

ETL Using Apache Airflow, Aws Glue and Pyspark

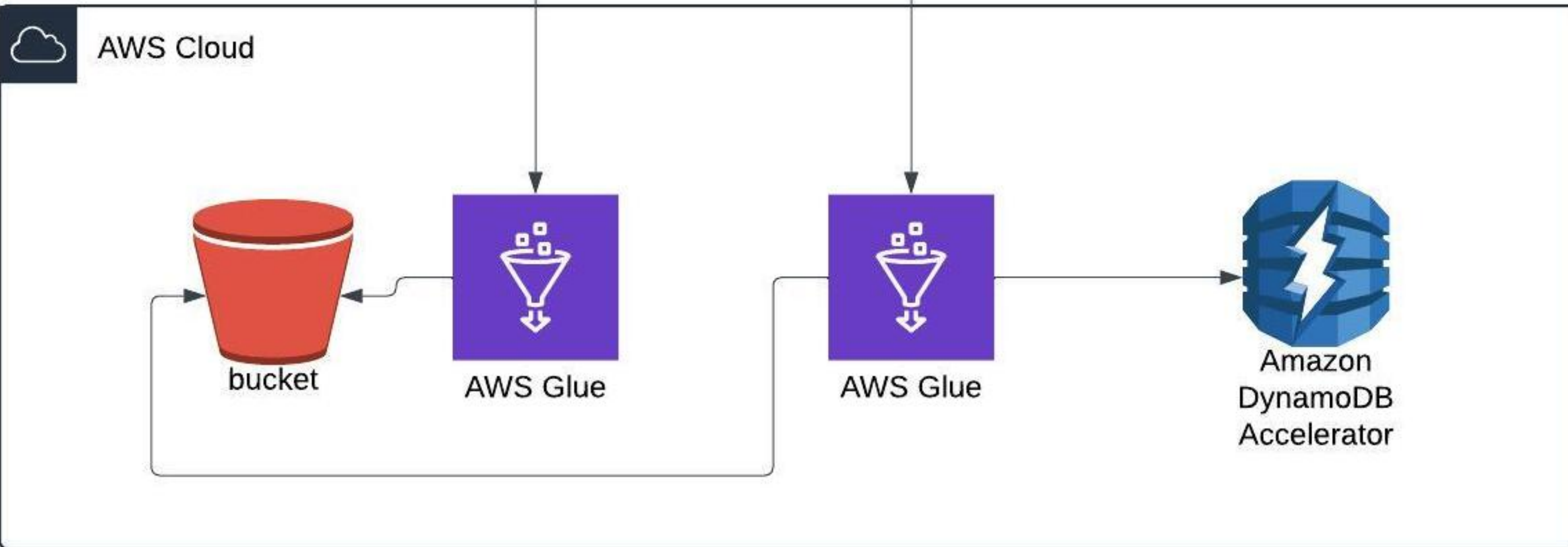
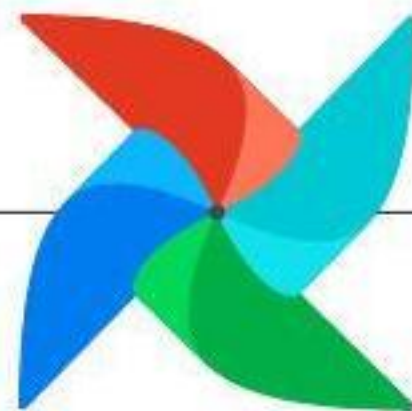


Objective

- Trigger: Use Apache Airflow to orchestrate two AWS Glue jobs (PySpark and Python-based) for processing data stored in Amazon S3.
- Transform: Perform data transformations through the AWS Glue jobs as per business requirements.
- Load: Store the processed data into Amazon DynamoDB for real-time or transactional use cases.
- Organize: Move the processed files within Amazon S3 to a designated location for archiving or further processing.

Tools Used

- Apache Airflow: For workflow orchestration and scheduling of Glue job triggers.
- AWS Glue: For data transformation using PySpark and Python scripts.
- Amazon S3: For storing raw and processed data.
- Amazon DynamoDB: As the target database for storing processed data.
- Python: For scripting and implementing transformation logic within Glue jobs.
- AWS SDK (boto3): For interacting with AWS services programmatically.



