# ECV-BD-311

# Building a Data Lake with AWS Glue and Amazon S3

## 2018.02.23

## Version 1.1

# Agenda

# About this lab

## Scenario

The following procedures help you set up a data lake that could store and analyze data that addresses the challenges of dealing with massive volumes of heterogeneous data. A data lake allows organizations to store all their data—structured and unstructured—in one centralized repository. Because data can be stored as-is, there is no need to convert it to a predefined schema. This tutorial walks you define a database, configure a crawler to explore data in an Amazon S3 bucket, create a table, transform the CSV file into Parquet, create a table for the Parquet data, and query the data with Amazon Athena.

## AWS Glue introduction

What is AWS Glue?

AWS Glue is a fully managed data catalog and ETL (extract, transform, and load) service that simplifies and automates the difficult and time-consuming tasks of data discovery, conversion, and job scheduling. AWS Glue crawls your data sources and constructs a data catalog using pre-built classifiers for popular data formats and data types, including CSV, Apache Parquet, JSON, and more. It is significantly reducing the time and effort that it takes to derive business insights quickly from an Amazon S3 data lake by discovering the structure and form of your data. Also automatically crawls your Amazon S3 data, identifies data formats, and then suggests schemas for use with other AWS analytic services.

## Amazon S3 introduction

What is Amazon S3?

Amazon S3 is a highly durable, cost-effective object start that supports Open Data Formats while decoupling storage from compute, and it works with all the AWS analytic services.

## The workshop's region will be in 'Virginia'

## Prerequisites

+ Sign-in a AWS account, and make sure you have select N.Virginia region.

# Lab tutorial

## 1.1 Create IAM role

1.1. On the **service** menu, click **IAM**.

1.2. In the navigation pane, choose **Roles.**

1.3. Click **Create role**.

1.4. For role type, choose **AWS Service**, find and choose **Glue**, and choose **Next: Permissions**.

1.5. On the **Attach permissions policy** page, search and choose **AWSS3FullAccess, AWSGlueServiceRole**, and choose **Next: Review**.

