

# MySQL to Delta Lake Pipeline Using AWS Glue and PySpark



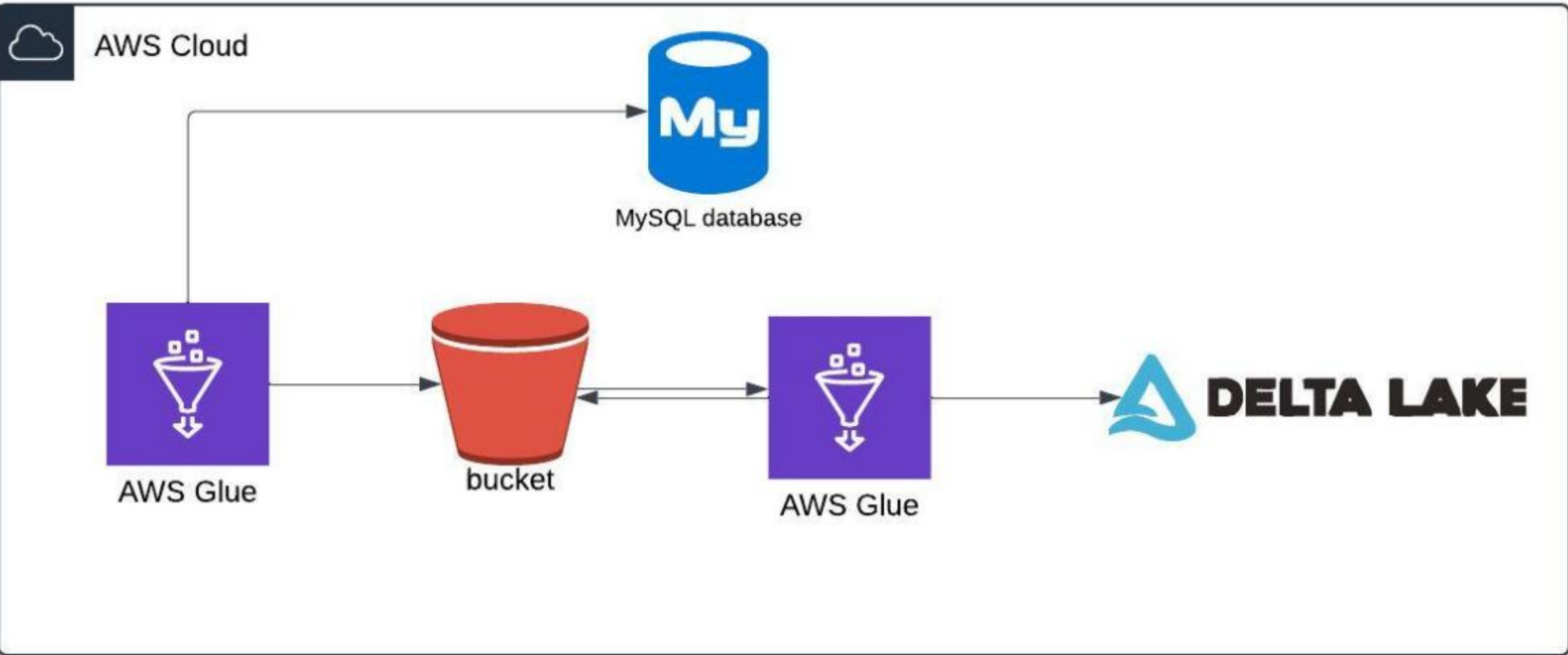
# Objective

- Extract: Use an AWS Glue job with runtime arguments to dynamically extract data from a MySQL database and store it in Amazon S3.
- Transform: Leverage a PySpark-based AWS Glue job to process the extracted data and save it into a Delta Lake format in S3 for optimized storage and querying.