DAGs



















All 2Active 0Paused 2Running 0Failed 0

Filter DAGs by tag

Search DAGs

Auto-refresh



 DAG 	Owner 	Runs 	Schedule	Last Run  	Next Run  	Recent Tasks 	Actions	Links
<div><div></div>process-songs-metrics</div>	airflow	<div><div></div><div></div><div></div></div>	<div>* / 5 * * * *<div></div></div>		2025-01-16, 12:05:00 	<div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	<div><div></div></div>
<div><div></div>sample_etl_pipeline</div>	airflow	<div><div></div><div></div><div></div></div>	<div>None<div></div></div>			<div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	<div><div></div></div>





16/01/2025 12:15



All Run Types



All Run States



Clear Filters

Auto-refresh



25

Press **shift** + **/** for Shortcuts

deferred

failed

queued

removed

restarting

running

scheduled

shutdown

skipped

success

up_for_reschedule

up_for_retry

upstream_failed

no_status



DAG

process-songs-metrics

Details

Graph

Gantt

<> Code

Event Log

Run Duration

Task Duration

Calendar

Layout:

Left -> Right



check_files

trigger_spark_job_task

wait_for_spark_job_completion_task

trigger_python_job_task

wait_for_python_job_completion_task

move_files

skip_execution

check_files

skip_execution

trigger_spark_job_task

wait_for_spark_job_completion_task

trigger_python_job_task

wait_for_python_job_completion_task

move_files



dag-glue-workflow.py > ...

```
1  from airflow import DAG
2  from airflow.operators.python_operator import PythonOperator, BranchPythonOperator
3  from airflow.operators.dummy_operator import DummyOperator
4  from datetime import datetime, timedelta
5  import boto3
6  import logging
7  import time
8
9  default_args = {
10     'owner': 'airflow',
11     'depends_on_past': False,
12     'start_date': datetime.now() - timedelta(days=1),
13     'email_on_failure': False,
14     'email_on_retry': False,
15     'retries': 1,
16     'retry_delay': timedelta(minutes=15)
17 }
18
19 dag = DAG('process-songs-metrics',
20           default_args=default_args,
21           description='Trigger Glue job when new files are uploaded to S3 and manage output',
22           schedule_interval='*/5 * * * *',
23           catchup=False)
24
25 def check_files_in_s3(prefix):
26     s3 = boto3.client('s3')
27     bucket_name = 'gd-aws-de-labs'
28
29     response = s3.list_objects_v2(Bucket=bucket_name, Prefix=prefix)
30     contents = response.get('Contents', [])
31     logging.info(f"Contents in {prefix}: {contents}")
32
```



```
22     schedule_interval='*/5 * * * *',
23     catchup=False)
24
25 def check_files_in_s3(prefix):
26     s3 = boto3.client('s3')
27     bucket_name = 'gd-aws-de-labs'
28
29     response = s3.list_objects_v2(Bucket=bucket_name, Prefix=prefix)
30     contents = response.get('Contents', [])
31     logging.info(f"Contents in {prefix}: {contents}")
32
33     for obj in contents:
34         if obj['Size'] > 0:
35             logging.info(f"Non-empty file found in {prefix}: {obj['Key']}")
36             return True
37
38     logging.info(f"No non-empty files found in {prefix}")
39     return False
40
41 def check_all_files(**kwargs):
42     logging.info("Checking for files in S3 prefixes")
43     user_streams = check_files_in_s3('spotify_data/user-streams/')
44     songs = check_files_in_s3('spotify_data/songs/')
45     users = check_files_in_s3('spotify_data/users/')
46
47     logging.info(f"user_streams: {user_streams}, songs: {songs}, users: {users}")
48
49     if user_streams and songs and users:
50         logging.info("All directories have files, proceeding with Spark job")
51         return 'trigger_spark_job_task'
52     else:
53         logging.info("One or more directories are missing files, skipping execution")
```



```
41 def check_all_files(**kwargs):
42     logging.info("Checking for files in S3 prefixes")
43     user_streams = check_files_in_s3('spotify_data/user-streams/')
44     songs = check_files_in_s3('spotify_data/songs/')
45     users = check_files_in_s3('spotify_data/users/')
46
47     logging.info(f"user_streams: {user_streams}, songs: {songs}, users: {users}")
48
49     if user_streams and songs and users:
50         logging.info("All directories have files, proceeding with Spark job")
51         return 'trigger_spark_job_task'
52     else:
53         logging.info("One or more directories are missing files, skipping execution")
54         return 'skip_execution'
55
56 def wait_for_glue_job_completion(job_name, client, poll_interval=60):
57     while True:
58         response = client.get_job_runs(JobName=job_name, MaxResults=1)
59         job_runs = response.get('JobRuns', [])
60
61         if job_runs and job_runs[0]['JobRunState'] in ['RUNNING', 'STARTING', 'STOPPING']:
62             logging.info(f"Glue job {job_name} is still running. Waiting for it to finish...")
63             time.sleep(poll_interval)
64         else:
65             logging.info(f"Glue job {job_name} has finished.")
66             break
67
68 def trigger_glue_job(job_name, **kwargs):
69     client = boto3.client('glue', region_name='us-east-2')
70     logging.info(f"Checking if Glue job {job_name} is running...")
71     wait_for_glue_job_completion(job_name, client)
```