1.6. On the **Review** page, enter the following detail:

> **Role name: AWSGlueServiceRoleDefault**

1.7. Click **Create role**.

1.8. Choose **Roles** page, select the role **AWSGlueServiceDefault** you just created.

1.9. On the **Permissions** tab, choose the link **add inline policy** to create an inline policy.

1.10. On the JSON tab, paste in the following policy:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "logs:CreateLogGroup",
        "logs:CreateLogStream",
        "logs:PutLogEvents",
        "logs:DescribeLogStreams"
      ],
      "Resource": [
        "arn:aws:logs:*:*:*"
      ]
    }
  ]
```

}

1.11. Click **Review policy**.

1.12. On the Review policy, enter policy name: **AWSCloudWatchLogs**.

1.13. Click **Create policy**.

1.14. Now confirm you have policies as below figure.

| Policy name ▼ | Policy type ▼ | |
|---|---|---|
| ⬛ AmazonS3FullAccess | AWS managed policy | ✖ |
| ⬛ AWSGlueConsoleFullAccess | AWS managed policy | ✖ |
| AWSCloudWatchLogs | Inline policy | ✖ |

Figure1: IAM role policies

## 1.2 Add Crawler

1.15. On the **Services** menu, click **AWS Glue**.

1.16. In the console, choose **Add database**. In the **Database name**, type **nycitytaxi**, and choose **Create**.

1.17. Choose **Crawlers** in the navigation pane, choose **Add crawler**. Add type Crawler name **nytaxicrawler**, and choose **Next**.

1.18. On the **Add a data store** page, choose **S3** as data store.

1.19. Select **Specified path in my account**.

1.20. Enter data store path **s3://aws-bigdata-blog/artifacts/glue-data-lake/data/**, and choose **Next**.

1.21. On **Add another data store** page, choose **No**, and choose **Next**.

1.22. Select Choose an existing IAM role, and choose the role **AWSGlueServiceRoleDefault** you just created in the drop-down list, and choose **Next**.

1.23. For **Frequency**, choose **Run on demand**, and choose **Next**.

1.24. For **Database**, choose **nycitytaxi**, and choose **Next**.

1.25. Review the steps, and choose **Finish**.

1.26. The crawler is ready to run. Choose **Run it now**.

1.27. When the crawler has finished, one table has been added. Choose **Tables** in the left navigation pane, and then choose **data** to confirmed.
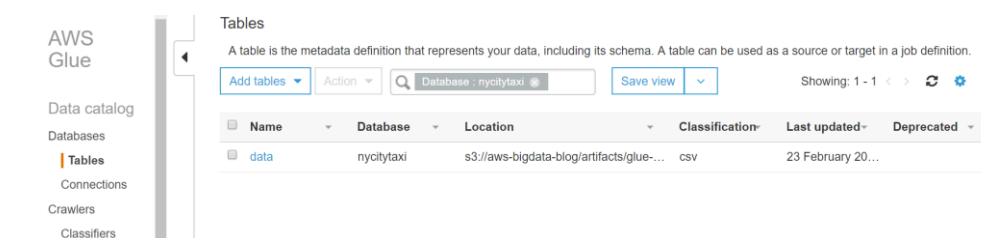
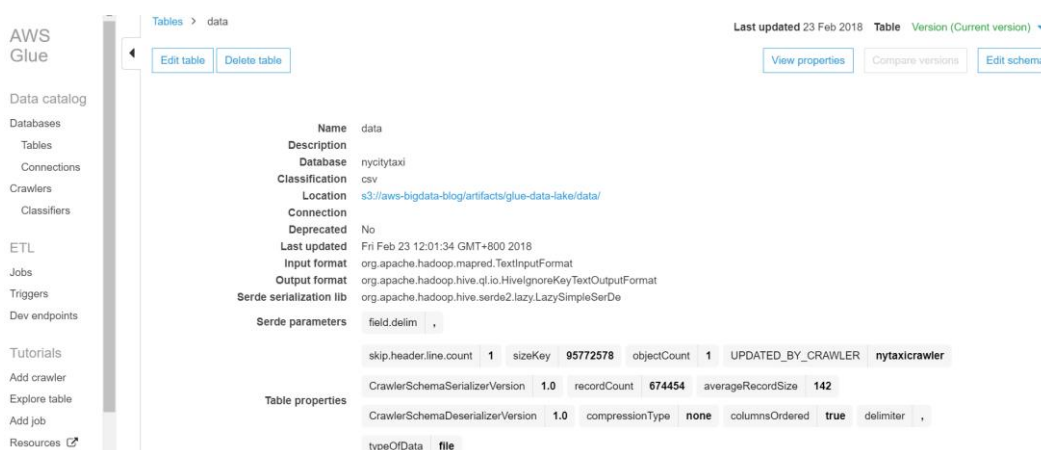Figure2: AWS Glue table has been added



Figure3: table information

## 1.3 Transform the Data from CSV to Parquet Format

1.28. In the navigation pane, under **ETL**, choose **Jobs**, and then choose **Add job**.

1.29. On the Job properties, enter the following details:

**Name: nytaxi-csv-parquet**

**IAM role:** choose **AWSGlueServiceRoleDefault**

1.30. For **This job runs**, select **A proposed script generated by AWS Glue**.

1.31. Choose **Next**.

1.32. Choose **data** as the data source, and choose **Next**.

1.33. Choose **Create tables in your data target**.

1.34. For Data store, choose Amazon S3, and choose **Parque**t as the format.

1.35. For **Target path**, select S3 bucket with **s3://aws-glue-result-xxxx** to store the results.

1.36. Verify the schema mapping, and choose **Finish**.

1.37. View the job. This screen provides a complete view of the job and allows you to edit, click **Save**, and choose **Run job**. This steps may be waiting around 10
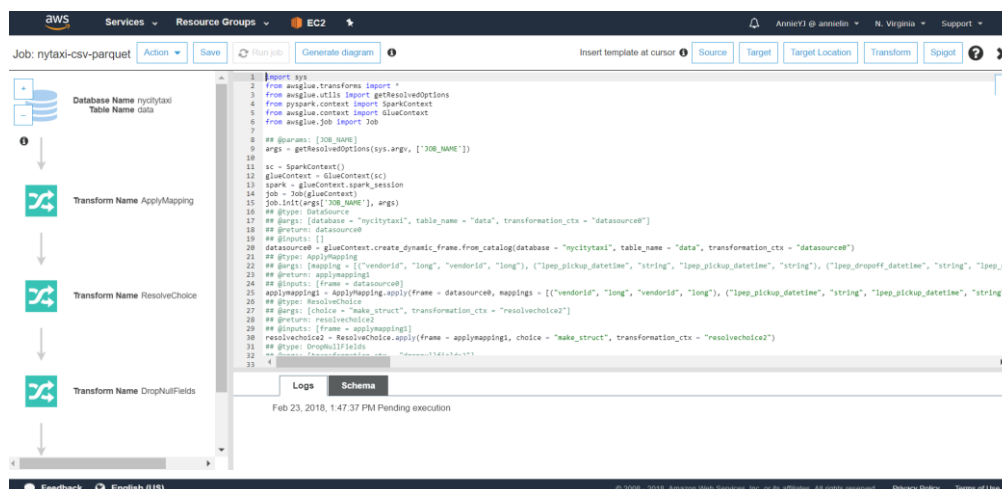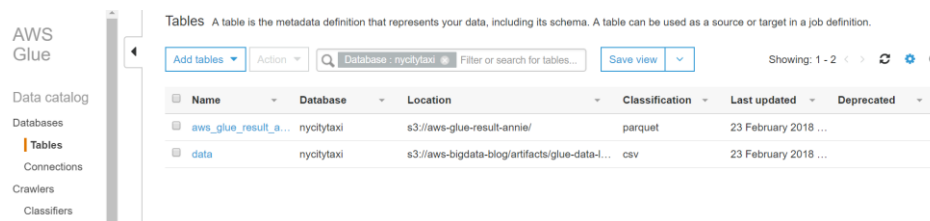
minutes.



Figure4: job runing

## 1.4 Add the Parquet Table and Crawler

1.38. When the job has finished, add a new table for the Parquet data using a crawler.

1.39. In the navigation pane, choose **Add crawler**. Add type Crawler name **nytaxiparquet**, and choose **Next**.

1.40. Choose S3 as the **Data store**.

1.41. Include path choose **s3://aws-glue-result-xxxx** to store data.

1.42. Choose **Next.**

1.43. On **Add another data store** page, choose **No**, and choose **Next**.

1.44. Select **Choose an existing IAM role**, and choose the role **AWSGlueServiceRoleDefault** you just created in the drop-down list, and choose **Next**.

1.45. For **Frequency**, choose **Run on demand**, and choose **Next**.

1.46. For **Database**, choose **nycitytaxi**, and choose **Next**.

1.47. Review the steps, and choose **Finish**.

1.48. The crawler is ready to run. Choose **Run it now**.

1.49. After the crawler has finished, there are two tables in the **nycitytaxi** database: a table for the raw CSV data and a table for the transformed Parquet data.

Figure 5: table for transformed parquet data

## 1.5 Analyze the Data with Amazon Athena

1.50.  On the **Services** menu, click **Athena**.

1.51.  On the **Query Editor** tab, choose the database **nycitytaxi**.
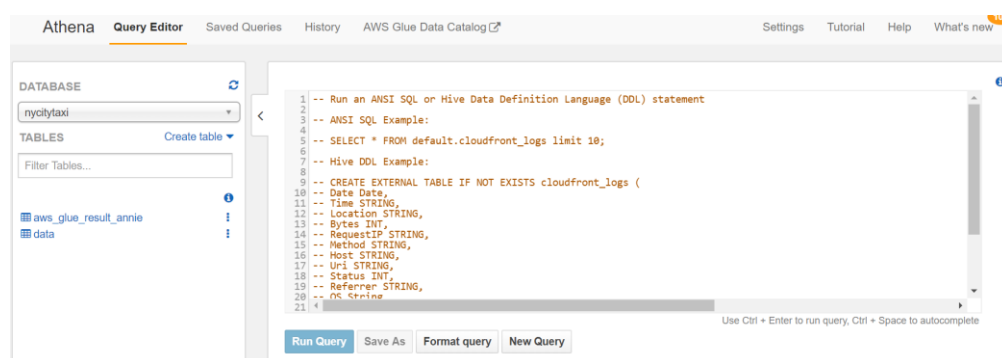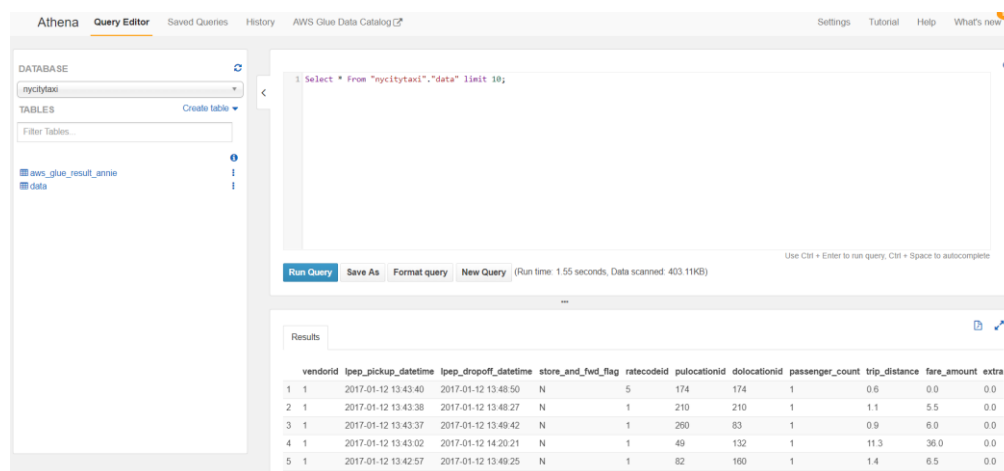


Figure 6: Athena query editor

1.52.  Choose the **aws_glue_result_xxxx** table.

1.53.  Query the data, type below standard SQL:

```
Select * From "nycitytaxi"."data" limit 10;
```

1.54.  Choose **Run Query**.



Figure 7: Athena run query

# Conclusion

Congratulations! You now have learned how to:

- Build data lake using AWS Glue and Amazon S3.

- Crawler your data to Amazon S3 by AWS Glue.

- Analysis through Amazon Athena service.