# Tools Used

- MySQL: As the source database for extracting raw data.
- AWS Glue: For performing ETL operations, including data extraction and transformation.
- PySpark: For writing custom transformation logic in the Glue job to process data.
- Amazon S3: As the storage layer for both extracted data and Delta Lake files.
- Delta Lake: For storing processed data in a format optimized for data lakes and analytics.
- AWS Glue Runtime Arguments: For passing dynamic parameters like table names, partitions, or filters during the Glue job execution.
- Python: For scripting transformation logic in PySpark and configuring Glue jobs.
- CloudWatch: For monitoring Glue job execution and troubleshooting issues.





```
-> departure_airport VARCHAR(10),
-> arrival_airport VARCHAR(10),
-> scheduled_departure_time DATETIME,
-> actual_departure_time DATETIME,
-> scheduled_arrival_time DATETIME,
-> actual_arrival_time DATETIME,
-> gate VARCHAR(10),
-> terminal VARCHAR(5),
-> created_at DATETIME
-> );
```

Query OK, 0 rows affected (0.26 sec)

```
mysql> LOAD DATA LOCAL INFILE "flights-data/flight_logistics.csv" INTO TABLE flights_db.flight_logistics FIELDS TERMINATED BY ',' ENCLOSED BY
'"'"' LINES TERMINATED BY '\n' IGNORE 1 ROWS;
```

Query OK, 24966 rows affected, 25029 warnings (4.57 sec)

Records: 25000 Deleted: 0 Skipped: 34 Warnings: 25029

```
mysql> LOAD DATA LOCAL INFILE "flights-data/passenger_experience.csv" INTO TABLE flights_db.passenger_experience FIELDS TERMINATED BY ',' ENCL
OSED BY '"'"' LINES TERMINATED BY '\n' IGNORE 1 ROWS;
```

Query OK, 24966 rows affected, 34 warnings (2.87 sec)

Records: 25000 Deleted: 0 Skipped: 34 Warnings: 34

```
mysql> LOAD DATA LOCAL INFILE "flights-data/flights.csv" INTO TABLE flights_db.flights FIELDS TERMINATED BY ',' ENCLOSED BY '"'"' LINES TERMINAT
ED BY '\n' IGNORE 1 ROWS;
```

Query OK, 24966 rows affected, 34 warnings (4.60 sec)

Records: 25000 Deleted: 0 Skipped: 34 Warnings: 34

```
mysql> SHOW tables;
```

```
+-----+
| Tables_in_flights_db |
+-----+
| flight_logistics      |
| flights               |
| passenger_experience   |
+-----+
```

3 rows in set (6.91 sec)

```
mysql> █
```



```
1 import boto3, sys, csv, pymysql, io, json, logging
2 from datetime import date
3 from botocore.exceptions import ClientError
4 from awsglue.utils import getResolvedOptions
5
6 logging.basicConfig(level=logging.INFO)
7 logger = logging.getLogger()
8
9 args = getResolvedOptions(sys.argv, ["table_name", "delta_value"])
10 table_name = args["table_name"]
11 delta_value = args["delta_value"]
12
13 def get_rds_credentials(secret_name, region_name):
14     session = boto3.session.Session()
15     client = session.client(service_name='secretsmanager', region_name=region_name)
16     try:
17         get_secret_value_response = client.get_secret_value(SecretId=secret_name)
18     except ClientError as e:
19         logger.error(f"Unable to retrieve secret: {e}")
20         return None
21
22     if 'SecretString' not in get_secret_value_response:
23         logger.error("Secret does not contain a string.")
24         return None
25
26     secret = get_secret_value_response['SecretString']
27     credentials = json.loads(secret)
28     return credentials
29
30 credentials = get_rds_credentials("flights_db", "us-east-1")
31
32 user_name = credentials['username']
```

```

30 credentials = get_rds_credentials("flights_db", "us-east-1")
31
32 user_name = credentials['username']
33 host = credentials['host']
34 password = credentials['password']
35 db_name = "flights_db"
36
37 s3_bucket = "gd-aws-de-labs"
38 s3_key = f'raw_landing_zone/{db_name}/{table_name}/data.csv'
39
40 def main():
41     try:
42         connection = pymysql.connect(
43             host=host,
44             user=user_name,
45             password=password,
46             database=db_name,
47             cursorclass=pymysql.cursors.DictCursor
48         )
49
50         if table_name=='flights':
51             sql = f"SELECT * FROM {table_name} where date(created_at)='{delta_value}'"
52
53         else:
54             sql = f"SELECT * FROM {table_name}"
55
56         with connection.cursor() as cursor:
57
58             cursor.execute(sql)
59             result = cursor.fetchall()
60
61     except:
62         print("Error: {}".format(sys.exc_info()[0]))
63         return
64
65     csv_data = convert_to_csv(result)
66     upload_csv(csv_data, s3_bucket, s3_key)
67
68 if __name__ == '__main__':
69     main()

