

Media Stream Analytics

Datasets & Schema Summary

Ad Revenue

Column	Type	Description
channel_id	int	Unique Channel ID
channel_name	varchar	Name of the channel
date	varchar	Date of revenue
ad_revenue	double	Revenue from ads

Channel Metadata

Column	Type	Description
channel_id	int	Unique Channel ID
channel_name	varchar	Name of the channel
genre	varchar	Genre of the channel
language	varchar	Language of the channel
launch_year	int	Year of launch

Demographics

Column	Type	Description
user_id	int	Unique User ID
gender	varchar	Gender

Column	Type	Description
age_group	varchar	Age group
region	varchar	Region
subscription_type	varchar	Subscription type

Viewership Logs (Real-Time)

Column	Type	Description
user_id	int	Unique User ID
channel_name	varchar	Channel name
channel_id	varchar	Channel ID
timestamp	varchar	Timestamp
duration	int	Duration watched
region	varchar	Region
subscription	varchar	Subscription type
device	varchar	Device used
platform	varchar	Platform used
is_live	boolean	Live stream indicator
genre	varchar	Genre
ads_watched	int	Number of ads watched

Column	Type	Description
ad_revenue	double	Revenue from ads
engagement	int	Engagement score
buffer_count	int	Buffering incidents
completion_pct	int	Completion percentage
session_id	varchar	Session ID
show_name	varchar	Name of the show

Step 1: EC2 Setup for Data Generation

Launch EC2 Instance

- Go to EC2 Console → Launch Instance
- **AMI:** Amazon Linux 2
- **Instance Type:** t2.micro (or higher)
- **Key Pair:** Create or use existing
- **Security Group:** Allow SSH (port 22)
- **SSH ssh -i your-key.pem ec2-user@your-ec2-public-ip**
- Upload and Run Data Generator
- nano generate_media_data.py
- python3 generate_media_data.py

```
import csv
import random
from datetime import datetime, timedelta

# -----
# Global Lists and Mappings
# -----
CHANNEL_NAMES = [
    "Star Plus", "Colors TV", "Sony Entertainment Television", "Zee TV", "Star Bharat",
    "Sony SAB", "Dangal TV", "Zee Anmol", "Colors Rishtey", "DD National",
    "Star Sports 1", "Star Sports 2", "Star Sports 3", "Star Sports Select 1", "Star Sports
Select 2",
    "Sony Sports Ten 1", "Sony Sports Ten 2", "Sony Sports Ten 3", "DD Sports",
    "Eurosport India",
    "Aaj Tak", "ABP News", "India TV", "News18 India", "Republic Bharat",
    "Times Now", "CNN-News18", "NDTV India", "Zee News", "India Today",
    "Sun TV", "Star Vijay", "Colors Tamil", "Zee Tamil", "KTV",
    "Star Maa", "Zee Telugu", "ETV Telugu", "Gemini TV", "Asianet",
    "Surya TV", "Zee Kannada", "Colors Kannada", "Zee Marathi", "Zee Bangla",
    "ETV Marathi", "DD Sahyadri", "DD Bangla", "News18 Kerala", "News18 Tamil Nadu"
]

GENRES = ["Entertainment", "News", "Sports", "Kids", "Music", "Movies"]
LANGUAGES = ["Hindi", "English", "Tamil", "Telugu", "Malayalam", "Kannada",
    "Marathi", "Bengali"]
REGIONS = ["North", "South", "East", "West", "Central", "Northeast", "Pan-India"]
```

```

DEVICES = ["Mobile", "Tablet", "Smart TV", "Laptop"]
PLATFORMS = ["Android", "iOS", "Web", "FireTV", "Roku"]

CHANNEL_ATTRIBUTES = {
    # genre, language pairs
    "Star Sports 1": ("Sports", "Hindi"), "Sony Sports Ten 1": ("Sports", "English"),
    "DD Sports": ("Sports", "Hindi"), "Star Plus": ("Entertainment", "Hindi"),
    "Sony Entertainment Television": ("Entertainment", "Hindi"), "Zee TV":
    ("Entertainment", "Hindi"),
    "Colors TV": ("Entertainment", "Hindi"), "SAB TV": ("Entertainment", "Hindi"),
    "Aaj Tak": ("News", "Hindi"), "NDTV India": ("News", "English"),
    "Republic Bharat": ("News", "Hindi"), "Sun TV": ("Entertainment", "Tamil"),
    "ETV Telugu": ("Entertainment", "Telugu"), "Asianet": ("Entertainment", "Malayalam"),
    "Colors Marathi": ("Entertainment", "Marathi"), "Zee Bangla": ("Entertainment",
    "Bengali"),
    "Times Now": ("News", "English"), "CNN-News18": ("News", "English"),
    "Zee News": ("News", "Hindi"), "ABP News": ("News", "Hindi"),
    "Star Vijay": ("Entertainment", "Tamil"), "KTV": ("Movies", "Tamil"),
}
# -----
# Dataset Generators
# -----


def generate_channel_metadata():
    print("Writing channel meta data.....")
    channel_metadata = []
    channel_id_map = {}
    for i, name in enumerate(CHANNEL_NAMES):
        channel_id = f"CH{i+1:03d}"
        genre, language = CHANNEL_ATTRIBUTES.get(name, (random.choice(GENRES),
random.choice(LANGUAGES)))
        launch_year = random.randint(2000, 2022)
        channel_metadata.append([channel_id, name, genre, language, launch_year])
        channel_id_map[name] = channel_id
    with open("channel_metadata.csv", "w", newline="") as f:
        writer = csv.writer(f)
        writer.writerow(["channel_id", "channel_name", "genre", "language",
"launch_year"])
        writer.writerows(channel_metadata)
    print("Writing channel meta data completed .....")
    return channel_id_map


def generate_demographics(num_users=5000):
    demographics = []
    user_ids = set()
    for _ in range(num_users):

```

```

user_id = f"U{10000 + random.randint(0, 49999)}"
gender = random.choice(["Male", "Female", "Other"])
age_group = random.choice(["<18", "18-25", "26-35", "36-45", "46-60", "60+"])
region = random.choice(REGIONS)
subscription_type = random.choice(["Free", "Basic", "Premium"])
demographics.append([user_id, gender, age_group, region, subscription_type])
user_ids.add(user_id)

with open("demographics.csv", "w", newline="") as f:
    writer = csv.writer(f)
    writer.writerow(["user_id", "gender", "age_group", "region", "subscription_type"])
    writer.writerows(demographics)

return user_ids

def generate_ad_revenue(channel_id_map, days=7):
    ad_revenue = []
    for name, cid in channel_id_map.items():
        for i in range(days):
            date = (datetime.now() - timedelta(days=i)).strftime('%Y-%m-%d')
            revenue = round(random.uniform(10000, 100000), 2)
            ad_revenue.append([cid, name, date, revenue])
    with open("ad_revenue.csv", "w", newline="") as f:
        writer = csv.writer(f)
        writer.writerow(["channel_id", "channel_name", "date", "ad_revenue"])
        writer.writerows(ad_revenue)

def generate_viewership_logs(channel_id_map, user_ids, num_records=10000):
    logs = []
    start_time = datetime.now() - timedelta(days=7)
    user_id_list = list(user_ids)
    for _ in range(num_records):
        user_id = random.choice(user_id_list)
        channel_name = random.choice(CHANNEL_NAMES)
        channel_id = channel_id_map[channel_name]
        timestamp = start_time + timedelta(seconds=random.randint(0, 604800))
        duration = random.randint(1, 180)
        region = random.choice(REGIONS)
        subscription = random.choice(["Free", "Basic", "Premium"])
        device = random.choice(DEVICES)
        platform = random.choice(PLATFORMS)
        is_live = random.choice([True, False])
        genre = CHANNEL_ATTRIBUTES.get(channel_name,
                                       (random.choice(GENRES),))[0]
        ads_watched = random.randint(0, 5)
        ad_revenue = round(ads_watched * random.uniform(10, 200), 2)
        engagement = round(random.uniform(0.2, 1.0), 2)
        buffer_count = random.randint(0, 3)
        completion_pct = round(random.uniform(20, 100), 2)
        logs.append([user_id, channel_id, timestamp, duration, region,
                    subscription, device, platform, is_live, genre,
                    ads_watched, ad_revenue, engagement, buffer_count,
                    completion_pct])
    return logs

