

Querying Data by Using Athena

Lab overview and objectives

Sofía is happy with the proof of concept (POC) that you created, and the team is learning how to use it. However, one of the data scientists, Mary, would like to perform more advanced queries on comma-separated values (CSV) data. She often works with large datasets that are stored across multiple CSV files. She would like the ability to query multiple files in the same S3 bucket and aggregate the data from these files into the same database. She also wants to be able to transform the dataset metadata, such as column names and data type declarations. This information is part of the dataset's schema when the data is stored as a relational database.

Through some discovery, you learn that Amazon Athena provides this functionality. You explore Athena and its capabilities to see if you can address Mary's needs. Athena is an interactive query service that you can use to query data that is stored in Amazon S3. Athena stores data about the data sources that you query. You can store your queries for reuse, and you can share them with other users.

In this lab, you will learn how to use Athena and AWS Glue to query data that is stored in Amazon S3.

After completing this lab, you should be able to do the following:

- Use the Athena query editor to create an AWS Glue database and table.
- Define the schema for the AWS Glue database and associated tables by using the Athena bulk add columns feature.
- Configure Athena to use a dataset that is located in Amazon S3.
- Optimize Athena queries against a sample dataset.
- Create views in Athena to simplify data analysis for other users.
- Create Athena named queries by using AWS CloudFormation.

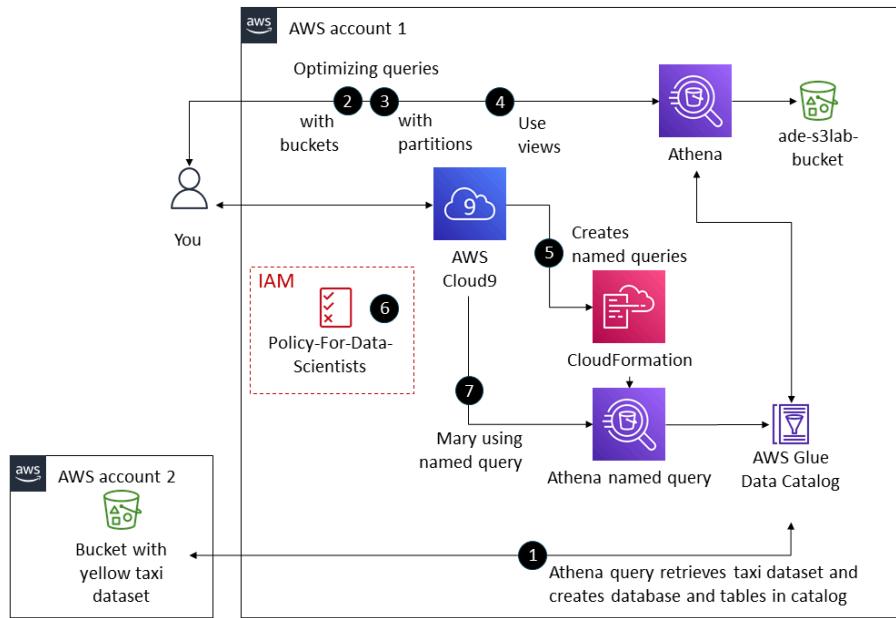
Scenario

In this lab, you will assume the role of a member of the data science team. You will build a POC using AWS Glue and Athena to analyze the data that's stored in Amazon S3. You want to experiment with Athena by accessing raw data in an S3 bucket, building an AWS Glue database, and querying this database by using Athena.

You also want to determine if you can scale this solution so that others on the team have access. Mary is part of the IAM group in AWS for the data science team and has similar access to the S3 bucket, AWS Glue database, and Athena. However, you limited her access to prevent unnecessary creation of resources. This follows the principle of least privilege, where an administrator limits access to AWS services and resources based upon what permissions are necessary to perform a specific job or task. In one task, you will use an IAM user to test the ability for a data science team member to run Athena managed queries. The IAM user was created for you, and the user belongs to an IAM group that has a policy attached to define permissions.

You ask Mary if she has a few sample datasets that you can use for your experiments. Mary provides access to a dataset that contains taxi trip data. You will use the data for your POC.

By the *end* of the lab, you will have created the additional architecture that is shown in the following diagram. The table after the diagram provides a detailed explanation of the architecture in relation to the tasks that you will complete in this lab.



Note: currently, AWS Cloud9 is not available in AWS services so we can use other IDE, for example, VScode to fix or rewrite your program.

Task 1: Creating and querying an AWS Glue database and table in Athena

Your first task is to use SQL statements to define a schema for a table that can be used to work with the sample CSV data.

In this task, you will do the following:

- Specify an S3 bucket for query results.
- Create an AWS Glue database by using the Athena query editor.
- Create a table in the AWS Glue database and import data.
- Preview the data in the AWS Glue table.

4. In the AWS Management Console, in the search box to the right of Services, search for and choose Athena to open the Athena console.