```







```
40 def main():
68     logger.error(f'Error: {str(e)}')
69     sys.exit(1)
70
71     finally:
72         connection.close()
73
74 def convert_to_csv(data):
75     if not data:
76         return ""
77     csv_file = io.StringIO()
78     fieldnames = data[0].keys()
79     writer = csv.DictWriter(csv_file, fieldnames=fieldnames)
80     writer.writeheader()
81     for row in data:
82         writer.writerow(row)
83     return csv_file.getvalue()
84
85 if __name__ == '__main__':
86     main()
87
```



## my-sql-extraction-job

Last modified on 06/01/2025, 14:41:06

Actions ▼

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

### Script [Info](#)

```
1 import boto3, sys, csv, pymysql, io, json, logging
2 from datetime import date
3 from botocore.exceptions import ClientError
4 from awsglue.utils import getResolvedOptions
5
6 logging.basicConfig(level=logging.INFO)
7 logger = logging.getLogger()
8
9 args = getResolvedOptions(sys.argv, ["table_name", "delta_value"])
10 table_name = args["table_name"]
11 delta_value = args["delta_value"]
12
13 ▼ def get_rds_credentials(secret_name, region_name):
14     session = boto3.session.Session()
15     client = session.client(service_name='secretsmanager', region_name=region_name)
```

Python

Ln 1, Col 1



Errors: 0



Warnings: 0







my-sql-extraction-job

Last modified on 06/01/2025, 14:41:06

Actions ▼

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

Job parameters [Info](#)

Key

Value - *optional*

Q --table\_name X

Q flights X

Remove

Q --delta\_value X

Q 2024-05-29 X

Remove

Add new parameter

You can add 48 more parameters.

Tags

No tags associated with the resource.

Add new tag

You can add 50 more tags.





# my-sql-extraction-job

Last modified on 06/01/2025, 14:41:06

Actions ▼

Save

Run

- Script
- Job details
- Runs
- Data quality
- Schedules
- Version Control
- Upgrade analysis - *preview*

## Job parameters [Info](#)

Key	Value - <i>optional</i>	
<input type="text" value="--table_name"/>	<input type="text" value="passenger_experience"/>	<div>Remove</div>
<input type="text" value="--delta_value"/>	<input type="text" value="null"/>	<div>Remove</div>

Add new parameter

You can add 48 more parameters.

## Tags

No tags associated with the resource.

Add new tag

You can add 50 more tags.







# my-sql-extraction-job

Last modified on 06/01/2025, 14:41:06

Actions ▼

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

## Job parameters [Info](#)

Key

Value - *optional*

Q --table\_name X

Q flight\_logistics X

Remove

Q --delta\_value X

Q null X

Remove

Add new parameter

You can add 48 more parameters.

## Tags

No tags associated with the resource.

Add new tag

You can add 50 more tags.



## Amazon S3



### General purpose buckets

Directory buckets

Table buckets

Access Grants

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

### Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

## flights\_db/

Copy S3 URI

Objects

Properties

### Objects (3) Info



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

1

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	flight_logistics/	Folder	-	-	-
<input type="checkbox"/>	flights/	Folder	-	-	-
<input type="checkbox"/>	passenger_experience/	Folder	-	-	-





glue-spark > spark-transactional-delta-lake.py > ...