```

```

session_id = f"SID{random.randint(100000, 999999)}"
show_name = f>Show_{random.randint(1, 200)}"
logs.append([
    session_id, user_id, channel_id, channel_name, show_name, genre,
    timestamp.strptime('%Y-%m-%d %H:%M:%S'),
    duration, region, subscription, device, platform, is_live, ads_watched,
    ad_revenue,
    engagement, buffer_count, completion_pct
])
with open("viewership_logs.csv", "w", newline="") as f:
    writer = csv.writer(f)
    writer.writerow([
        "session_id", "user_id", "channel_id", "channel_name", "show_name", "genre",
        "timestamp",
        "duration_minutes", "region", "subscription_type", "device", "platform", "is_live",
        "ads_watched", "ad_revenue", "engagement_score", "buffer_count",
        "completion_percentage"
    ])
    writer.writerows(logs)

# -----
# Master Execution
# -----
if __name__ == "__main__":
    channel_id_map = generate_channel_metadata()
    user_ids = generate_demographics()
    generate_ad_revenue(channel_id_map)
    generate_viewership_logs(channel_id_map, user_ids, num_records=10000) #
Adjust as needed

```

- Upload Generated CSVs to S3
- aws s3 cp ad_revenue.csv s3://tbsm-calm/ice/
- aws s3 cp channel_metadata.csv s3://tbsm-calm/ice/
- aws s3 cp demographics.csv s3://tbsm-calm/ice/

ice/

Copy S3 URI

Objects Properties

Objects (3)

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

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	ad_revenue.csv	csv	August 1, 2025, 13:18:38 (UTC+05:30)	13.4 KB	Standard
<input type="checkbox"/>	channel_metadata.csv	csv	August 1, 2025, 13:18:39 (UTC+05:30)	2.0 KB	Standard
<input type="checkbox"/>	demographics.csv	csv	August 1, 2025, 13:18:40 (UTC+05:30)	160.5 KB	Standard

Actions Create folder Upload

```
ec2-user@ip-172-31-10-45:~ x + v
Microsoft Windows [Version 10.0.22621.5624]
(c) Microsoft Corporation. All rights reserved.

C:\Users\yashvardhan.tekavade>cd Downloads

C:\Users\yashvardhan.tekavade\Downloads>ssh -i "tbsm-khora.pem" ec2-user@ec2-54-153-1
aws.com
'_
~\_ #####
~~ \#####
~~ \###|
~~ \#/ ___ https://aws.amazon.com/linux/amazon-linux-2023
~~ V~' '-->
~~ / \
~~ ._. / \
~~ _/ _/
~~ /m/'

Last login: Fri Aug 1 07:45:24 2025 from 210.212.71.1
[ec2-user@ip-172-31-10-45 ~]$ ls
ad_revenue.csv channel_metadata.csv demographics.csv gen.py viewership_logs.csv
```

Step 2: AWS Glue Setup

Create Glue Catalog Database

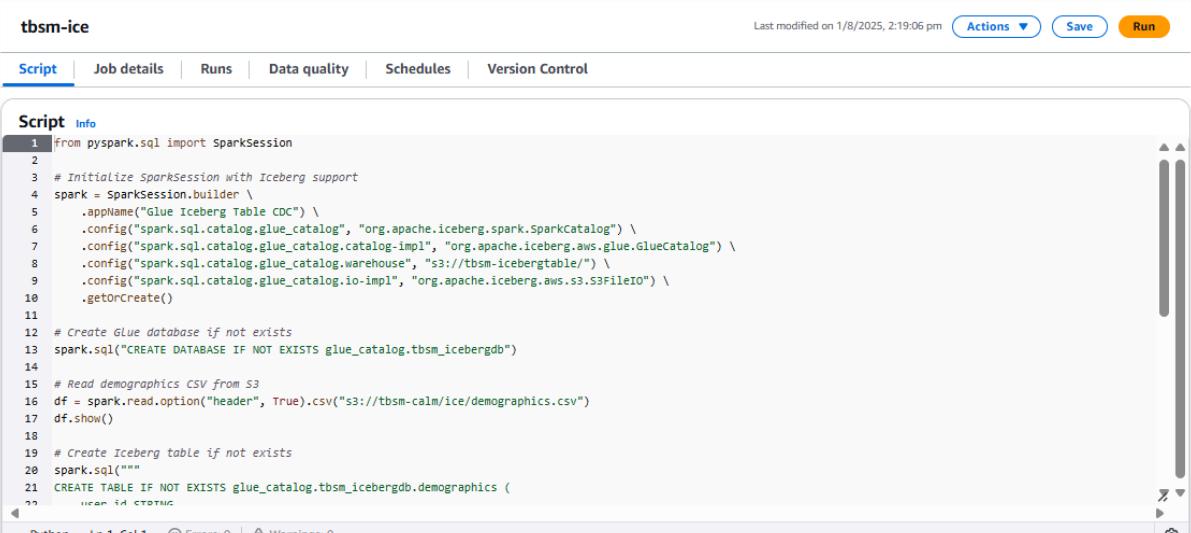
- Go to AWS Glue → Databases → Create
- Name: tbsm_icebergdb

Create Glue Job

- Language: Python
- Glue Version: 4.0 or 5.0
- IAM Role: Must have access to S3 and Glue
- Job Parameters:

Key: --datalake-formats

Value: iceberg



The screenshot shows the AWS Glue Job Editor interface. At the top, it says "tbsm-ice" and "Last modified on 1/8/2025, 2:19:06 pm". Below that are tabs for "Actions", "Save", and "Run". The "Script" tab is selected. The code in the editor is as follows:

```
from pyspark.sql import SparkSession
# Initialize SparkSession with Iceberg support
spark = SparkSession.builder \
    .appName("Glue Iceberg Table CDC") \
    .config("spark.sql.catalog.glue_catalog", "org.apache.iceberg.spark.SparkCatalog") \
    .config("spark.sql.catalog.glue_catalog.catalog-impl", "org.apache.iceberg.aws.glue.GlueCatalog") \
    .config("spark.sql.catalog.glue_catalog.warehouse", "s3://tbsm-icebergtable/") \
    .config("spark.sql.catalog.glue_catalog.io-impl", "org.apache.iceberg.aws.s3.S3FileIO") \
    .getOrCreate()
# Create Glue database if not exists
spark.sql("CREATE DATABASE IF NOT EXISTS glue_catalog.tbsm_icebergdb")
# Read demographics CSV from S3
df = spark.read.option("header", True).csv("s3://tbsm-calm/ice/demographics.csv")
df.show()
# Create Iceberg table if not exists
spark.sql("""
CREATE TABLE IF NOT EXISTS glue_catalog.tbsm_icebergdb.demographics (
    user_id STRING
""")
Python  Ln 1, Col 1  Errors: 0  Warnings: 0
```