dag-glue-workflow.py > ...

```
56 def wait_for_glue_job_completion(job_name, client, poll_interval=60):
65     logging.info(f"Glue job {job_name} has finished.")
66     break
67
68 def trigger_glue_job(job_name, **kwargs):
69     client = boto3.client('glue', region_name='us-east-2')
70     logging.info(f"Checking if Glue job {job_name} is running...")
71     wait_for_glue_job_completion(job_name, client)
72     logging.info(f"Triggering Glue job: {job_name}")
73     response = client.start_job_run(JobName=job_name)
74
75 def wait_for_spark_job_completion(**kwargs):
76     glue_job_name = 'calculate_metrics_etl'
77     client = boto3.client('glue', region_name='us-east-2')
78     wait_for_glue_job_completion(glue_job_name, client)
79
80 def wait_for_python_job_completion(**kwargs):
81     glue_job_name = 'insert_metrics_dynamo'
82     client = boto3.client('glue', region_name='us-east-1')
83     wait_for_glue_job_completion(glue_job_name, client)
84
85 def move_files_to_archived(**kwargs):
86     s3 = boto3.client('s3')
87     bucket = 'gd-aws-de-labs'
88     source_prefix = 'spotify_data/user-streams/'
89     dest_prefix = 'spotify_data/user-streams-archived/'
90
91     response = s3.list_objects_v2(Bucket=bucket, Prefix=source_prefix)
92     for obj in response.get('Contents', []):
93         source_key = obj['Key']
94         dest_key = source_key.replace(source_prefix, dest_prefix)
```



```

dag-glue-workflow.py > ...
    wait_for_python_job_completion(**kwargs):
83     wait_for_glue_job_completion(glue_job_name, client)
84
85     def move_files_to_archived(**kwargs):
86         s3 = boto3.client('s3')
87         bucket = 'gd-aws-de-labs'
88         source_prefix = 'spotify_data/user-streams/'
89         dest_prefix = 'spotify_data/user-streams-archived/'
90
91         response = s3.list_objects_v2(Bucket=bucket, Prefix=source_prefix)
92         for obj in response.get('Contents', []):
93             source_key = obj['Key']
94             dest_key = source_key.replace(source_prefix, dest_prefix)
95
96             s3.copy_object(Bucket=bucket, CopySource={'Bucket': bucket, 'Key': source_key}, Key=dest_key)
97             s3.delete_object(Bucket=bucket, Key=source_key)
98
99     check_files = BranchPythonOperator(
100         task_id='check_files',
101         python_callable=check_all_files,
102         provide_context=True,
103         dag=dag
104     )
105
106     trigger_spark_job_task = PythonOperator(
107         task_id='trigger_spark_job_task',
108         python_callable=trigger_glue_job,
109         op_args=['calculate_metrics_etl'],
110         provide_context=True,
111         dag=dag
112     )
113

```

```
111 dag = dag
112 )
113
114 wait_for_spark_job_completion_task = PythonOperator(
115     task_id='wait_for_spark_job_completion_task',
116     python_callable=wait_for_spark_job_completion,
117     provide_context=True,
118     dag=dag
119 )
120 trigger_python_job_task = PythonOperator(
121     task_id='trigger_python_job_task',
122     python_callable=trigger_glue_job,
123     op_args=['insert_metrics_dynamo'],
124     provide_context=True,
125     dag=dag
126 )
127
128 wait_for_python_job_completion_task = PythonOperator(
129     task_id='wait_for_python_job_completion_task',
130     python_callable=wait_for_python_job_completion,
131     provide_context=True,
132     dag=dag
133 )
134
135 move_files = PythonOperator(
136     task_id='move_files',
137     python_callable=move_files_to_archived,
138     provide_context=True,
139     dag=dag
140 )
141
142 skip_execution = DummyOperator(
```



```
122     python_callable=trigger_glue_job,
123     op_args=['insert_metrics_dynamo'],
124     provide_context=True,
125     dag=dag
126 )
127
128 wait_for_python_job_completion_task = PythonOperator(
129     task_id='wait_for_python_job_completion_task',
130     python_callable=wait_for_python_job_completion,
131     provide_context=True,
132     dag=dag
133 )
134
135 move_files = PythonOperator(
136     task_id='move_files',
137     python_callable=move_files_to_archived,
138     provide_context=True,
139     dag=dag
140 )
141
142 skip_execution = DummyOperator(
143     task_id='skip_execution',
144     dag=dag
145 )
146
147 # Setup the task dependencies
148 check_files >> [trigger_spark_job_task, skip_execution]
149 trigger_spark_job_task >> wait_for_spark_job_completion_task >> trigger_python_job_task >> wait_for_python_job_completion_
150
```

(import) DummyOperator: Any

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

spotify_data/

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"/>	songs/	Folder	-	-	-
<input type="checkbox"/>	user-streams/	Folder	-	-	-
<input type="checkbox"/>	users/	Folder	-	-	-



AWS Glue



Getting started

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Data Catalog tables

Data connections

Workflows (orchestration)

Zero-ETL integrations [New](#)

▶ **Data Catalog**

▶ **Data Integration and ETL**

▶ **Legacy pages**

What's New

Documentation

AWS Marketplace

☒ Enable compact mode

☒ Enable new navigation

AWS Glue Studio [Info](#)

Create job [Info](#)



Author in a visual interface focused on data flow.