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, to_date
3 from delta.tables import DeltaTable
4 import logging
5 import boto3
6 import sys
7 from awsglue.utils import getResolvedOptions
8
9 logging.basicConfig(level=logging.INFO)
10 logger = logging.getLogger(__name__)
11
12 args = getResolvedOptions(sys.argv, ["table_name"])
13 table_name = args["table_name"]
14
15 spark = SparkSession.builder \
16     .appName("Delta Lake Upsert Data Aggregations") \
17     .getOrCreate()
18
19 s3_bucket_path = "s3://gd-aws-de-labs/"
20 input_path = f"{s3_bucket_path}/raw_landing_zone/flights_db/{table_name}/"
21 output_path = f"{s3_bucket_path}/lakehouse-dwh/{table_name}"
22
23 primary_key = "flight_id"
24
25 df = spark.read.format("csv").option("header", "true").option("inferSchema", "true").load(input_path)
26
27 if table_name=='flights':
28     df = df.withColumn("flight_date", to_date(col("created_at")))
29
30 logger.info(f"DataFrame schema for table {table_name}: {df.dtypes}")
31
32 delta_table_exists = DeltaTable.isDeltaTable(spark, output_path)
```

```
25 df = spark.read.format("csv").option("header", "true").option("inferSchema", "true").load(input_path)
26
27 if table_name=='flights':
28     df = df.withColumn("flight_date", to_date(col("created_at")))
29
30 logger.info(f"DataFrame schema for table {table_name}: {df.dtypes}")
31
32 delta_table_exists = DeltaTable.isDeltaTable(spark, output_path)
33
34 if delta_table_exists:
35     delta_table = DeltaTable.forPath(spark, output_path)
36
37     logger.info(f"Delta table schema: {delta_table.toDF().dtypes}")
38
39     merge_condition = " AND ".join([f"target.{key} = source.{key}" for key in primary_key])
40     delta_table.alias("target").merge(
41         df.alias("source"),
42         merge_condition
43     ).whenMatchedUpdateAll().whenNotMatchedInsertAll().execute()
44 else:
45     if table_name=='flights':
46         df.write.format("delta").partitionBy("flight_date").mode("overwrite").save(output_path)
47     else:
48         df.write.format("delta").mode("overwrite").save(output_path)
49
50 s3 = boto3.client('s3')
51 bucket_name = s3_bucket_path.split('/')[2]
52 input_prefix = f"raw_landing_zone/{table_name}/"
53 archive_prefix = f"archived/{table_name}/"
54
55 try:
```



```
46         df.write.format("delta").partitionBy("flight_date").mode("overwrite").save(output_path)
47     else:
48         df.write.format("delta").mode("overwrite").save(output_path)
49
50     s3 = boto3.client('s3')
51     bucket_name = s3_bucket_path.split('/')[2]
52     input_prefix = f"raw_landing_zone/{table_name}/"
53     archive_prefix = f"archived/{table_name}/"
54
55     try:
56         objects = s3.list_objects_v2(Bucket=bucket_name, Prefix=input_prefix).get('Contents', [])
57         for obj in objects:
58             source_key = obj['Key']
59             destination_key = source_key.replace(input_prefix, archive_prefix)
60             copy_source = {'Bucket': bucket_name, 'Key': source_key}
61             s3.copy_object(CopySource=copy_source, Bucket=bucket_name, Key=destination_key)
62             s3.delete_object(Bucket=bucket_name, Key=source_key)
63         logger.info("Files moved to archive successfully.")
64     except Exception as e:
65         logger.error(f"Error moving files to archive: {e}")
66
67
68     spark.stop()
69
```



testing

Last modified on 06/01/2025, 14:15:40

Actions ▼

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

Script [Info](#)