Code :-

```
from pyspark.sql import SparkSession

# Initialize SparkSession with Iceberg support
spark = SparkSession.builder \
    .appName("Glue Iceberg Table CDC") \
    .config("spark.sql.catalog.glue_catalog", "org.apache.iceberg.spark.SparkCatalog") \
    .config("spark.sql.catalog.glue_catalog.warehouse", "s3://tbsm-icebergtable/") \
    .config("spark.sql.catalog.glue_catalog.io-impl", "org.apache.iceberg.aws.s3.S3FileIO") \
    .getOrCreate()
```

```
.config("spark.sql.catalog.glue_catalog.catalog-impl",
"org.apache.iceberg.aws.glue.GlueCatalog") \
.config("spark.sql.catalog.glue_catalog.warehouse", "s3://tbsm-icebergtable/") \
.config("spark.sql.catalog.glue_catalog.io-impl",
"org.apache.iceberg.aws.s3.S3FileIO") \
.getOrCreate()

# Create Glue database if not exists
spark.sql("CREATE DATABASE IF NOT EXISTS glue_catalog.tbsm_icebergdb")

# Read demographics CSV from S3
df = spark.read.option("header", True).csv("s3://tbsm-calm/ice/demographics.csv")
df.show()

# Create Iceberg table if not exists
spark.sql("""
CREATE TABLE IF NOT EXISTS glue_catalog.tbsm_icebergdb.demographics (
    user_id STRING,
    gender STRING,
    age_group STRING,
    region STRING,
    subscription_type STRING
)
USING ICEBERG
""")

# Register DataFrame as temporary view
df.createOrReplaceTempView("updates")
```

```

# Perform MERGE (upsert) into Iceberg table

spark.sql("""
MERGE INTO glue_catalog.tbsm_icebergdb.demographics AS target
USING updates AS source
ON target.user_id = source.user_id
WHEN MATCHED THEN UPDATE SET *
WHEN NOT MATCHED THEN INSERT *
""")
```

print("✅ Merge operation completed successfully.")

Run job

Run status	Retries	Start time (Local)	End time (Local)	Duration	Capacity (DPUs)	Worker type	Glue version
Succeeded	0	08/01/2025 14:19:11	08/01/2025 14:20:31	1 m 13 s	10 DPUs	G.1X	5.0
Failed	0	08/01/2025 14:07:51	08/01/2025 14:08:42	42 s	10 DPUs	G.1X	5.0
Failed	0	08/01/2025 14:05:35	08/01/2025 14:06:39	52 s	10 DPUs	G.1X	5.0

Run details	Input arguments (9)	Logs	Run insights	Metrics	Troubleshooting analysis - preview	Spark UI
Job name	Start time (Local)	Glue version	Last modified on (Local)			
tbsm-ice	08/01/2025 14:19:11	5.0	08/01/2025 14:20:31			
Id	End time (Local)	Worker type	Log group name			
jr_4d43596bf6af6a4519a12ff33bc74e6bc7e85fe3	08/01/2025 14:20:31	G.1X	/aws-glue/jobs			
d3cc3ba77a8888593441086e						
Run status	Start-up time	Max capacity	Number of workers			
Succeeded	7 seconds	10 DPUs	10			
Retry attempt number	Execution time	Execution class	Timeout			
Initial run	1 minute 13 seconds	Standard	480 minutes			
Trigger name	Security configuration	Cloudwatch logs	Usage profile			

Table created in iceberg s3

Step 3: S3 Event Notification

Enable Event Notification

- Go to S3 → tbsm-calm → Properties
- Scroll to **Event Notifications**
- Create:
 - Name: trigger-on-upload
 - Event type: PUT
 - Prefix: ice/demographics.csv
 - Destination: Lambda Function

Edit event notification [Info](#)

To enable notifications, you must first add a notification configuration that identifies the events you want Amazon S3 to publish and the destinations where you want Amazon S3 to send the notifications.

General configuration

Event name
trigger-on-upload

Prefix - optional
Limit the notifications to objects with key starting with specified characters.
ice/demographics.csv

Suffix - optional
Limit the notifications to objects with key ending with specified characters.
.jpg

Event types
Specify at least one event for which you want to receive notifications. For each group, you can choose an event type for all events, or you can choose one or more individual events.

Object creation

All object create events
s3:ObjectCreated*

Put
s3:ObjectCreated:Put

Post
s3:ObjectCreated:Post

Copy

Event notifications (1)

Send a notification when specific events occur in your bucket. [Learn more](#)

Name	Event types	Filters	Destination type	Destination
trigger-on-upload	All object create events	ice/demographics.csv	Lambda function	tbsm-trigger-airflow-dag

Amazon EventBridge

For additional capabilities, use Amazon EventBridge to build event-driven applications at scale using S3 event notifications. [Learn more](#) or [see EventBridge pricing](#)

Send notifications to Amazon EventBridge for all events in this bucket

Off

Step 4: Lambda Setup

Create Lambda Function

- Runtime: Python 3.10
- Name: trigger-airflow-dag
- Permissions:
 - MWAAFullAccess
 - AmazonS3ReadOnlyAccess

The screenshot shows the 'Create function' wizard in the AWS Lambda console. The 'Basic information' section is visible, containing fields for Function name (tbsm-trigger-airflow-dag), Runtime (Python 3.10), and Architecture (x86_64). Other tabs like Code, Test, Monitor, Configuration, Aliases, and Versions are present at the top.

Execution role

Role name: tbsm-trigger-airflow-dag-role-bh7xtm0k

Resource summary

To view the resources and actions that your function has permission to access, choose a service.

Amazon CloudWatch Logs

By action | By resource

Resource: arn:aws:logs:ap-southeast-2:008673239246:*

Actions: Allow: logs:CreateLogGroup

tbsm-trigger-airflow-dag-role-bh7xtm0k [Info](#) [Delete](#)

Summary		Edit
Creation date	August 01, 2025, 14:49 (UTC+05:30)	ARN
Last activity	-	arn:aws:iam::008673239246:role/service-role/tbsm-trigger-airflow-dag-role-bh7xtm0k
		Maximum session duration 1 hour

[Permissions](#) [Trust relationships](#) [Tags](#) [Last Accessed](#) [Revoke sessions](#)

Permissions policies (3) [Info](#) You can attach up to 10 managed policies.

Filter by Type		
Search	All types	
<input type="checkbox"/> Policy name	Type	Attached entities
<input checked="" type="checkbox"/> AmazonS3FullAccess	AWS managed	302
<input type="checkbox"/> AWSLambdaBasicExecutionRole-50f47db3-655d-42...	Customer managed	1
<input type="checkbox"/> MWAAFullAccess	Customer managed	3

[C](#) [Simulate](#) [Remove](#) [Add permissions](#)

```

import json
import boto3
import urllib3
from datetime import datetime

MWAA_ENV_NAME = "tbsm-MyAirflowEnvironment"
DAG_NAME = "trigger_upsert_glue_job"

def lambda_handler(event, context):
    print("📦 Received S3 Event:")
    print(json.dumps(event, indent=2))

    mwaa = boto3.client('mwaa')

    try:
        # 1. Create MWAA CLI token
        resp = mwaa.create_cli_token(Name=MWAA_ENV_NAME)
        web_token = resp['CliToken']
        web_server_hostname = resp['WebServerHostname']
        print("🔒 CLI token and hostname retrieved successfully.")

        # 2. Generate Airflow CLI trigger command
        execution_date = datetime.utcnow().isoformat()
        cli_command = f"args trigger -e {execution_date} {DAG_NAME}"
        trigger_url = f"https://{web_server_hostname}/aws_mwaa/cli"

        print(f"🚀 Triggering DAG: {DAG_NAME} at {execution_date}")

        # 3. Send request using urllib3
        http = urllib3.PoolManager()
        headers = {
            "Authorization": f"Bearer {web_token}",