Visual ETL



Author using an interactive code notebook.

Notebook



Author code with a script editor.

Script editor

▶ Example jobs [Info](#)

Create example job

Your jobs (2) [Info](#)



Actions ▼

Run job

🔍 *Filter jobs by property*

< 1 > ⚙

<input type="checkbox"/>	Job name ▼	Type	Created by	Last modified ▼	AWS Glue version ▼
<input type="checkbox"/>	insert_metrics_dynamo	Python shell	Script	16/01/2025, 16:45:59	
<input type="checkbox"/>	calculate_metrics_etl	Glue ETL	Script	16/01/2025, 16:41:18	5.0



```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, to_date, count, sum as _sum, avg, approx_count_distinct, expr, rank, desc
3 from pyspark.sql.window import Window
4
5 spark = SparkSession.builder \
6     .appName("Spotify Advanced KPI Processing") \
7     .getOrCreate()
8
9 # Define S3 bucket and file paths
10 bucket_name = 'nl-aws-de-labs'
11 songs_file_path = f's3a://{bucket_name}/spotify_data/songs/'
12 users_file_path = f's3a://{bucket_name}/spotify_data/users/'
13 user_streams_path = f's3a://{bucket_name}/spotify_data/user-streams/'
14
15 # Read the CSV files into DataFrames
16 songs_df = spark.read.csv(songs_file_path, header=True, inferSchema=True)
17 users_df = spark.read.csv(users_file_path, header=True, inferSchema=True)
18 user_streams_df = spark.read.csv(user_streams_path + "*", header=True, inferSchema=True)
19
20 # Extract the date from listen_time to use as report_date
21 user_streams_df = user_streams_df.withColumn("report_date", to_date(col("listen_time")))
22
23 # Ensure there are no null keys where joins are going to happen
24 user_streams_df = user_streams_df.filter(user_streams_df["track_id"].isNotNull())
25 songs_df = songs_df.filter(songs_df["track_id"].isNotNull())
26
27 # Calculate KPIs for each song on a daily basis
28 song_kpis_df = user_streams_df.groupBy("track_id", "report_date").agg(
29     count("*").alias("total_listens"),
30     approx_count_distinct("user_id").alias("unique_users"),
31     _sum(expr("unix_timestamp(listen_time)").alias("total_listening_time"),
32     avg(expr("unix_timestamp(listen_time)").alias("avg_listening_time_per_user"))
```



```

27 # Calculate KPIs for each song on a daily basis
28 song_kpis_df = user_streams_df.groupBy("track_id", "report_date").agg(
29     count("*").alias("total_listens"),
30     approx_count_distinct("user_id").alias("unique_users"),
31     _sum(expr("unix_timestamp(listen_time)").alias("total_listening_time"),
32     avg(expr("unix_timestamp(listen_time)").alias("avg_listening_time_per_user")
33 )
34
35 # Join with songs data to get genre
36 song_kpis_with_details_df = song_kpis_df.join(songs_df, "track_id")
37
38 # Window specification for ranking songs within each genre by total_listens
39 windowSpec = Window.partitionBy("report_date", "track_genre").orderBy(desc("total_listens"))
40 ranked_songs_df = song_kpis_with_details_df.withColumn("rank", rank().over(windowSpec))
41
42 # Filter for top 3 songs per genre per day
43 top_songs_per_genre = ranked_songs_df.filter(ranked_songs_df.rank <= 3)
44
45 # Find the top 5 genres based on total listens across all songs and all days
46 genre_window = Window.partitionBy("report_date").orderBy(desc("total_listens"))
47 top_genres = top_songs_per_genre.withColumn("genre_rank", rank().over(genre_window)) \
48     .filter(col("genre_rank") <= 5)
49
50 final_df = top_genres.select(
51     col("report_date"),
52     col("track_id"),
53     col("track_name"),
54     col("artists"),
55     col("track_genre"),
56     col("total_listens"),
57     col("unique_users"),
58     col("total_listening_time"),