5. Specify an S3 bucket for query results.

When you use Athena, you must first specify an S3 bucket to hold the results for any queries that you run. A bucket was created for you. Complete the following steps to configure and use it.

- In the Athena console, choose **Explore the query editor**.
- Choose the **Settings** tab.
- In the **Query result and encryption settings** section, choose **Manage**.
- For **Location of query result**, choose **Browse S3**.
- Choose the bucket that was created for you, and then select **Choose**.
- Keep the defaults for the other settings on the page, and choose **Save**.

Note: At this time, you aren't required to specify an **Expected bucket owner**.

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.

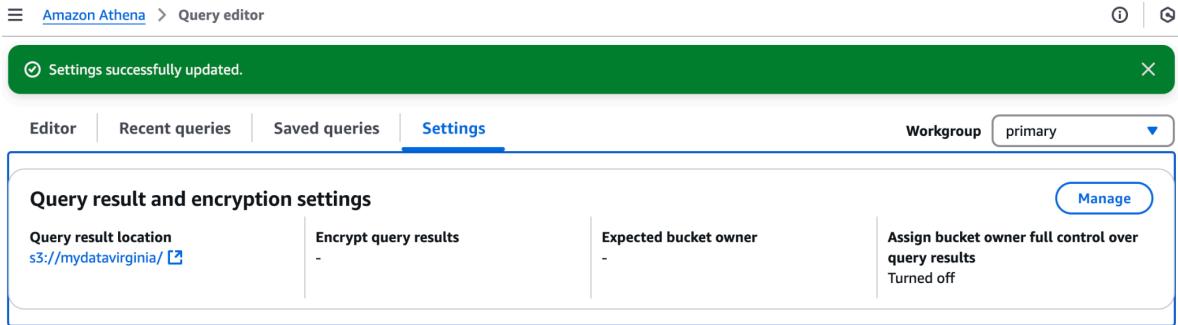
X View Browse S3

① 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.

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



6. Create an AWS Glue database by using the Athena query editor.

- Choose the **Editor** tab.
- In the **Query 1** section, enter the following SQL command:

`CREATE DATABASE taxidata;`

- Choose Run.
- The message Query successfully displays, and an AWS Glue database named taxidata is created.

Tip: To confirm that the database was created in AWS Glue, open a new tab and navigate to the AWS Glue console. In the navigation pane, choose **Databases**. The *taxidata* database is listed on the page.

A screenshot of the Amazon Athena Query Editor. On the left, the "Data" sidebar shows the "Data source" set to "AwsDataCatalog" and the "Database" set to "taxidata". The "Tables (0)" and "Views (0)" sections are shown below. On the right, the "Query 1" section contains the SQL command `CREATE DATABASE taxidata;`. A message at the top of the editor states: "Athena now supports typeahead code suggestions to speed up SQL query development. Typeahead suggestions are turned on by default. You can change this setting in query editor preferences." There is also an "Edit preferences" button.

Next, you will create a table in the AWS Glue database and import data.

7. In the Athena console, in the **Tables and views** section on the left, choose **Create > S3 bucket data**, and configure the following:

- **Table name:** Enter yellow
- **Description:** Enter Table for taxi data
- **Database configuration:** Select **Choose an existing database**, and then choose **taxidata** from the dropdown list.
- **Location of input dataset:** Copy the following link into the field:

s3://datafortable

Note: s3 location may be different.

- **Encryption:** Keep the default setting, which is not selected.
- **Data format:** Choose **CSV**.
- In the **Column details** section, choose **Bulk add columns**.
Note: This feature provides the ability to quickly add metadata to the table, such as column names and data type declarations.
- In the pop-up box, copy and paste the following text, and then choose **Add**:

```
vendor string,
pickup timestamp,
dropoff timestamp,
count int,
distance int,
ratecode string,
storeflag string,
pulocid string,
dolocid string,
paytype string,
fare decimal,
extra decimal,
mta_tax decimal,
tip decimal,
tolls decimal,
surcharge decimal,
total decimal
```

The column names and data types populate in the **Column details** section.

Amazon Athena > [Query editor](#) > Create table from S3 bucket data

Create table from S3 bucket data [Info](#)

Table details

Table name

yellow

Table name must be from 1-128 characters and must be unique. Valid characters are a-z, A-Z, 0-9, _(underscore). Table names tend to correspond to the directory where the data will be stored.

Description - optional

Table for taxi data

Table description must be from 1-1024 characters. 1005 characters remaining.

Database configuration [Info](#)

Choose an existing database or create a new database

Choose to access an existing database or to create a new database in order to create a new table. Athena stores the table schema in the AWS Glue Data Catalog.