```
1 import sys
2 import logging
3 import boto3
4 from awsglue.transforms import *
5 from awsglue.utils import getResolvedOptions
6 from pyspark.context import SparkContext
7 from awsglue.context import GlueContext
8 from awsglue.job import Job
9 from awsgluedq.transforms import EvaluateDataQuality
10 from pyspark.sql.functions import col, to_date
11
12 # Logging setup
13 logging.basicConfig(level=logging.INFO)
14 logger = logging.getLogger(__name__)
15
```

Python

Ln 1, Col 1

✖ Errors: 0

⚠ Warnings: 0







testing

Last modified on 06/01/2025, 14:15:40

Actions ▼

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

## Job parameters [Info](#)

Key

Value - *optional*

🔍 --table\_name



🔍 flights



Remove

Add new parameter

You can add 49 more parameters.

## Tags

No tags associated with the resource.

Add new tag

You can add 50 more tags.





✔ **Successfully started job**  
Successfully started job testing. Navigate to [Run details](#) for more details.

testing

Last modified on 06/01/2025, 15:02:46

Actions ▼

Save

Run

Script

**Job details**

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

## Job parameters [Info](#)

Key

Value - *optional*

Q --table\_name



Q flight\_logistics



Remove

Add new parameter

You can add 49 more parameters.

## Tags

No tags associated with the resource.







testing

Last modified on 06/01/2025, 15:03:12

Actions ▼

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

## Job parameters [Info](#)

Key

Value - *optional*

Q --table\_name



Q passenger\_experience



Remove

Add new parameter

You can add 49 more parameters.

## Tags

No tags associated with the resource.

Add new tag

You can add 50 more tags.





## Amazon S3



### General purpose buckets

Directory buckets

Table buckets

Access Grants

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

### ▼ Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

## Objects (6) [Info](#)



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

< 1 >



<input type="checkbox"/>	Name ▲	Type ▼	Last modified ▼	Size ▼	Storage class ▼
<input type="checkbox"/>	<a href="#">archived/</a>	Folder	-	-	-
<input type="checkbox"/>	<a href="#">glue-scripts/</a>	Folder	-	-	-
<input type="checkbox"/>	<a href="#">lakehouse-dwh/</a>	Folder	-	-	-
<input type="checkbox"/>	<a href="#">raw_landing_zone/</a>	Folder	-	-	-
<input type="checkbox"/>	<a href="#">spotify_data/</a>	Folder	-	-	-
<input type="checkbox"/>	<a href="#">streams/</a>	Folder	-	-	-





## Amazon S3

### General purpose buckets

Directory buckets

Table buckets

Access Grants

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

### Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

## Objects Properties

### Objects (3) Info

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](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	flight_logistics/	Folder	-	-	-
<input type="checkbox"/>	flights/	Folder	-	-	-
<input type="checkbox"/>	passenger_experience/	Folder	-	-	-

## Add data source



### Data source

Choose the source of data to be crawled.

Delta Lake



### Connection - optional

Select a connection to access the data sources below.



Clear selection

Add new connection 

### Include delta lake table paths

Browse for or enter an existing S3 path.

s3://gd-aws-de-labs/lakehouse-dwh/flights/

Remove

Add new delta table path

### Create tables for querying

- ☒ Create Native tables  
Allow integration with query engines that support querying of the Delta transaction log directly.
- ☐ Create Symlink tables  
Create a symlink manifest folder with manifest files partitioned by the partition keys, based on specified configuration parameters.

Cancel

Add a Delta Lake data source





- Data Catalog
  - Databases
  - Tables
  - Stream schema registries
  - Schemas
  - Connections
  - Crawlers
  - Classifiers
  - Catalog settings
- ▼ Data Integration and ETL
  - Zero-ETL integrations [New](#)
  - ETL jobs
    - Visual ETL
    - Notebooks
    - Job run monitoring
  - Interactive Sessions
  - Data classification tools
    - Sensitive data detection
    - Record Matching
  - Triggers
  - Workflows (orchestration)
    - Blueprints
  - Security configurations

✓ **Crawler successfully starting**  
The following crawler is now starting: "testing"

## Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

### Crawlers (1) [Info](#)

Last updated (UTC)  
January 6, 2025 at 09:41:54