```

```
43 top_songs_per_genre = ranked_songs_df.filter(ranked_songs_df.rank <= 3)
44
45 # Find the top 5 genres based on total listens across all songs and all days
46 genre_window = Window.partitionBy("report_date").orderBy(desc("total_listens"))
47 top_genres = top_songs_per_genre.withColumn("genre_rank", rank().over(genre_window)) \
48 | .filter(col("genre_rank") <= 5)
49
50 final_df = top_genres.select(
51     col("report_date"),
52     col("track_id"),
53     col("track_name"),
54     col("artists"),
55     col("track_genre"),
56     col("total_listens"),
57     col("unique_users"),
58     col("total_listening_time"),
59     col("avg_listening_time_per_user")
60 )
61
62 # Write the final DataFrame to the same S3 bucket as a CSV file
63 output_path = f's3a://{bucket_name}/spotify_data/output/song_kpis/'
64 final_df.write.mode("overwrite").csv(output_path, header=True)
65 spark.stop()
```



```

1  import boto3
2  import csv
3  from io import StringIO
4  from decimal import Decimal
5  import logging
6
7  # Setting up logging
8  logging.basicConfig(level=logging.INFO, format='%(asctime)s - %(levelname)s - %(message)s')
9
10 # Initialize a boto3 client
11 s3_client = boto3.client('s3')
12 dynamodb = boto3.resource('dynamodb')
13 table = dynamodb.Table('track_level_reports')
14
15 # S3 bucket and file path
16 bucket_name = 'gd-aws-de-labs'
17 file_path = 'spotify_data/output/song_kpis/'
18
19 # Get the latest file from the output folder (variable) bucket_name: Literal['gd-aws-de-labs']
20 response = s3_client.list_objects_v2(Bucket=bucket_name, Prefix=file_path)
21 all_files = response['Contents']
22 latest_file = max(all_files, key=lambda x: x['LastModified'])['Key']
23
24 # Download the latest CSV file content
25 csv_obj = s3_client.get_object(Bucket=bucket_name, Key=latest_file)
26 body = csv_obj['Body'].read().decode('utf-8')
27 csv_data = StringIO(body)
28 csv_reader = csv.DictReader(csv_data)
29
30 # Function to upsert into DynamoDB
31 def upsert_item(row):

```

```

30 # Function to upsert into DynamoDB
31 def upsert_item(row):
32     try:
33         # Attempt to parse and insert/update data in DynamoDB
34         track_id = row['track_id']
35         report_date = row['report_date']
36         total_listens = int(row['total_listens'])
37         unique_users = int(row['unique_users'])
38         total_listening_time = Decimal(str(row['total_listening_time']))
39         avg_listening_time_per_user = Decimal(str(row['avg_listening_time_per_user']))
40
41         table.update_item(
42             Key={
43                 'track_id': track_id,
44                 'report_date': report_date
45             },
46             UpdateExpression='SET total_listens = :tl, unique_users = :uu, total_listening_time = :tlt, avg_listening_time_p
47             ExpressionAttributeValues={
48                 ':tl': total_listens,
49                 ':uu': unique_users,
50                 ':tlt': (variable) avg_listening_time_per_user: Decimal
51                 ':alt': avg_listening_time_per_user
52             },
53             ReturnValues="UPDATED_NEW"
54         )
55         logging.info(f"Successfully processed record for track_id={track_id}, report_date={report_date}")
56     except ValueError as e:
57         logging.error(f"Error processing row {row}: {e}")
58
59 # Process each row in the CSV file
60 for row in csv_reader:
61     upsert_item(row)

```



```

31 def upsert_item(row):
32     table.update_item(
33         Key={
34             'track_id': track_id,
35             'report_date': report_date
36         },
37         UpdateExpression='SET total_listens = :tl, unique_users = :uu, total_listening_time = :tlc, avg_listening_time_per_user = :alt',
38         ExpressionAttributeValues={
39             ':tl': total_listens,
40             ':uu': unique_users,
41             ':tlc': total_listening_time,
42             ':alt': avg_listening_time_per_user
43         },
44         ReturnValues="UPDATED_NEW"
45     )
46     logging.info(f"Successfully processed record for track_id={track_id}, report_date={report_date}")
47 except ValueError as e:
48     logging.error(f"Error processing row {row}: {e}")
49
50 # Process each row in the CSV file
51 for row in csv_reader:
52     upsert_item(row)
53
54 # Close the StringIO object
55 csv_data.close()
56
57
58

```

DynamoDB

Tables

Dashboard

Tables

Explore items

PartiQL editor

Backups

Exports to S3

Imports from S3

Integrations

New

Reserved capacity

Settings

DAX

Clusters

Subnet groups

Parameter groups

Events

The track_level_reports table was created successfully.

Tables (1)

Info

Find tables

Any tag key

Any tag value

<

1

>

Name

Status

Partition key

Sort key

Indexes

Replication Regions

Deletion protection

track_level_reports

Active

track_id (S)

report_date (S)

0

0

Off

CloudShell

Feedback

© 2025, Amazon Web Services, Inc. or its affiliates.

Privacy

Terms

Cookie preferences



DynamoDB



Dashboard

Tables

Explore items

PartiQL editor

Backups

Exports to S3

Imports from S3

Integrations [New](#)

Reserved capacity

Settings

▼ DAX

Clusters

Subnet groups

Parameter groups

Events

Any tag value ▼

Find tables



1



track_level_reports



Protect your DynamoDB table from accidental writes and deletes



When you turn on point-in-time recovery (PITR), DynamoDB backs up your table data automatically so that you can restore to any given second in the preceding 35 days. Additional charges apply. [Learn more](#)

Edit PITR

General information [Info](#)

Partition key

track_id (String)

Capacity mode

On-demand

Alarms

✓ No active alarms

Resource-based policy [Info](#)

⊖ Not active

Sort key

report_date (String)

Table status

✓ Active

Point-in-time recovery (PITR) [Info](#)