- Create a database
- Choose an existing database

taxidata

Dataset [Info](#)

Location of input data set

s3://mydatavirginia

X

View

Browse S3

Input the path to the data set you want to process on Amazon S3. For example if your data is stored at s3://input-data-set/logs/1.csv, please enter s3://input-data-set/logs/. If your data is already partitioned, e.g. s3://input-data-set/logs/year=2004/month=12/day=11/ just input the base path s3://input-data-set/logs/

Data format [Info](#)

Table type

Apache Hive

File format

CSV

SerDe library

org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe

SerDe properties - optional

Name

field.delim

Value

,

Remove

Add SerDe property

Column details

Column name must be from 1-128 characters. Valid characters are a-z, A-Z, 0-9, _(underscore). Certain advanced column types (namely, structs) are not exposed in this interface.

Column name	Column type	Description - optional	Remove
vendor	string	Enter description	Remove
pickup	timestamp	Enter description	Remove
dropoff	timestamp	Enter description	Remove
count	int	Enter description	Remove
distance	int	Enter description	Remove

Column name	Column type	Description - optional	Remove
ratecode	string	Enter description	Remove
storeflag	string	Enter description	Remove
pulocid	string	Enter description	Remove
dolocid	string	Enter description	Remove
paytype	string	Enter description	Remove
fare	decimal	Enter description	Remove
Precision - optional	0	Scale - optional	

Column name	Column type	Description - optional
extra	decimal	<input type="text" value="Enter description"/>
Precision - optional		Scale - optional
0		
Column name	Column type	Description - optional
mta_tax	decimal	<input type="text" value="Enter description"/>
Precision - optional		Scale - optional
0		
Column name	Column type	Description - optional
tip	decimal	<input type="text" value="Enter description"/>
Precision - optional		Scale - optional
0		
Column name	Column type	Description - optional
tolls	decimal	<input type="text" value="Enter description"/>
Precision - optional		Scale - optional
0		
Column name	Column type	Description - optional
surcharge	decimal	<input type="text" value="Enter description"/>
Precision - optional		Scale - optional
0		
Column name	Column type	Description - optional
total	decimal	<input type="text" value="Enter description"/>
Precision - optional		Scale - optional
0		
Add a column Bulk add columns		

► **Table properties - optional** Info

► **Partition details - optional** Info

Partitions are a way to group specific information together. Column name must be from 1-128 characters. Valid characters are a-z, A-Z, 0-9, _(underscore). Certain advanced column types (namely, structs) are not exposed in this interface.

► **Bucketing - optional**

Preview table query

The preview table query will be populated in the query editor. Creating tables allows you to be ready for real-time querying in the query editor.

```

8   `storeflag` string,
9   `pulocid` string,
10  `dolocid` string,
11  `paytype` string,
12  `fare` decimal,
13  `extra` decimal,
14  `mta_tax` decimal,
15  `tip` decimal,
16  `tolls` decimal,
17  `surcharge` decimal,
18  `total` decimal
19 ) COMMENT "Table for taxi data"
--
```

- In the Preview table query section, view the preview table query text, which matches the following text:

```

CREATE EXTERNAL TABLE IF NOT EXISTS `taxidata`.`yellow` (
  `vendor` string,
  `pickup` timestamp,
  `dropoff` timestamp,
  `count` int,
  `distance` int,
  `ratecode` string,
  `storeflag` string,
  `pulocid` string,
  `dolocid` string,
  `paytype` string,
  `fare` decimal,
  `extra` decimal,
  `mta_tax` decimal,
  `tip` decimal,
  `tolls` decimal,
  `surcharge` decimal,
  `total` decimal
) COMMENT "Table for taxi data"

ROW FORMAT SERDE
  'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe'
WITH SERDEPROPERTIES ('field.delim' = ',')
STORED AS INPUTFORMAT
  'org.apache.hadoop.mapred.TextInputFormat'
OUTPUTFORMAT
  'org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat'
LOCATION 's3://datafortable/'

TBLPROPERTIES ('classification' = 'csv');

```

Note that the data is not encrypted. Serde is a data serialization format, and the SERDE property specifies that the data must be comma delimited, which it is.

- Choose Create table.

The message Query successful displays.

The screenshot shows the AWS Athena console interface. At the top, there are tabs for 'Editor' (which is selected), 'Recent queries', 'Saved queries', and 'Settings'. To the right, a 'Workgroup' dropdown is set to 'primary'. A message banner at the top states: 'Athena now supports typeahead code suggestions to speed up SQL query development. Typeahead suggestions are turned on by default. You can change this setting in query editor preferences.' Below the banner, the 'Data' section is visible, showing 'Data source' (AwsDataCatalog), 'Catalog' (None), and 'Database' (taxidata). The 'Tables and views' section shows a table named 'yellow' with 15 columns listed. The 'Query 1' tab is active, displaying the SQL code for creating the external table:

```

1 CREATE EXTERNAL TABLE IF NOT EXISTS `taxidata`.`yellow` (
2   `vendor` string,
3   `pickup` timestamp,
4   `dropoff` timestamp,
5   `count` int,
6   `distance` int,
7   `ratecode` string,
8   `storeflag` string,
9   `pulocid` string,
10  `dolocid` string,
11  `paytype` string,
12  `fare` decimal,
13  `extra` decimal,
14  `mta_tax` decimal,
15  `tip` decimal,

```

Below the code, the 'SQL' section shows 'Ln 1, Col 1'. At the bottom of the query area are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Reuse query results'.