Action ▼

Run

Create crawler

View and manage all available crawlers.

< 1 > 

<input type="checkbox"/>	Name ▼	State ▼	Schedule	Last run ▼	Last run t... ▼	Log	Table chan...
<input type="checkbox"/>	<a href="#">testing</a>	 Running	-	-	-	-	-



Zero-ETL integrations [New](#)

▼ **Data Catalog**

Databases

Tables

Stream schema registries

Schemas

Connections

**Crawlers**

Classifiers

Catalog settings

▼ **Data Integration and ETL**

Zero-ETL integrations [New](#)

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Interactive Sessions

Data classification tools

Sensitive data detection

Record Matching

Triggers

Workflows (orchestration)

## Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

### Crawlers (1) [Info](#)

Last updated (UTC)  
January 6, 2025 at 09:51:13



Action ▼

Run

Create crawler

View and manage all available crawlers.

< 1 > 

<input type="checkbox"/>	Name ▼	State ▼	Schedule	Last run ▼	Last run t... ▼	Log	Table chan...
<input type="checkbox"/>	<a href="#">testing</a>	✓ Ready		✓ Succeeded	January 6, 2...	<a href="#">View log</a> 	2 created, ...







Zero-ETL integrations [New](#)

▼ **Data Catalog**

**Databases**

Tables

Stream schema registries

Schemas

Connections

Crawlers

Classifiers

Catalog settings

▼ **Data Integration and ETL**

Zero-ETL integrations [New](#)

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Interactive Sessions

Data classification tools

Sensitive data detection

Record Matching

Triggers

Workflows (orchestration)

## Databases (2)

Last updated (UTC)  
January 6, 2025 at 09:35:49



Edit

Delete

Add database

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

< 1 > 

<input type="checkbox"/>	Name ▲	Description ▼	Location URI ▼	Created on (UTC) ▼
<input type="checkbox"/>	<a href="#">flights_data</a>	-	-	January 6, 2025 at 09:35:47
<input type="checkbox"/>	<a href="#">mobile_network_aggregations</a>	-	-	December 29, 2024 at 07:34:24





Zero-ETL integrations [New](#)

▼ **Data Catalog**

[Databases](#)

Tables

Stream schema registries

Schemas

Connections

Crawlers

Classifiers

Catalog settings

▼ **Data Integration and ETL**

Zero-ETL integrations [New](#)

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Interactive Sessions


Data classification tools

Sensitive data detection

Record Matching

Triggers

Workflows (orchestration)

 **Announcing new optimization features for Apache Iceberg tables**

Optimize storage for Apache Iceberg tables with automatic snapshot retention and orphan file deletion. [Learn more](#)



## flights\_data

Last updated (UTC)  
January 6, 2025 at 09:49:55



Edit

Delete

### Database properties

Name

flights\_data

Description

-

Location

-

Created on (UTC)

January 6, 2025 at 09:35:47

### Tables (3)

Last updated (UTC)  
January 6, 2025 at 09:55:08



Delete

Add tables using crawler

Add table

View and manage all available tables.



Filter tables



1



Name



Database



Location



Classific...



Depreca...



View data

Data quality

Column st



flight\_logistics

flights\_data

s3://gd-aws-de-

Parquet

-

Table data

View data qualit

View statis



flights

flights\_data

s3://gd-aws-de-

delta

-

Table data

View data qualit

View statis



passenger\_expei

flights\_data

s3://gd-aws-de-

Parquet

-

Table data

View data qualit

View statis







Editor

Recent queries

Saved queries

Settings

Workgroup

primary



## Data



## Data source

AwsDataCatalog



## catalogue

None



## Database

flights\_data



## Tables and views

Create



Filter tables and views

## ▼ Tables (3)



1

+ flight\_logistics



+ flights



Query 1 : X

Query 2 : X

Query 3 : X

✖ Query 4 : X

✔ Query 5 : X

( + ) ▼

1 SELECT \* FROM "flights\_data"."flight\_logistics" limit 10;

SQL Ln 1, Col 1



Run again

Explain

Cancel

Clear

Create ▼

☐ Reuse query results  
up to 60 minutes ago 

passenger\_experience

Views (0)

Query results

Query status

Completed

Time in queue: 87 ms

Run time: 467 ms

Data scanned: 479.20 KB

Results (10)

Copy

Download results

Search rows

#	flight_id	fuel_usage	maintenance_check	baggage_errors	safety_incidents
1	FL1000-676	14032.87942444	0	5	0
2	FL1001-278	17636.57421086	0	2	2
3	FL1001-382	12101.33745187	0	5	0
4	FL1003-509	8001.70091412	0	2	1
5	FL1004-893	18434.16753342	0	5	1
6	FL1005-220	13874.24824317	1	4	1
7	FL1005-904	7375.96111243	1	3	0
8	FL1006-188	9460.68282960	1	3	1