⊖ Off

► Additional info



16/01/2025 12:23 

All Run Types 

All Run States 


Clear Filters

 Auto-refresh 

25 

Press **shift** + **/** for Shortcuts

deferredfailedqueuedremovedrestartingrunningscheduledshutdownskippedsuccessup_for_rescheduleup_for_retryupstream_failedno_status



Duration

00:01:06

00:00:33

00:00:00

check_files

trigger_spark_job_task

wait_for_spark_job_completion_task

trigger_python_job_task


wait_for_python_job_completion_task


move_files


skip_execution


DAG


process-songs-metrics


 Details


 Graph


 Gantt

 Code


 Event Log

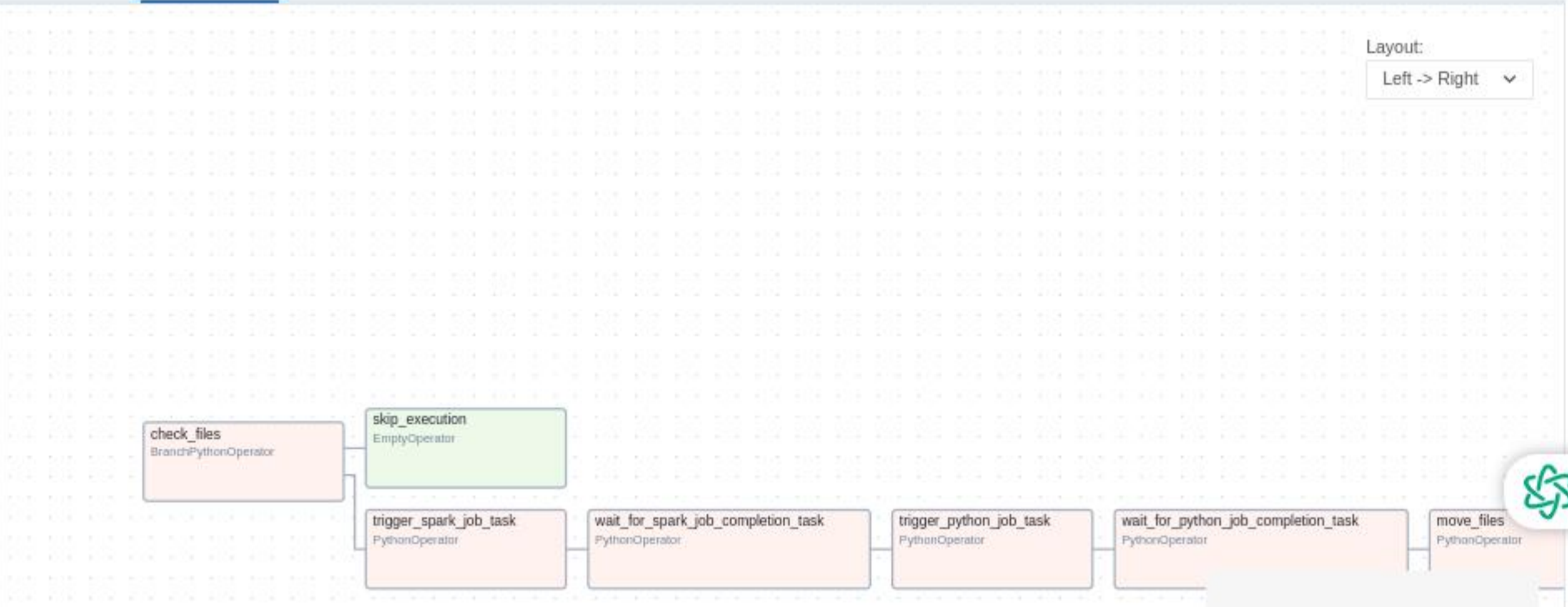
 Run Duration

 Task Duration

 Calendar

Layout:

Left -> Right 



```
graph LR; A[check_files  
BranchPythonOperator] --> B[skip_execution  
EmptyOperator]; A --> C[trigger_spark_job_task  
PythonOperator]; C --> D[wait_for_spark_job_completion_task  
PythonOperator]; D --> E[trigger_python_job_task  
PythonOperator]; E --> F[wait_for_python_job_completion_task  
PythonOperator]; F --> G[move_files  
PythonOperator];
```




AWS Glue

Getting started

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Data Catalog tables

Data connections

Workflows (orchestration)

Zero-ETL integrations [New](#)

► **Data Catalog**

► **Data Integration and ETL**

► **Legacy pages**

What's New [↗](#)

Documentation [↗](#)

AWS Marketplace

☒ Enable compact mode

☒ Enable new navigation

calculate_metrics_etl

Last modified on 16/01/2025, 16:41:18

Actions ▼

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Upgrade analysis - *preview*

Job runs (1/1) [Info](#)

Last updated (UTC)
January 16, 2025 at 12:24:19



View details

Stop job run



Troubleshoot with AI

Table View

Card View



Filter job runs by property



1



Run status ▼

Retries ▼

Start time (Local) ▼

End time (Local) ▼

Duration ▼

Capacit... ▼

Wo



Running

0

01/16/2025 17:53:16

-

54 s

3 DPUs

G.