The screenshot shows the 'General purpose buckets' page in the AWS S3 console. At the top, there are tabs for 'General purpose buckets' (selected) and 'Directory buckets'. A message banner at the top left says: '▶ Account snapshot - updated every 24 hours [All AWS Regions]'. It also mentions 'Storage lens provides visibility into storage usage and activity trends. Metrics don't include directory buckets.' and a link to 'Learn more'. Below the banner, the 'General purpose buckets' section shows two buckets: 'datafortable' and 'mydatavirginia'. Each bucket entry includes the name, AWS Region (US East (N. Virginia)), the creation date (February 11, 2025, 14:26:43 (UTC+07:00) and February 11, 2025, 13:26:00 (UTC+07:00)), and a link to 'View analyzer for us-east-1'.

The screenshot shows the 'Objects (1)' page for the 'datafortable' bucket in the AWS S3 console. At the top, there are tabs for 'Objects' (selected), 'Metadata', 'Properties', 'Permissions', 'Metrics', 'Management', and 'Access Points'. Below the tabs, there are buttons for 'Actions' (with options like 'Copy S3 URI', 'Copy URL', 'Download', 'Open', 'Delete', 'Actions', and 'Create folder'), and a 'Upload' button. A note below the buttons states: 'Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions.' A link to 'Learn more' is provided. The main table lists the single object: 'generated_taxi_data.csv'. The table columns are Name, Type, Last modified, Size, and Storage class. The object details are: Name is 'generated_taxi_data.csv', Type is 'csv', Last modified is 'February 11, 2025, 14:31:16 (UTC+07:00)', Size is '12.3 KB', and Storage class is 'Standard'.

11/

[Copy S3 URI](#)

Objects

Properties

Objects (12)

[Copy S3 URI](#) [Copy URL](#) [Download](#) [Open](#) [Delete](#) [Actions](#)

[Create folder](#) [Upload](#)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

Name	Type	Last modified	Size	Storage class
6da59261-1688-41cc-942e-0606b6aea3dc.csv	csv	February 11, 2025, 14:44:20 (UTC+07:00)	13.2 KB	Standard
6da59261-1688-41cc-942e-0606b6aea3dc.csv.metadata	metadata	February 11, 2025, 14:44:20 (UTC+07:00)	766.0 B	Standard

Now that you have a table with the taxi data, you can write queries to retrieve data from the data source in Amazon S3.

8. Preview data in the AWS Glue table.

- In the Data section on the left, for Database, choose taxidata.
- In the Tables section, to the right of the yellow table, choose the ellipsis (three dot) icon, and then choose Preview Table.
The Results section displays the first 10 records of the table.
- Before continuing, at the top of the query editor, close the open queries by choosing the X icon on each query tab. When prompted, choose Close query.
Note: You can only have 10 active queries in Athena at any given time, so you must manage these accordingly.

Task 1 summary

In this task, you created an AWS Glue database and table by using the Athena query editor. You connected an Amazon S3 dataset to the table and defined the schema for the table by using the bulk add columns feature. After you created the table, you learned how to preview the data by using the preview table feature.

Mary can now use more advanced SQL queries and custom column data types for CSV files that are stored securely in Amazon S3.

Capture your result

- First 10 rows of the read Athena query from yellow.

Task 2: Optimizing Athena queries by using buckets

When you work with large datasets that are spread out across multiple files, two major goals are to optimize query performance and to minimize cost. Cost for Athena is based on usage, which is the amount of data that is scanned, and prices vary based on your Region.

Three possible strategies that you can use to minimize your costs and improve performance are compressing data, bucketizing data, and partitioning data.

- Compressing data: Compress your data by using one of the open standards for file compression (such as Gzip or tar). Compression results in a smaller size for the dataset when it's stored in Amazon S3.
- The cardinality of your data also effects how you should optimize your queries. For more information, see [Cardinality \(SQL Statements\)](#).
- **Cardinality in SQL** refers to either the "**number of unique values**" in a column of a table or the "**size of the relationship between two tables**" in a relational database.
 - **High Cardinality:** A column with many unique values, such as **ID numbers, phone numbers, or email addresses**.
 - **Low Cardinality:** A column with many repeated values, such as **gender (Male/Female) or marital status (Married/Single)**.
- There are two options to optimize based upon high or low cardinality. These are:
 - **Bucketizing data:** For high cardinality with data, store records in distinct buckets (not to be confused with S3 buckets) based upon a shared value in a specific field. Consider bucketizing data as part of the preprocessing phase of your data pipeline. In this lab, data for a single month will be bucketized separately from the original dataset, which contains data for an entire year. This strategy will help to optimize performance.
 - **Partitioning data:** You can also use partitions to improve performance and reduce cost. Partitioning is used with low-cardinality data, which means that fields have few unique or distinct values.