```

```
"Content-Type": "text/plain"
}

response = http.request(
    "POST",
    trigger_url,
    headers=headers,
    body=cli_command.encode("utf-8"),
    timeout=10.0
)

# 4. Handle response
status = response.status
response_body = response.data.decode("utf-8")

print(f"🌐 HTTP Status: {status}")
print(f"🖨 Response Body:\n{response_body}")

return {
    'statusCode': status,
    'body': json.dumps({
        'message': f"DAG trigger attempted with status {status}",
        'response': response_body
    })
}

except Exception as e:
    print(f"❌ Exception occurred: {str(e)}")
    return {
        'statusCode': 500,
        'body': json.dumps({'error': str(e)})
    }
```

```

lambda_function.py
1 import json
2 import boto3
3 import urllib3
4 from datetime import datetime
5
6 MWAAS_ENV_NAME = "tbsm-MyAirflowEnvironment"
7 DAG_NAME = "trigger_upsert_glue_job"
8
9 def lambda_handler(event, context):
10     print("Received S3 Event:")
11     print(json.dumps(event, indent=2))
12
13     mwaas = boto3.client('mwaas')
14
15     try:
16         # 1. Create MWAAS CLI token
17         resp = mwaas.create_cli_token(Name=MWAAS_ENV_NAME)
18         web_token = resp['cliToken']
19         web_server_hostname = resp['WebServerHostname']
20         print(f"CLI token and hostname retrieved successfully.")
21
22         # 2. Generate Airflow CLI trigger command
23         execution_date = datetime.utcnow().isoformat()
24         cli_command = f"dag trigger -e {execution_date} {DAG_NAME}"
25         trigger_url = f"https:///{web_server_hostname}/aws_mwaas/cli"
26
27         print(f"Triggering DAG: {DAG_NAME} at {execution_date}")
28
29         # 3. Send request using urllib3
30         http = urllib3.PoolManager()
31         headers = {
32             "Authorization": f"Bearer {web_token}",
33             "Content-Type": "text/plain"
34         }
35

```

Create new test event

Event Name
test1

Maximum of 25 characters consisting of letters, numbers, dots, hyphens and underscores.

Event sharing settings
 Private
 This event is only available in the Lambda Console and to the event creator. You can configure a total of ten. [Learn more](#)

Shareable
 This event is available to IAM users within the same account who have permissions to access and use shareable events.
[Learn more](#)

Template - optional
S3 Put

Event JSON

```

1 {
2     "Records": [
3         {
4             "eventVersion": "2.0",
5             "eventSource": "aws:s3",
6             "awsRegion": "us-east-1",
7             "eventTime": "1970-01-01T00:00:00.000Z",
8             "eventName": "ObjectCreated:Put",
9             "userIdentity": {
10                 "principalId": "EXAMPLE"
11             },
12             "requestParameters": {
13                 "sourceIPAddress": "127.0.0.1"
14             },
15             "responseElements": {
16                 "x-amz-request-id": "EXAMPLE123456789",
17                 "x-amz-id-2": "EXAMPLE123/5678abcdefghijklambdaisawesome/mnc"
18             },
19             "s3": {

```

```

 7  def lambda_handler(event, context):
12      web_token = response['CliToken']
13      mwaa_webserver = response['WebServerHostname']
14
15      # Build the Airflow CLI URL to trigger the DAG
16      trigger_url = f"https://{{mwaa_webserver}}/aws_mwaa/cli"
17
18      # The airflow command to trigger the DAG (normally sent to the MWAA CLI endpoint)
19      cli_command = "airflow dags trigger {{DAG_NAME}}"
20
21      # For demonstration, returning the URL and token,
22      # but actual triggering would need an authenticated HTTP POST with this token

```

PROBLEMS OUTPUT CODE REFERENCE LOG TERMINAL Execution Results

Status: Succeeded
Test Event Name: test1

Response:

```

{
  "statusCode": 200,
  "body": "{\"message\": \"Trigger command for DAG run_glue_demographics_etl ready\", \"cli_token\": \"eyJ0exAiOiJKViQiLCJhbGciOiJSUzIiNj9...\"}

```

Step 5: Apache Airflow (MWAA) Setup

MWAA Environment

- Create MWAA environment in VPC
- Attach S3 bucket for DAGs

Upload DAG to /dags/ folder in MWAA S3 bucket

Amazon MWAA > Environments > tbsm-MyAirflowEnvironment > Edit

tbsm-MyAirflowEnvironment is being created. This takes 20-30 minutes.

Specify details

Environment details

Airflow version: 2.10.3 (Latest)

Weekly maintenance window start (UTC): Sunday 00:30

DAG code in Amazon S3

S3 Bucket: s3://tbsm-encore

DAGs folder: s3://tbsm-encore/dag

Plugins file - optional

Need to create new MWAA vpc around it

CloudFormation > Stacks

Stacks (7)

Stack name	Status	Created time	Description
MWAA-VPC-Airflow2-TK	CREATE_IN_PROGRESS	2025-08-01 15:18:45 UTC+0530	This template deploys a VPC, with a pair of public and private subnets spread across two Availability Zones. It deploys an internet gateway, with a default route on the public subnets. It deploys a pair of NAT gateways (one in each AZ), and default routes for them in the private subnets.
sm-MWAA-VPC	CREATE_COMPLETE	2025-08-01 15:03:47 UTC+0530	This template deploys a VPC, with a pair of public and private subnets spread across two Availability Zones. It deploys an internet gateway, with a default route on the public subnets. It deploys a pair of NAT gateways (one in each AZ), and default routes for them in the private subnets.

sm-MWAA-VPC

Resources (22)

Logical ID	Physical ID	Type	Status	Module
DefaultPrivateRoute1	rtb-02e7dc90d3f2b96bc0.0.0/0	AWS::EC2::Route	CREATE_COMPLETE	-
DefaultPrivateRoute2	rtb-08a1e08a9638fde09j0.0.0/0	AWS::EC2::Route	CREATE_COMPLETE	-
DefaultPublicRoute	rtb-089733ac72b19e1f9j0.0.0/0	AWS::EC2::Route	CREATE_COMPLETE	-
InternetGateway	igw-05f0ef293d995c281	AWS::EC2::InternetGateway	CREATE_COMPLETE	-
InternetGatewayAttachment	IGWVpc-0cf5747025ac8078b	AWS::EC2::VPCGatewayAttachment	CREATE_COMPLETE	-

Select VPC

Configure advanced settings

Networking

[Virtual private cloud \(VPC\)](#) | [Info](#)

Defines the networking infrastructure setup of your Airflow environment. An environment needs 2 private subnets in different availability zones. To create a new VPC with private subnets, choose Create MWAA VPC. [Learn more](#)

[Create MWAA VPC](#)

- [vpc-0cf5747025ac8078b](#)
VPC
- [vpc-01fdada2d208280c6](#)
VPC
- [vpc-025cc9b28b367cd7f](#)
MWAA-VPC-Airflow2-TK
- [vpc-0bc3981d59038a53e](#)
Default

Additional setup required. Your Airflow UI can only be accessed by secure login behind your VPC. Choose this option if your Airflow UI is only accessed within a corporate network and you do not require a public repository for webserver requirements installation. IAM must be used to handle user authentication.

