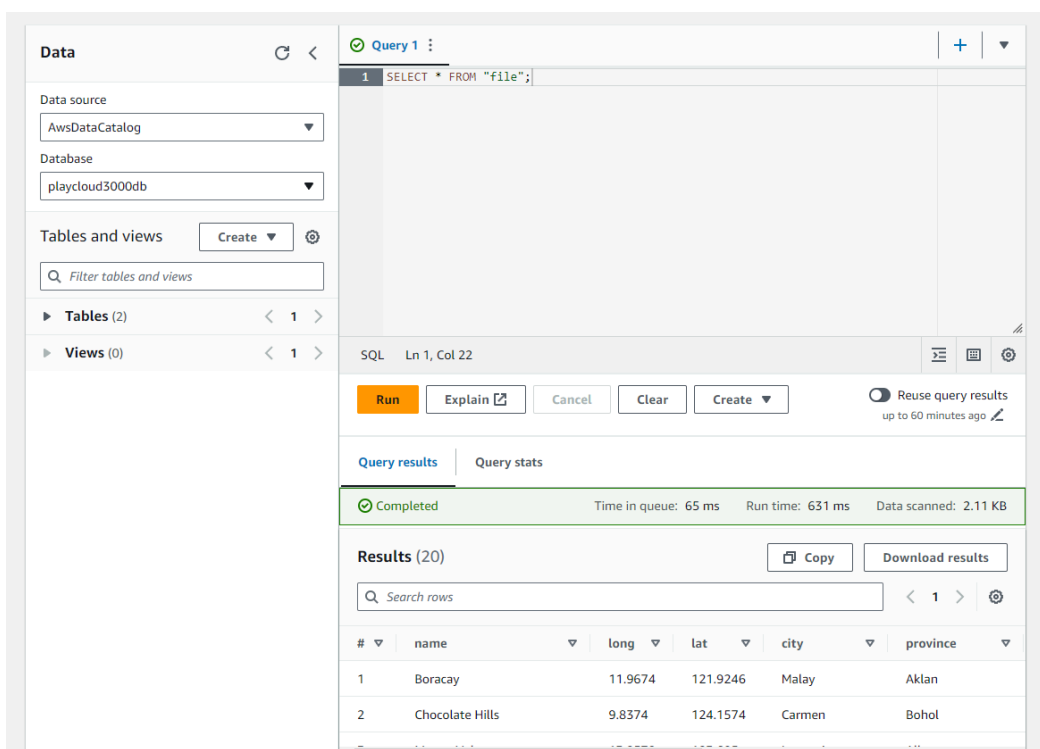4. **Write and run SQL queries** to analyze your data. You can use the standard SQL syntax.

- **To view all records:**
  - Copy & Paste. Then, **Run** the following SQL query and check the results afterward.

```
SELECT * FROM "file";
```



- You can click on **Clear** to clear the current cell contents and click on the plus button ( **+** ) to add a new cell besides the current cell

- **Filtering Records with a WHERE Clause**
  - Copy & Paste. Then, **Run** the following SQL query and check the results afterward.

```
SELECT * FROM "file"
WHERE province = 'Bohol';
```



- **Sorting Results**
  - Copy & Paste. Then, **Run** the following SQL query and check the results afterward.

```
SELECT name
FROM "file"
ORDER BY name DESC;
```

**Results** (20)　　　　　　　　　　　　　　　　□ Copy　　　Download results

| # ▽ | name | ▽ |
|---|---|---|
| 1 | Vigan | |
| 2 | Tubbataha Reef | |
| 3 | Taal Volcano | |
| 4 | Siargao Island | |
| 5 | Rizal Park | |
| 6 | Palawan Underground River | |
| 7 | Pagsanjan Falls | |
| 8 | Mount Apo | |
| 9 | Mayon Volcano | |
| 10 | Intramuros | |
| 11 | Hundred Islands | |
| 12 | Enchanted River | |
| 13 | Coron | |
| 14 | Chocolate Hills | |
| 15 | Camiguin Island | |
| 16 | Boracay | |
| 17 | Batanes | |
| 18 | Banaue Rice Terraces | |

- **Limiting Results**
  - Copy & Paste. Then, **Run** the following SQL query and check the results afterward.

```
SELECT name
FROM "file"
ORDER BY name
DESC LIMIT 5;
```

# That's it! Congratulations!

You just learned how to use AWS Glue to create a data catalog (database and tables) for organizing data from Amazon S3 and query data directly from S3 using Amazon Athena. This lab serves as a foundational step into the world of cloud-based data analytics, empowering you to explore more complex data analytics scenarios.

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!