Run details

Input arguments (10)

Continuous logs

Run insights

Metrics

Stop job run



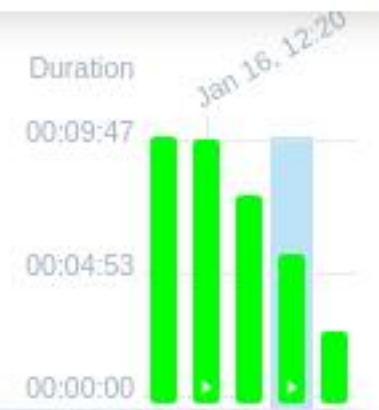
Job name

Start time (Local)

Glue version

Last modified on (Local)



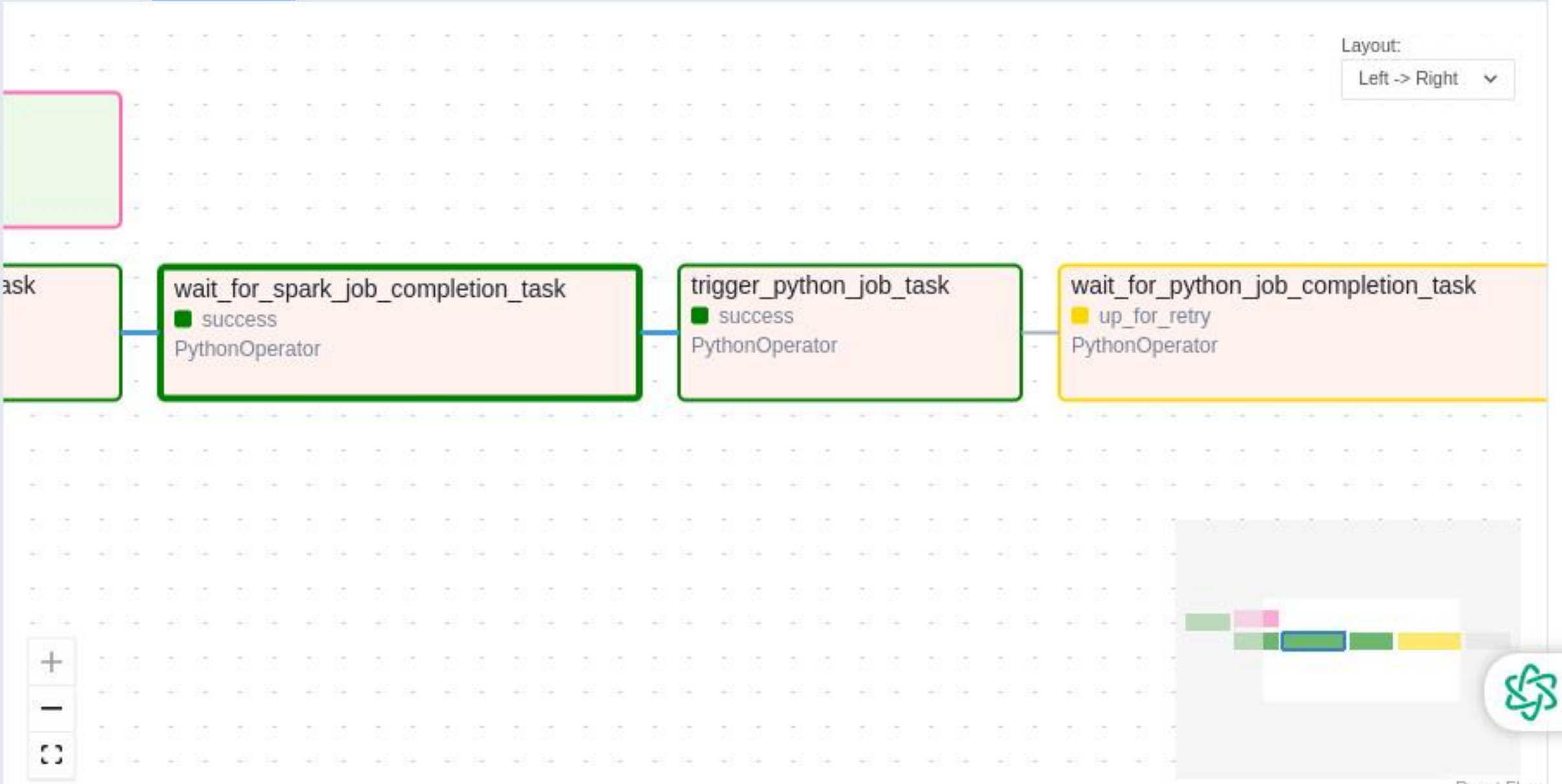


check_files	<div><div></div><div></div><div></div><div></div><div></div></div>
trigger_spark_job_task	<div><div></div><div></div><div></div><div></div><div></div></div>
wait_for_spark_job_completion_task	<div><div></div><div></div><div></div><div></div><div></div></div>
trigger_python_job_task	<div><div></div><div></div><div></div><div></div><div></div></div>
wait_for_python_job_completion_task	<div><div></div><div></div><div></div><div></div><div></div></div>
move_files	<div><div></div><div></div><div></div><div></div><div></div></div>
skip_execution	<div><div></div><div></div><div></div><div></div><div></div></div>

process-songs-metrics / 2023-01-16, 12:27:11 UTC / wait_for_spark_job_completion_task