- Public network (Internet accessible)
Your Airflow UI can be accessed by secure login over the Internet. Choose this option if your Airflow UI is accessed outside of a corporate network. IAM must be used to handle user authentication.

For private network access, the Airflow web server is reached via a VPC endpoint inside your VPC. Connecting to the endpoint requires additional setup. [Learn more about VPC endpoints](#)

[Security group\(s\)](#) | [Info](#)

A VPC security group is required to allow traffic between your environment and your web server.

[Create new security group](#)

Allow MWAA to create a VPC security group with inbound and outbound rules based on your selection for web server access.

[Existing security group\(s\)](#)

You can choose 1 or more existing security groups to configure the inbound and outbound rules for your environment.

- [sg-00eff3a3165d63ffd](#)
Security Group for Amazon MWAA Environment tbsm-MyAirflowEnvironment airflow-security-group-tbsm-MyAirflowEnvironment-4n13c5

Max 5 security groups

AIRLFOW DAG CODE

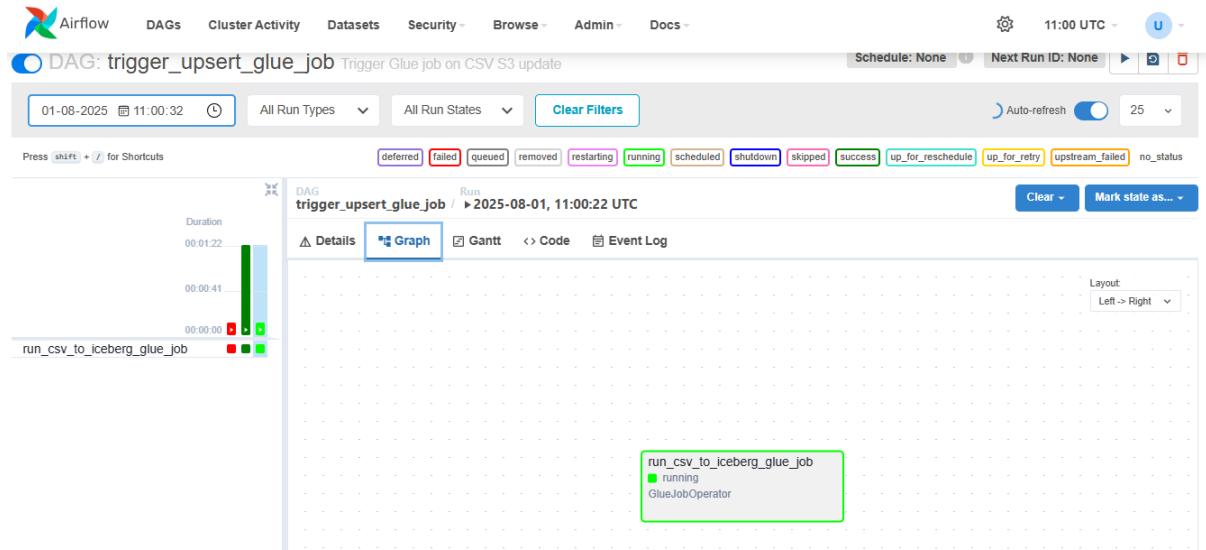
```
GNU nano 8.3      run_glue_demographics_etl.py          from airflow import DAG
from airflow.providers.amazon.aws.operators.glue import AwsGlueJobOperator>from
datetime import datetime

with DAG(
    dag_id='run_glue_demographics_etl',
    schedule_interval=None,
    start_date=datetime(2023, 1, 1),
    catchup=False,
    tags=['iceberg', 'glue', 'etl']
) as dag:

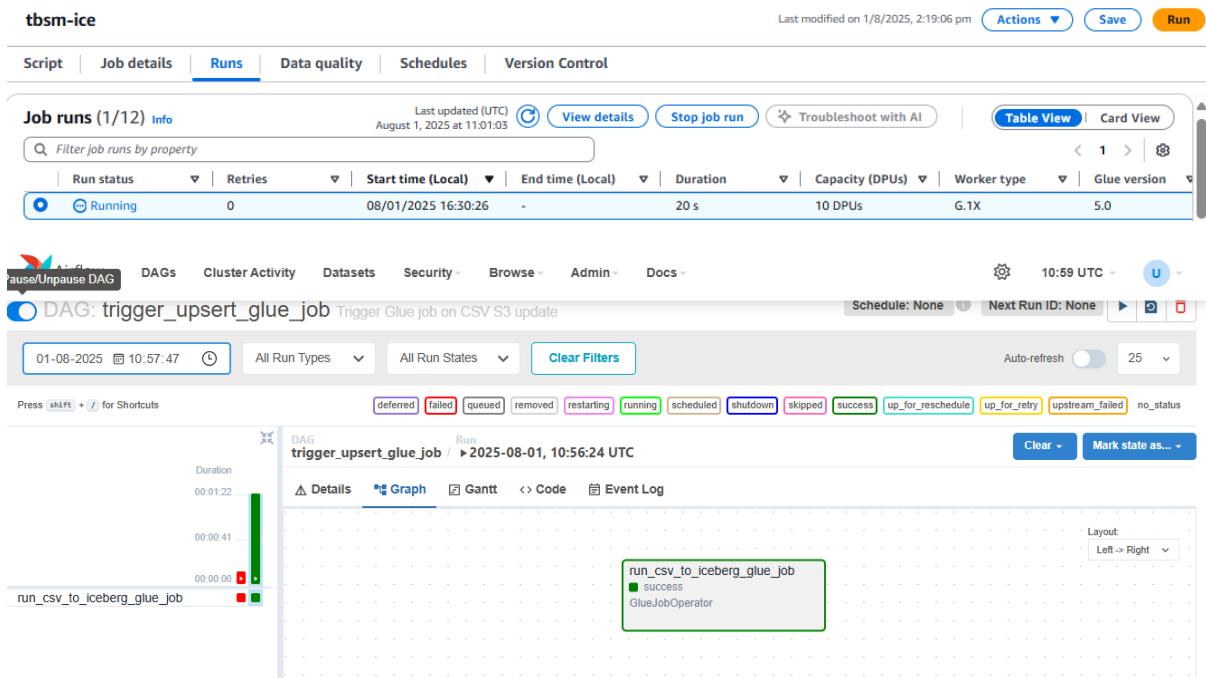
    run_glue_job = AwsGlueJobOperator(
        task_id='run_iceberg_merge_job',
        job_name='tbsm-ice',
        script_location='s3://aws-glue-assets-008673239246-ap-southeast>
aws_conn_id='aws_default',
        region_name='ap-southeast-2'
    )
```

The screenshot shows the AWS MWAH (Amazon Managed Workflows for Apache Airflow) environment details page. The top navigation bar includes the AWS logo, search bar, and user information (Asia Pacific (Sydney) and user28@tac-loud). The main content area displays the environment configuration:

