

We can use DeltaLake file format with AWS Glue by using the Open Source DeltaLake library JAR file and configuring the Glue Job options. The JAR file can be downloaded from the Maven repository and stored in the S3 to use in the Glue job. In addition to configuring both Python lib path and Jar lib path to this S3 location, the Job parameters option in the Glue job needs to be configured as shown below.

Python lib path	s3://aws-glue-pr/jar/delta-core_2.11-0.6.1.jar
Jar lib path	s3://aws-glue-pr/jar/delta-core_2.11-0.6.1.jar
Other lib path	-
	--conf
Job parameters	spark.delta.logStore.class=org.apache.spark.sql.delta.storage.S3SingleDriverLogStore --conf spark.sql.extensions=io.delta.sql.DeltaSparkSessionExtension

We can use below sample code snippet to write data files in DeltaLake format

```

1 from awsglue.context import GlueContext
2 from pyspark.context import SparkContext
3
4 # create/get spark + glue context
5 spark = GlueContext(SparkContext.getOrCreate()).sparkSession
6
7 # folder within S3 for the delta table
8 s3_path = f"s3://aws-glue-pr/users"
9
10 # initially prepopulate the table with some data
11 users_initial = [
12     { 'user_id': 1, 'name': 'Gina Burch', 'gender': 'f' },
13     { 'user_id': 2, 'name': 'Francesco Coates', 'gender': 'm' },
14     { 'user_id': 3, 'name': 'Saeed Wicks', 'gender': 'm' },
15     { 'user_id': 4, 'name': 'Raisa Oconnell', 'gender': 'f' },
16     { 'user_id': 5, 'name': 'Josh Copeland', 'gender': 'm' },
17     { 'user_id': 6, 'name': 'Kaiden Williamson', 'gender': 'm' }
18 ]
19
20 spark.createDataFrame(users_initial) \
21     .write.format("delta").mode("overwrite").save(s3_path)
22
23 # load and print results via Spark API
24 print("DF reading after initial load:")
25 spark.read.format("delta").load(s3_path).orderBy("user_id").show()

```

This would generate standard parquet files with transaction logs in JSON format. Once written the data can be queried through Athena by generating the manifest file from the written DeltaLake data.

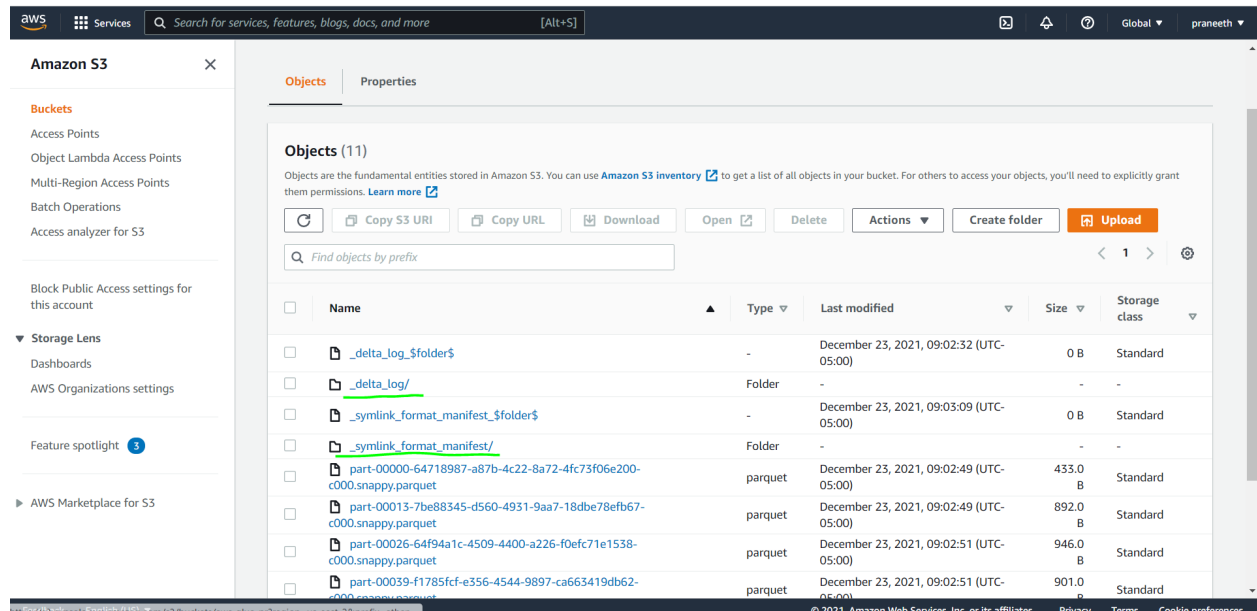
The manifest files can be generated using the below API calls

```
27 deltaTable = DeltaTable.forPath(spark, s3_path)
28 deltaTable.generate("symlink_format_manifest")
```