In this task, you will experiment with bucketizing data to optimize Athena queries. You will complete the following actions:

- Create a table named jan to hold bucketized data.
- Compare how long it takes to run a query against bucketized data for January 2025 and how long it takes to query the entire dataset for the January 2025 data.

9. Create a table for the January 2025 data.

- In the Athena query editor, choose the **Editor** tab.
- To open a new query tab, choose the plus icon on the right side of the query section.
- Copy and paste the following text into the new query tab, and then choose **Run**:

```
CREATE TABLE jan AS  
SELECT *  
FROM yellow WHERE pickup  
BETWEEN TIMESTAMP '2025-01-01 00:00:00'  
AND TIMESTAMP '2025-02-01 00:00:01'
```

The message Query successful displays. Athena created a new table named jan, which contains the taxi data for only the month of January 2025. The table is listed in the Tables section on the left.

Analysis: In this step, you created a table to store data for only January by using Serde. The fares in this dataset all have pickup timestamps, and almost all the timestamps are unique for each record (meaning the records have high cardinality), but also have January in the month column, so you can use bucketizing. The january data is stored in a separate file in Amazon S3 from the file with the full dataset. In the following steps, you will investigate whether bucketizing these records in a separate file improves performance compared to querying for January records in the full dataset.

10. Run the following query on the yellow table, which has data for the entire year. The data is not divided into monthly buckets.

```
SELECT count (count) AS "Number of trips" ,  
       sum (total) AS "Total fares" ,  
       pickup AS "Trip date"  
FROM yellow WHERE pickup  
between TIMESTAMP '2025-01-01 00:00:00'  
and TIMESTAMP '2025-02-01 00:00:01'  
GROUP BY pickup;
```

The message Query successful displays. Note the run time and amount of data scanned for the query.

Now, you will compare that by running a query against the jan table.

11. Run the following query on the jan table.

```
SELECT count (count) AS "Number of trips" ,  
       sum (total) AS "Total fares" ,  
       pickup AS "Trip date"  
FROM jan  
GROUP BY pickup;
```

Note the run time and amount of data scanned for the query. As you can see, much less data is scanned when the query is performed against the table that only contains the January data.

Analysis: Putting data in separate buckets works when you have data that has a high degree of cardinality. Cardinality refers to the number of distinct records in a database. In this example, the pickup field has a high degree of cardinality because every trip has a specific date and time.

Task 2 summary

In this task, you created a table for the January data. You learned how to optimize queries by using various strategies, such as dividing data into multiple buckets and partitioning data.

To accomplish this task, you first ran a query on data that was not divided into buckets and used this query as the baseline. Then, you ran a query on the January data, which was divided into a separate bucket, and compared the results to the baseline.

Capture your result

- The result and time from no. 11 and no 12.

Task 3: Optimizing Athena queries by using partitions

If you are interested in querying a field with low cardinality, which means that the field has few unique values, you would partition the data instead of using distinct buckets. In some cases, your data will be partitioned by another process. To find the most efficient approach, you will try to partition the data by using the *paytype* field in a specific query in Athena.

Because the set of possible values is limited, *paytype* is an excellent column to use to create partitions. In this task, you will use the CREATE TABLE AS function to partition the data. You will also specify a columnar storage format.

You can store data in Athena with the Apache Parquet or Optimized Row Columnar (ORC) formats. Columnar storage formats compress the data, which will further reduce costs for your queries.

In this task, you will experiment with using partitions with your queries. You will complete the following steps:

- Create a table that is based on the Apache Parquet format and uses the *paytype* field as a partition.
- Compare the time it takes to run a query against the *yellow* database for records with *paytype* 1 (credit card) to how long it takes to query the partitioned table that was built with the Apache Parquet format.

To begin, you will create a partition by running a query to select the data that you want to use for a partition and to indicate a storage format.

12. To create a new table called *taxidata.creditcard* that is partitioned for *paytype* = 1 (credit card transactions), run the following query in a new query tab:

```
CREATE TABLE taxidata.creditcard
WITH (
    format = 'PARQUET'
) AS
SELECT * from "yellow"
WHERE paytype = 'Credit';
```

Now, you will compare the performance of running queries on the nonpartitioned data in the *yellow* table and the partitioned data in the *creditcard* table.