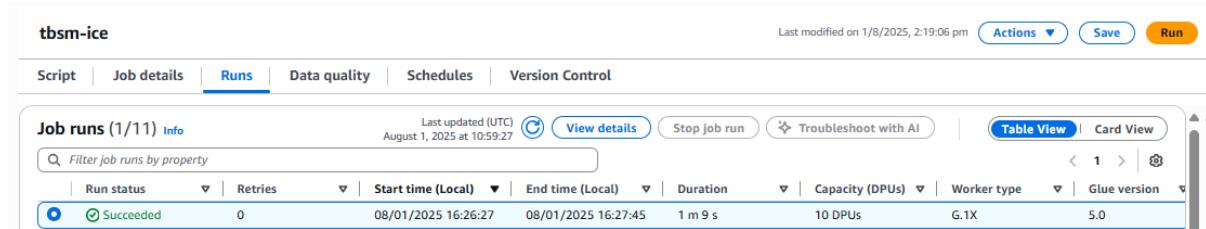
- Details**:
 - Status: Available
 - ARN: arn:aws:airflow:ap-southeast-2:008673239246:environment/tbsm-MyAirflowEnvironment
- Airflow UI**: A link to the Airflow UI at 248a24c3-f0de-42e5-bffb-40160ccb0c02-vpce.c6.ap-southeast-2.airflow.amazonaws.com.
- Weekly maintenance window start (UTC)**: Sunday 00:30
- Last update**:
 - Status: SUCCESS
 - Created at: 2025-08-01T09:42:53.000Z
- Worker replacement strategy**: Forced
- DAG code in Amazon S3**:
 - S3 Bucket: tbsm-encore
 - Plugins file: -
- DAGs folder**: dag
- Requirements file**: -



After running airflow it should trigger glue



After airflow successfully you can see glue job completed



Done with airflow and lambda

Step 6: End-to-End Flow Validation

1. **Login to EC2**
2. Modify or regenerate demographics.csv
3. Upload to S3: s3://tbsm-calm/ice/demographics.csv
4. S3 triggers **Lambda**
5. Lambda triggers **Airflow DAG**
6. Airflow runs **Glue Job**
7. Glue performs **MERGE into Iceberg table**
8. Iceberg table stored in: s3://tbsm-icebergtable/
9. Data available in Athena or Spark SQL

>>START HERE

-open ec2 login
-make changes to csv file
-old file

```
U32047,Male,<18,East,Premium
U30279,Other,46-60,Central,Pre
U36000,Female,46-60,North,Free
U33235,Male,36-45,West,Free
U13573,Male,60+,North,Premium
```

-new file

```
U30279,Male,46-60,West,Premium
U32047,Male,<18,East,Premium
U30279,Other,46-60,Central,Pre
U36000,Female,46-60,North,Free
U33235,Male,36-45,West,Free
U13573,Male,60+,North,Premium
U42900,Other,18-25,West,Basic
```

Upload file

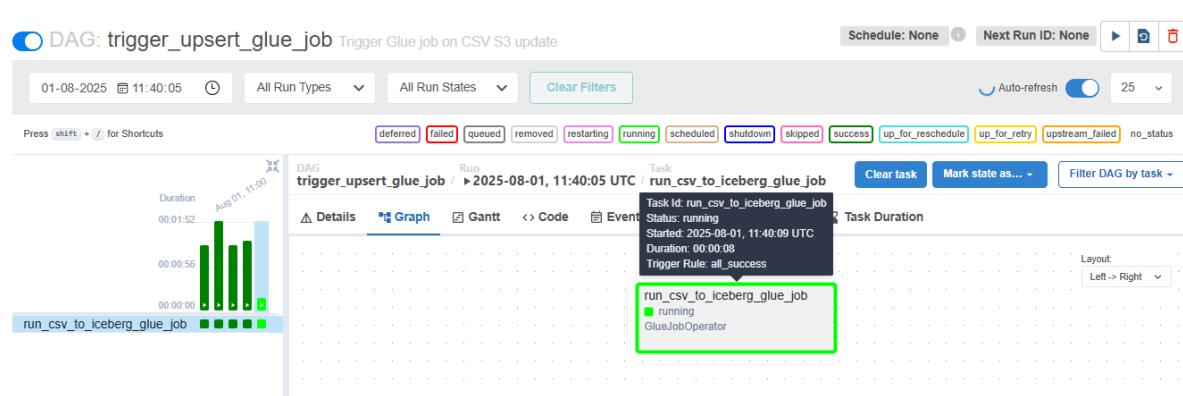
```
[ec2-user@ip-172-31-10-45 ~]$ aws s3 cp demographics.csv s3://tbsm-calm/ice/
upload: ./demographics.csv to s3://tbsm-calm/ice/demographics.csv
```

It will trigger the lambda

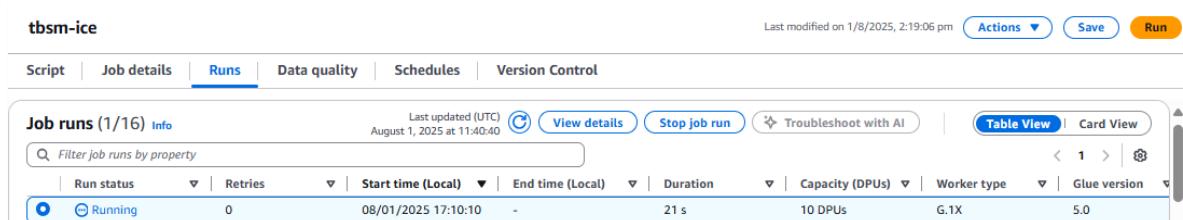
```
④ Executing function: succeeded (logs [2])
▼ Details

{
  "statusCode": 200,
  "body": "{\"message\": \"DAG trigger attempted with status 200\", \"response\": \"{}\\\"stderr\\\"\\n:\\\"\\nVzvz9sb2nhbC9axhb093lysb2nhbC9saWlvc1leG9MydXv592xrX1LxByN8yZthZ2LjPcmzsb3cvWcm1jcy92dGF0c2HfbG9nZzVylN850je4NCBS2v1dnVmKsM8axJmb693M1dhcm5pbmc51fRoZS81vXlpY8tZXrYwMpdnfساWhd69y1idpb0wgymJgZGwcmVjYXX1L2C8pbiB9aduJgznwEx0d11G1G1Ghd9y1HmIbhR1cm4tb1vE9phmcu1CB2b3UgY2uHye5BaOlg1G5wdybie8zBzX9rau51Gnvbmzp2yvChRhp24gbwWcm1jci91c2VfcF0dgyv19tYxRj3C80p8lUcnV1Lgqvdx0LyL2xv2fsl2Zpcmzsb3cvLmxv2fsl2xpy19weXrob241jEx1.3pndgtc0c1j3a2rZxMVNjY2mxdyj9j25naW1cneFaa0uJl850j1058gd0R1cnVXXJuaW6n01Bz2wN8a0u12t1eSbyZ9zCwY9vJb25uSb0Ymgv1b18kZxBy2wNdgKLBC5b3ulg2zvdkxx1hvZvtkYxRhyNfzZ5zckwFwvjaGvev9j25uSb0pkh8Zwfkl1bqDvnb41tc6dyG3ncrfwaiUpd1gabu20yawm1sdyx0o1hr01R6mc0fuad1j7Mwgdmnywf1Bhqd0dwpm00121hvc3lpmx1LgvodMyL2xv1Zfsl2pcntzb1cvLavvY2Z1xpy19heahob241jEx1.3pdu1ts0fja2rDm2Zc2tf3fsNxqadvte3ht021bcSwet0fjigdrXVX3uaW5nG1BuaG1zG1zY2hcncfBax21icnhc21gNoy0kfkes8j3239yNlucry01GexXtzh1hpdggd611fRhnb1g2vchc3lgsirz25Bnb0Qsuecd1Lgebmkg70gkaf1LX3QUTZxtza0uLchmg511f5CB25byX8zYMKZC8p31B9d0lg3hakmLkvv9t1c08yWtz54K\\\"\\n\\\"stdout\\\"\\n:\\\"\\n"
}
```