Clear taskMark state as...Filter DAG by task

DetailsGraphGanttCodeEvent LogLogsXComTask Duration



aws

Search

[Alt+S]

United States (Ohio)

manishtripathi

Amazon S3

Buckets

gd-aws-de-labs

spotify_data/

Copy S3 URI

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

Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

spotify_data/

Objects

Properties

Refresh

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload





Objects (4)

Info

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 >

	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	 output/	Folder	-	-	-
<input type="checkbox"/>	 songs/	Folder	-	-	-
<input type="checkbox"/>	 user-streams/	Folder	-	-	-
<input type="checkbox"/>	 users/	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

Amazon S3

> Buckets

> gd-aws-de-labs

> spotify_data/

> output/

> song_kpis/

Info

Help

Feedback

song_kpis/

Copy S3 URI

Objects

Properties

Objects (2) Info

Refresh

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 > Settings

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	 _SUCCESS	-	January 16, 2025, 18:02:02 (UTC+05:30)	0 B	Standard
<input type="checkbox"/>	 part-00000-0a07deac-256e-4e8f-b33a-55c62205e7cc-c000.csv	csv	January 16, 2025, 18:02:02 (UTC+05:30)	2.4 KB	Standard



Press **shift** + **/** for Shortcuts

deferred failed queued removed restarting running scheduled shutdown skipped success up_for_reschedule up_for_retry upstream_failed no_status



DAG Run
process-songs-metrics2 / ▶ 2025-01-17, 06:34:36 UTC

Clear Mark state as...

Details Graph Gantt Code Event Log

Duration

00:02:34

00:01:17

00:00:00

check_files
trigger_spark_job_task
wait_for_spark_job_completion_task
trigger_python_job_task
wait_for_python_job_completion_task
move_files
skip_execution

Layout:

Left -> Right



+

-





[Alt+S]



United States (Ohio) ▼

manishtripathi ▼



DynamoDB

PartiQL editor



DynamoDB



Dashboard

Tables

Explore items

PartiQL editor

Backups

Exports to S3

Imports from S3

Integrations New

Reserved capacity

Settings

▼ DAX

Clusters

Subnet groups

Parameter groups

Events

PartiQL editor

Operations performed using the PartiQL editor might incur charges. [Learn more](#)

Tables (1)



< 1 > ⚙️

▶ track_level_reports ...

Query 1



Query 2



1 SELECT * FROM track_level_reports;

Run

Clear

Table view

JSON view



CloudShell

Feedback

© 2025, Amazon Web Services, Inc. or its affiliates.

Privacy

Terms

Cookie preferences



DynamoDB



Dashboard

Tables

Explore items

PartiQL editor

Backups

Exports to S3

Imports from S3

Integrations [New](#)

Reserved capacity

Settings

▼ DAX

Clusters

Subnet groups

Parameter groups

Events



Find items



1



avg_listening_time_per_user ▾

report_date ▾

total_listening_time

1719314040

2024-06-25

5157942120

1719324855.6666667

2024-06-25

5157974567

1719325085.3333333

2024-06-25

5157975256

1719311941

2024-06-25

5157935823

1719319292

2024-06-25

5157957876

1719294071

2024-06-25

5157882213

1719337041

2024-06-25

5158011123

1719307529

2024-06-25

5157922587

1719339753

2024-06-25

5158019259

1719298982.25

2024-06-25

6877195929

1719340525.75

2024-06-25

6877362103





DynamoDB

[Dashboard](#)[Tables](#)[Explore items](#)[PartiQL editor](#)[Backups](#)[Exports to S3](#)[Imports from S3](#)[Integrations](#) [New](#)[Reserved capacity](#)[Settings](#)

▼ DAX

[Clusters](#)[Subnet groups](#)[Parameter groups](#)[Events](#)

< 1 >



▼	total_listening_time	▼	unique_users	▼	track_id	▼	total_lis
	5157942120		3		7wcaleaDTt...		3
	5157974567		3		1067ZkQkZ...		3
	5157975256		3		3a2jFwnts4...		3
	5157935823		3		1mGO8rwC...		3
	5157957876		3		5JyEUS5Ew...		3
	5157882213		3		4i3xoVyj8W...		3
	5158011123		3		2u1EtHkbp...		3
	5157922587		3		0pBgoq3dz...		3
	5158019259		3		499SvLeXf...		3
	6877195929		4		4686EWn7t...		4
	6877362103		4		19coiw9dD		4