13. To query the nonpartitioned data in the *yellow* table, run the following query in a new query tab:

```
SELECT sum (total), paytype FROM yellow
WHERE paytype = 'Credit' GROUP BY paytype;
```

The run time and data scanned values are similar to the following:

```
Time in queue: 70 ms
Run time: 497 ms
Data scanned: 12.25 KB
```

14. To query the partitioned data in the *creditcard* table, run the following query in a new query tab:

```
SELECT sum (total), paytype FROM creditcard
WHERE paytype = 'Credit' GROUP BY paytype;
```

The run time and data scanned values are similar to the following:

```
Time in queue: 69 ms
Run time: 418 sec
Data scanned: 0.20 KB
```

Time in queue

This refers to the amount of time a query has to wait in the queue before it starts running. Normally, if the system has other queries being processed or is under high load, Athena will queue your query.

Run time

This is the time taken to execute the query from start to finish.

Data scanned

This refers to the amount of data Athena needs to scan to process the query.

If the value is "-" (no data), it could mean:

- The query did not read actual data (e.g., commands like SHOW TABLES or DESCRIBE).
- No data matched the query's conditions.
- Partitions were used, and Athena did not need to scan all the data.
- The query might have been optimized to avoid reading actual files.

Analysis: Notice that the results of each query are the same, but the query on the partitioned data took significantly less time because it scanned less data. Remember that you are charged for the amount of data scanned for each query. So, since the last query scans less data due to the use of partitions, your cost is reduced.

Task 3 summary

In this task, you partitioned the dataset where the *paytype* field contained a specific value. Then, you ran queries on nonpartitioned and partitioned versions of the data and compared the results.

Capture your result

- The result and time from no. 13 and no 14.

Task 4: Using Athena views

You might have noticed that Athena can create views. Creating and using views with data can help to simplify analysis because you can hide some of the complexity of queries from users.

Athena supports running only one SQL statement at a time, but you can use views to combine data from various tables. You can also use views to optimize query performance by experimenting with different ways to retrieve data and then saving the best query as a view.

In this task, you will do the following:

- Create a view to calculate the total dollar value of taxi fares that were paid with a credit card.
- Create a view to calculate the total dollar value of fares that were paid with cash.
- Retrieve all records from each of these views.
- Create a new view that joins the data from these two views.
- Preview the results of the join.

15. To create a view for the total dollar value of fares that were paid with a *credit card*, run the following query:

```
CREATE VIEW cctrips AS
    SELECT "sum"("fare") "CreditCardFares"
        FROM yellow
    WHERE ("paytype"='Credit');
```

The run time and data scanned values are similar to the following:

```
Time in queue: 83 ms
Run time: 513 ms
Data scanned: -
```

16. To create a view for the total dollar value of fares that were paid with cash, run the following query:

```
CREATE VIEW cashtrips AS
    SELECT "sum"("fare") "CashFares"
        FROM yellow
    WHERE ("paytype"='Digital');
```

The run time and data scanned values are similar to the following:

Time in queue: 47 ms

Run time: 599 ms

Data scanned: -

17. To select all records from the *cctrips* view, run the following query:

```
Select * from cctrips;
```

The run time and data scanned values are similar to the following:

Time in queue: 108 ms

Run time: 650 ms

Data scanned: 12.25 KB

18. To select all records from the *cashtrips* view, run the following query:

```
Select * from cashtrips;
```

The run time and data scanned values are similar to the following:

Time in queue: 103 ms

Run time: 660 ms

Data scanned: 12.25 KB

19. To create a new view that joins the data, run the following query:

```
CREATE VIEW comparepay AS  
WITH  
    cc AS  
        (SELECT sum(fare) AS cctotal,  
         vendor  
        FROM yellow  
        WHERE paytype = 'Credit'  
        GROUP BY paytype, vendor),  
    cs AS  
        (SELECT sum(fare) AS cashtotal,  
         vendor, paytype  
        FROM yellow  
        WHERE paytype = 'Digital'  
        GROUP BY paytype, vendor)
```

```
SELECT cc.cctotal, cs.cashtotal  
FROM cc  
JOIN cs  
ON cc.vendor = cs.vendor;
```

The run time and data scanned values are similar to the following:

Time in queue: 36 ms

Run time: 547 ms

Data scanned: -

20. Preview the results from the join.

- In the **Views** section on the left, to the right of the **comparepay** view, choose the ellipsis icon, and then choose **Preview View**.
- If prompted about opening the query, choose **Open query**.
The results display and are similar to the following:

```
#    cctotal cashtotal  
1    496    612  
2    441    346
```

Task 4 summary

In this task, you learned how to create views in Athena. You created two views to calculate the total revenue of taxi fares that were paid with a credit card and paid with cash. You also learned how to use a view to join data from two other views. The skills that you learned in this task will help simplify analysis of data by using Athena.

Capture your result

- The result and time from no. 19.

Task 5: Creating Athena named queries by using CloudFormation

The data science team wants to share the queries that they built by using the taxi dataset. The team would like to share with other departments, but those departments use other accounts in AWS. Other departments have less experience with AWS, so the data science team would like to simplify the process of using Athena for other departments to use.