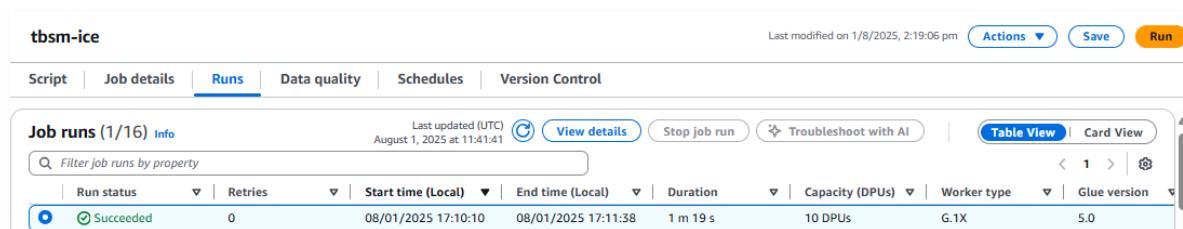
After that airflow will be trigger in running



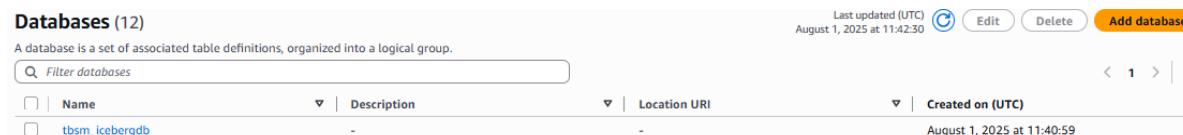
Which will trigger glue in running the job



It will be successfully run



Then you can see table database created



And iceberg s3 is created with folder named

Amazon S3

General purpose buckets

Directory buckets

Table buckets

Vector buckets [Preview](#)

Access Grants

Access Points (General Purpose Buckets, FSx file systems)

Access Points (Directory Buckets)

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

demographics/

Objects (2)

Actions ▾ Create folder Upload

Find objects by prefix

Name Type Last modified Size Storage class

Name	Type	Last modified	Size	Storage class
data/	Folder	-	-	-
metadata/	Folder	-	-	-

And in that parquet

Amazon S3

General purpose buckets

Directory buckets

Table buckets

Vector buckets [Preview](#)

Access Grants

Access Points (General Purpose Buckets, FSx file systems)

Access Points (Directory Buckets)

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

Feature spotlight 11

data/

Objects (5)

Actions ▾ Create folder Upload

Find objects by prefix

Name Type Last modified Size Storage class

Name	Type	Last modified	Size	Storage class
00000-4-08dd959c-19dd-4489-9543-7f71cabd969b-0-00001.parquet	parquet	August 1, 2025, 16:44:48 (UTC+05:30)	21.7 KB	Standard
00000-4-a0fa2e8a-f477-435b-80cf-088096cb775-0-00001.parquet	parquet	August 1, 2025, 16:24:20 (UTC+05:30)	21.7 KB	Standard
00000-4-a5c8d748-564f-434a-a6f2-5de231f2cc2-0-00001.parquet	parquet	August 1, 2025, 16:49:30 (UTC+05:30)	21.7 KB	Standard
00000-4-c19254d1-b259-475a-84be-5227d9ae1a0c-0-00001.parquet	parquet	August 1, 2025, 16:27:33 (UTC+05:30)	21.7 KB	Standard
00000-4-ce81db3a-2063-4846-bcf4-d338c96aaece-0-00001.parquet	parquet	August 1, 2025, 17:11:26 (UTC+05:30)	21.7 KB	Standard

Real-Time Streaming Setup: Viewership Logs

Data Flow:

viewership_logs.csv → Python Streamer → Kinesis Data Stream

- EMR Spark Streaming Job (PySpark)
 - Writes to S3 (raw + processed)
 - Writes to Snowflake (transformed)

Step 1: Kinesis Setup

Create Kinesis Stream

IAM Role for EC2

Attach a role with the following policies:

- AmazonKinesisFullAccess
 - AmazonS3FullAccess
 - CloudWatchLogsFullAccess (optional for logging)
-

- **2. Step 2: EC2 Setup for Streaming**
 - Launch EC2 Instance
 - AMI: Amazon Linux 2
 - Security Group: Allow SSH (port 22)
 - IAM Role: Attach the role created above
- **SSH into EC2**
 - ssh -i your-key.pem ec2-user@your-ec2-public-ip
 - Install Dependencies
 - sudo yum update -y
 - sudo yum install python3 -y
 - pip3 install boto3
- Upload and Run Streamer Script
 - nano generate_viewership_logs.py
 - python3 generate

This script should:

Generate JSON records

Send them to viewership-stream using boto3

```
import json
import boto3
import time

# Initialize Kinesis client
client = boto3.client('kinesis', region_name="ap-southeast-2")
stream_name = "viewership-stream"

# CSV Header (Must match order in the CSV file)
columns = [
    "session_id", "user_id", "channel_id", "channel_name", "show_name", "genre",
    "timestamp",
    "duration_minutes", "region", "subscription_type", "device", "platform",
    "is_live",
    "ads_watched", "ad_revenue", "engagement_score", "buffer_count",
    "completion_percentage"
]

# Read and process CSV file
with open("viewership_logs.csv") as f:
    lines = f.readlines()[1:] # Skip header

    for line in lines:
        # Convert line into dictionary
        values = line.strip().split(',')
        record = dict(zip(columns, values))

        # Print what is being streamed
        print("👉 Streaming Record →", json.dumps(record, indent=2))

        # Send to Kinesis
        response = client.put_record(
            StreamName=stream_name,
            Data=json.dumps(record),
            PartitionKey=record["user_id"]
        )

        # Optional: Print Kinesis response (Record ID)
        print("✅ Sent to Kinesis. SequenceNumber:", response['SequenceNumber'])
```

```
# Simulate real-time
time.sleep(0.01)
```

```
}

✓ Sent to Kinesis. SequenceNumber: 49665723011096496375117208002099357694618995169282752530
🕒 Streaming Record → {
  "session_id": "SID480667",
  "user_id": "U58068",
  "channel_id": "CH020",
  "channel_name": "Eurosport India",
  "show_name": "Show_60",
  "genre": "Music",
  "timestamp": "2025-07-29 07:31:07",
  "duration_minutes": "93",
  "region": "North",
  "subscription_type": "Basic",
  "device": "Smart TV",
  "platform": "iOS",
  "is_live": "False",
  "ads_watched": "4",
  "ad_revenue": "718.82",
  "engagement_score": "0.87",
  "buffer_count": "1",
  "completion_percentage": "90.01"
}
✓ Sent to Kinesis. SequenceNumber: 49665723011118797120315738625249355893628945935011676194
```

viewership-stream [Info](#) [Delete](#)

Data stream summary

Status Active	Capacity mode On-demand	ARN arn:aws:kinesis:ap-southeast-2:008673239246:stream/viewership-stream	Creation time August 01, 2025 at 17:25 GMT+5:30
Data retention period 1 day			

[Applications](#) [Monitoring](#) [Configuration](#) [Enhanced fan-out \(0\)](#) [Data viewer](#) [Data analytics - new](#) [Data stream sharing](#) [EventBridge Pipes](#)

Producers [Info](#)

Producers put records into Kinesis Data Streams.