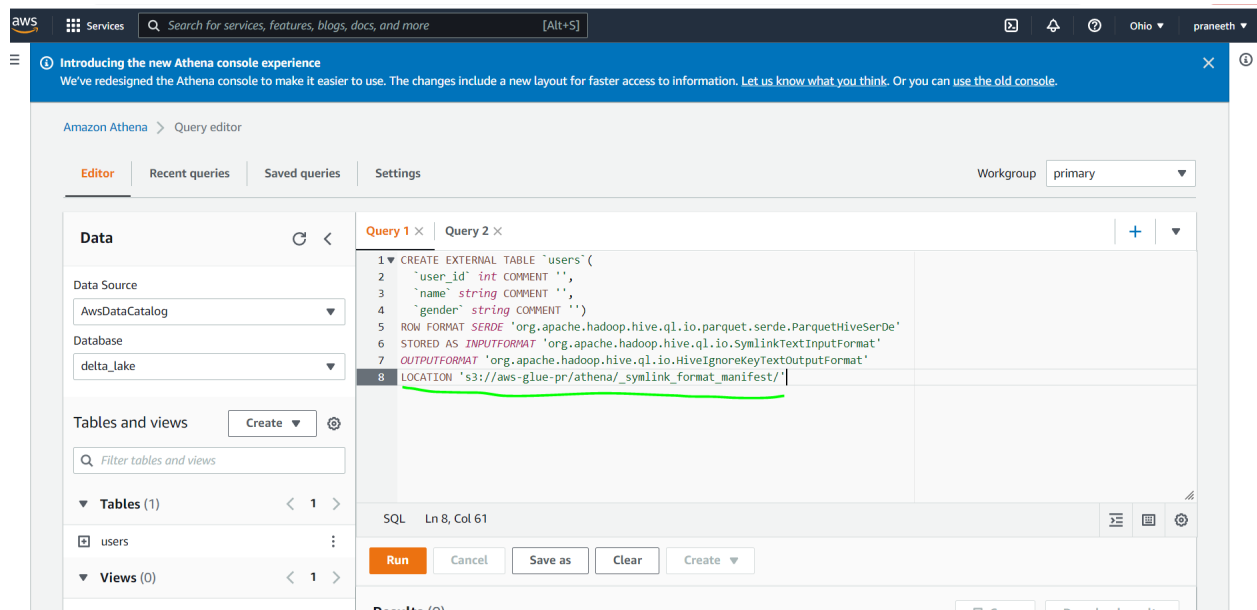
Logs

Schema

After this step the S3 data folder would look like the below picture



Here the **_delta_log** folder contains the DeltaLake transaction log files in JSON format and **_symlink_format_manifest** contains the manifest files. After this we can query the deltalake data in the Athena by creating the external tables pointing to the manifest files location as in the below image



The below image shows the query result of a DeltaLake folder

The screenshot displays the Amazon Athena Query Editor interface. The top navigation bar includes the AWS logo, a search bar, and user information (Ohio, praneeth). The main header shows 'Amazon Athena > Query editor' and a 'Workgroup' dropdown set to 'primary'. The left sidebar contains tabs for 'Editor', 'Recent queries', 'Saved queries', and 'Settings'. The 'Data' section on the left shows the 'Data Source' as 'AwsDataCatalog' and the 'Database' as 'delta_lake'. The 'Tables and views' section lists 'Tables (1)' with 'users' and 'Views (0)'. The main editor area shows 'Query 2' with the SQL statement 'select * from users;'. Below the query, a green status bar indicates 'Completed' with performance metrics: 'Time in queue: 0.165 sec', 'Run time: 0.747 sec', and 'Data scanned: 1.36 KB'. The 'Results (6)' section displays a table with 6 rows and 3 columns: 'user_id', 'name', and 'gender'. The results are as follows:

user_id	name	gender
6	Kaiden Williamson	m
3	Saeed Wicks	m
1	Gina Burch	f
4	Raisa Oconnell	f
5	Josh Copeland	m