Note: In this lab, you will focus on creating a CloudFormation template to create a named query in Athena. The template will not create the AWS Glue database or the tables and data within it.

You want to simplify the process to deploy queries. After some research, you determine that the best solution is to build a CloudFormation template for each query. The template can be shared with other departments and deployed as needed.

AWS CloudFormation is a service that helps you model and set up your AWS resources so that you can spend less time managing those resources and more time focusing on your applications that run in AWS. You create a template that describes all the AWS resources that you want (like Amazon EC2 instances or Amazon RDS DB instances), and CloudFormation takes care of provisioning and configuring those resources for you. You don't need to individually create and configure AWS resources and figure out what's dependent on what; CloudFormation handles that. The following scenarios demonstrate how CloudFormation can help.

21. Review and run the example query.

Mary provides you with the following example query:

```
SELECT distance, paytype, fare, tip, tolls, surcharge, total FROM yellow  
WHERE total >= 10.0 ORDER BY total DESC
```

The template for using CloudFormation to connect Athena services as below:

```
AWSTemplateFormatVersion: 2010-09-09  
Resources:  
  AthenaNamedQuery:  
    Type: AWS::Athena::NamedQuery  
    Properties:  
      Database: "taxidata"  
      Description: "A query that selects all fares over $100.00 (US)"  
      Name: "FaresOver10DollarsUS"  
      QueryString: >  
        SELECT distance, paytype, fare, tip, tolls, surcharge,  
        total  
        FROM yellow  
        WHERE total >= 10.0  
        ORDER BY total DESC  
Outputs:  
  AthenaNamedQuery:  
    Value: !Ref AthenaNamedQuery
```

22. Select CloudFormation and then click Create stack -> With new resource (standard)

The screenshot shows the AWS CloudFormation Stacks page. At the top, there is a blue banner with the text: "Introducing the new Hooks experience and capabilities. The Hooks' new console workflow simplifies how you author your Hooks by using Lambda functions or Guard domain-specific language (DSL). You can also learn more about the expanded evaluation targets." Below the banner, the page title is "CloudFormation > Stacks". There are buttons for "Delete", "Update", "Stack actions", and "Create stack". A dropdown menu for "Filter status" shows "With new resources (standard)" is selected. The main table lists two stacks:

Stack name	Status	Created time	Description
test	CREATE_COMPLETE	2025-02-11 22:56:05 UTC+0700	-
c145274a3763773l8920064t1w671186299989	CREATE_COMPLETE	2025-01-09 21:50:28 UTC+0700	associate Learner Lab template (academy)

23. Step 1: Prerequisite - Prepare template

Select "Choose an existing template"

The screenshot shows the "Create stack" wizard at Step 1: Prerequisite - Prepare template. On the left, a sidebar shows steps: Step 1 (selected), Step 2, Step 3, Step 4. The main area has a title "Create stack" and a section "Prerequisite - Prepare template". It says: "You can also create a template by scanning your existing resources in the IaC generator." Below are two options:

- Choose an existing template
Upload or choose an existing template.
- Build from Infrastructure Composer
Create a template using a visual builder.

And select "Specify template" after that upload your template (YAML file)

Specify template Info

This [GitHub repository](#) contains sample CloudFormation templates that can help you get started on new infrastructure projects. [Learn more](#)

Template source

Selecting a template generates an Amazon S3 URL where it will be stored. A template is a JSON or YAML file that describes your stack's resources and properties.

Amazon S3 URL

Provide an Amazon S3 URL to your template.

Upload a template file

Upload your template directly to the console.

Sync from Git

Sync a template from your Git repository.

Upload a template file

[Choose file](#)

JSON or YAML formatted file

S3 URL: Will be generated when template file is uploaded

[View in Infrastructure Composer](#)

Then, press “Next”

Specify template Info

This [GitHub repository](#) contains sample CloudFormation templates that can help you get started on new infrastructure projects. [Learn more](#)

Template source

Selecting a template generates an Amazon S3 URL where it will be stored. A template is a JSON or YAML file that describes your stack's resources and properties.

Amazon S3 URL

Provide an Amazon S3 URL to your template.

Upload a template file

Upload your template directly to the console.

Sync from Git

Sync a template from your Git repository.

Upload a template file

[Choose file](#)

athenaquery.cf.yaml

JSON or YAML formatted file

S3 URL: <https://s3.us-east-1.amazonaws.com/cf-templates-dv2sub6brjrd-us-east-1/2025-02-12T024747.675Z0jt-athenaquery.cf.yaml>

[View in Infrastructure Composer](#)

[Cancel](#)

[Next](#)

Step 2: Put your a stack name “queryforathena”

① Introducing the new Hooks experience and capabilities
The Hooks' new console workflow simplifies how you author your Hooks by using Lambda functions [2] or Guard domain-specific language (DSL) [2]. You can also learn more [2] about the expanded evaluation targets.