Amazon Kinesis Agent
Use a stand-alone Java software application to send data to the stream. [Learn more](#)

AWS SDK
Use AWS SDK for Java to develop producers. [Learn more](#)

Amazon Kinesis Producer Library (KPL)
Use KPL to develop producers. [Learn more](#)

[View in GitHub](#) [View in GitHub](#) [View in GitHub](#)

Partition key	Data	Approximate arrival timestamp	Sequence number
U29953	{"session_id": "SID572805", "user_id": "U29953", "...	August 01, 2025 at 17:39:57 GMT+5:30	4966572301107419562991867741422460598614...
U56062	{"session_id": "SID983604", "user_id": "U56062", "...	August 01, 2025 at 17:39:57 GMT+5:30	4966572301107419562991867741422944168942...
U40976	{"session_id": "SID234699", "user_id": "U40976", "...	August 01, 2025 at 17:39:58 GMT+5:30	4966572301107419562991867741423306846688...

STEP 3: Snowflake Setup

SQL Setup

```
CREATE OR REPLACE DATABASE TBSM_DB;
```

```
CREATE OR REPLACE SCHEMA TBSM_DB.PUBLIC;
```

```
CREATE OR REPLACE WAREHOUSE COMPUTE_WH
WITH WAREHOUSE_SIZE = XSMALL
AUTO_SUSPEND = 300
AUTO_RESUME = TRUE
INITIALLY_SUSPENDED = TRUE;
```

```
CREATE OR REPLACE TABLE TBSM_DB.PUBLIC.VIEWERSHIP_LOGS (
    session_id STRING,
    user_id STRING,
    channel_id STRING,
    channel_name STRING,
    show_name STRING,
    genre STRING,
    timestamp STRING,
    duration_minutes STRING,
    region STRING,
    subscription_type STRING,
```

```
device STRING,  
platform STRING,  
is_live STRING,  
ads_watched STRING,  
ad_revenue STRING,  
engagement_score STRING,  
buffer_count STRING,  
completion_percentage STRING  
);
```

```
CREATE OR REPLACE USER spark_user  
PASSWORD = 'StrongPassword@123'  
DEFAULT_ROLE = ACCOUNTADMIN  
DEFAULT_WAREHOUSE = COMPUTE_WH  
DEFAULT_NAMESPACE = TBSM_DB.PUBLIC  
MUST_CHANGE_PASSWORD = FALSE;
```

```
GRANT USAGE ON WAREHOUSE COMPUTE_WH TO USER spark_user;  
GRANT USAGE ON DATABASE TBSM_DB TO USER spark_user;  
GRANT USAGE ON SCHEMA TBSM_DB.PUBLIC TO USER spark_user;  
GRANT INSERT, SELECT ON ALL TABLES IN SCHEMA TBSM_DB.PUBLIC TO USER  
spark_user;
```

STEP 4: EMR Cluster Setup

Create EMR Cluster

- Release: EMR 6.x
- Applications: Spark, Hadoop
- EC2 Key Pair: Use existing .pem
- Networking: Public subnet + Auto-assign Public IP
- IAM Role:
 - AmazonKinesisReadOnlyAccess
 - AmazonS3FullAccess

- Snowflake access via Spark connector

STEP 5: S3 Buckets

- Script Location: s3://tbsm-encore/dag/kinesis_to_snowflake_s3.py
- Temp Location: s3://tbsm-encore/dag/temp/

STEP 6: Install Spark Connectors on EMR

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, expr, split
from pyspark.sql.types import *

# 1. Define the schema of the CSV-style records
schema = StructType([
    StructField("session_id", StringType()),
    StructField("user_id", StringType()),
    StructField("channel_id", StringType()),
    StructField("channel_name", StringType()),
    StructField("show_name", StringType()),
    StructField("genre", StringType()),
    StructField("timestamp", StringType()),
    StructField("duration_minutes", IntegerType()),
    StructField("region", StringType()),
    StructField("subscription_type", StringType()),
    StructField("device", StringType()),
    StructField("platform", StringType()),
    StructField("is_live", StringType()),
    StructField("ads_watched", IntegerType()),
    StructField("ad_revenue", FloatType()),
    StructField("engagement_score", FloatType()),
    StructField("buffer_count", IntegerType()),
    StructField("completion_percentage", FloatType())
])

# 2. Initialize Spark Session
print("🚀 Initializing Spark session...")
spark = SparkSession.builder \
    .appName("KinesisToS3AndSnowflake") \
    .getOrCreate()

print("✅ Spark session initialized.")
```

```

# 3. Read from Kinesis
print("⌚ Connecting to Kinesis...")
df_raw = spark.readStream \
    .format("kinesis") \
    .option("streamName", "viewership-stream") \
    .option("endpointUrl", "https://kinesis.ap-southeast-2.amazonaws.com") \
    .option("region", "ap-southeast-2") \
    .option("startingPosition", "LATEST") \
    .load()
print("✅ Connected to Kinesis.")

# 4. Convert and split raw records
df_string = df_raw.withColumn("data_string", expr("CAST(data AS STRING)"))
df_split = df_string.withColumn("fields", split(col("data_string"), ","))

# 5. Assign to structured schema
df_parsed = df_split.select(
    col("fields").getItem(0).alias("session_id"),
    col("fields").getItem(1).alias("user_id"),
    col("fields").getItem(2).alias("channel_id"),
    col("fields").getItem(3).alias("channel_name"),
    col("fields").getItem(4).alias("show_name"),
    col("fields").getItem(5).alias("genre"),
    col("fields").getItem(6).alias("timestamp"),
    col("fields").getItem(7).cast("int").alias("duration_minutes"),
    col("fields").getItem(8).alias("region"),
    col("fields").getItem(9).alias("subscription_type"),
    col("fields").getItem(10).alias("device"),
    col("fields").getItem(11).alias("platform"),
    col("fields").getItem(12).alias("is_live"),
    col("fields").getItem(13).cast("int").alias("ads_watched"),
    col("fields").getItem(14).cast("float").alias("ad_revenue"),
    col("fields").getItem(15).cast("float").alias("engagement_score"),
    col("fields").getItem(16).cast("int").alias("buffer_count"),
    col("fields").getItem(17).cast("float").alias("completion_percentage")
)
)

# 6. Function to write batch to S3 and Snowflake
def write_batch(batch_df, epoch_id):
    print(f"⌚ Processing batch {epoch_id}...")

    record_count = batch_df.count()
    print(f"📦 Records in batch: {record_count}")

    if record_count == 0:
        print("⚠️ Empty batch. Skipping writes.")
    return

```

```

try:
    # Snowflake options
    snowflake_options = {
        "sfURL": "vdhiuox-lib93136.snowflakecomputing.com",
        "sfUser": "YASHTEKAVADE",
        "sfPassword": "Yashvardhan@420",
        "sfDatabase": "TBSM_DB",
        "sfSchema": "PUBLIC",
        "sfWarehouse": "COMPUTE_WH",
        "sfRole": "ACCOUNTADMIN"
    }

    print("👉 Writing to Snowflake...")
    batch_df.write \
        .format("net.snowflake.spark.snowflake") \
        .options(**snowflake_options) \
        .option("dbtable", "VIEWERSHIP_LOGS") \
        .mode("append") \
        .save()
    print("✅ Write to Snowflake succeeded.")

except Exception as e:
    print("❌ Error writing to Snowflake:", e)

try:
    print("👉 Writing to S3...")
    batch_df.write \
        .format("csv") \
        .mode("append") \
        .save("s3://tbsm-encore/view/")
    print("✅ Write to S3 succeeded.")

except Exception as e:
    print("❌ Error writing to S3:", e)

# 7. Start the streaming query
print("📡 Starting streaming query...")
query = df_parsed.writeStream \
    .foreachBatch(write_batch) \
    .outputMode("append") \
    .option("checkpointLocation", "s3://tbsm-encore/dag/") \
    .trigger(processingTime="10 seconds") \
    .start()

print("✅ Streaming started.")

```