View Hooks overview X

Step 1 Create stack
Step 2 Specify stack details
Step 3 Configure stack options
Step 4 Review and create

Specify stack details

Provide a stack name

Stack name
Stack name must be 1 to 128 characters, start with a letter, and only contain alphanumeric characters. Character count: 0/128.

Parameters
Parameters are defined in your template and allow you to input custom values when you create or update a stack.
No parameters
There are no parameters defined in your template

Cancel Previous Next

Step 3: Configure stack options - using default values.

Step 1 Create stack
Step 2 Specify stack details
Step 3 Configure stack options
Step 4 Review and create

Configure stack options

Tags - optional
Tags (key-value pairs) are used to apply metadata to AWS resources, which can help in organizing, identifying, and categorizing those resources. You can add up to 50 unique tags for each stack.
No tags associated with the stack.
[Add new tag](#)
You can add 50 more tag(s)

Permissions - optional
Specify an existing AWS Identity and Access Management (IAM) service role that CloudFormation can assume.
IAM role - optional
Choose the IAM role for CloudFormation to use for all operations performed on the stack.
 [Remove](#)

Step 4: Review and create - using default values.

Step 1

- Create stack
- Specify stack details
- Configure stack options
- Review and create**

Review and create

Step 1: Specify template

Prerequisite - Prepare template

Template
Template is ready

Template

Template URL
<https://s3.us-east-1.amazonaws.com/cf-templates-dv2sub6brjrd-us-east-1/2025-02-12T024747.675Z0jt-athenaquery.cf.yaml>

Stack description
-

Finish

Stacks (2)

Filter status Active View nested

Stacks
queryforathena
2025-02-12 10:00:26 UTC+0700
CREATE_COMPLETE
c145274a3763773189200641w671186299989
2025-01-09 21:50:28 UTC+0700
CREATE_COMPLETE

queryforathena

Delete Update Stack actions ▾ Create stack ▾

Stack info

Events - updated Resources Outputs Parameters >

Overview

Stack ID	Description
arn:aws:cloudformation:us-east-1:671186299989:stack/queryforathena/827f51d0-e8ed-11ef-9d1e-0e5835081793	-
Status	Detailed status
CREATE_COMPLETE	-
Status reason	Root stack
User Initiated	-
Parent stack	Created time
-	2025-02-12 10:00:26 UTC+0700
	Updated time
	-
Deleted time	Drift status
-	NOT_CHECKED

23. Confirm the named query resource that the CloudFormation stack created.

```
aws athena list-named-queries
```

The output is similar to the following:

```
{
  "NamedQueryIds": [
    "644f6c10-bf57-48a2-a0ec-a3de179db511"
  ]
}
```

The named query ID is a unique identifier in AWS for the query that you created by using the CloudFormation stack.

- Copy the query ID to a text editor.
- To retrieve the details of the named query, including the SQL statement that is associated with it, run the following command. Replace <QUERY-ID> with the ID that you saved to a text editor.

```
aws athena get-named-query --named-query-id <QUERY-ID>
```

The output is similar to the following:

```
~ $ aws athena get-named-query --named-query-id d50815cd-10fb-4b5a-b110-d9d8a3d8652a
{
    "NamedQuery": {
        "Name": "FaresOver10DollarsUS",
        "Description": "A query that selects all fares over $100.00 (US)",
        "Database": "taxidata",
        "QueryString": "SELECT distance, paytype, fare, tip, tolls, surcharge, total FROM yellow WHERE total >= 10.0 ORDER BY total DESC\\n",
        "NamedQueryId": "d50815cd-10fb-4b5a-b110-d9d8a3d8652a",
        "WorkGroup": "primary"
    }
}
~ $ █
```

- To save the named query ID as a bash variable, run the following command. Replace <QUERY-ID> with the named query ID that you copied earlier. This will make it easier to use the ID in commands in later steps:

```
NQ=<QUERY-ID>
```

- Confirm the named query id was stored as the bash variable **NQ**. Run the command below.

```
echo $NQ
```

The result is similar to the following:

```
~ $ NQ=d50815cd-10fb-4b5a-b110-d9d8a3d8652a
~ $ echo $NQ
d50815cd-10fb-4b5a-b110-d9d8a3d8652a
~ $ █
```

Task 5 summary

In this task, you learned how to integrate an Athena named query into a CloudFormation template. You also learned how to validate and deploy a template to create a named query in a CloudFormation stack.

Many companies use multiple accounts with AWS to maintain separate development, testing, and production environments. Isolating these environments helps to ensure that teams follow best practices. Building queries in a development account and then testing them in a controlled account with production data can help to ensure that the query is designed as intended and generates the necessary insights. After validating the query, you can use DevOps best practices with CloudFormation to quickly move it to production so that the appropriate business stakeholders can reuse it without having to build from the beginning.

Capture your result