# ETL PIPELINE IN AWS CLOUD
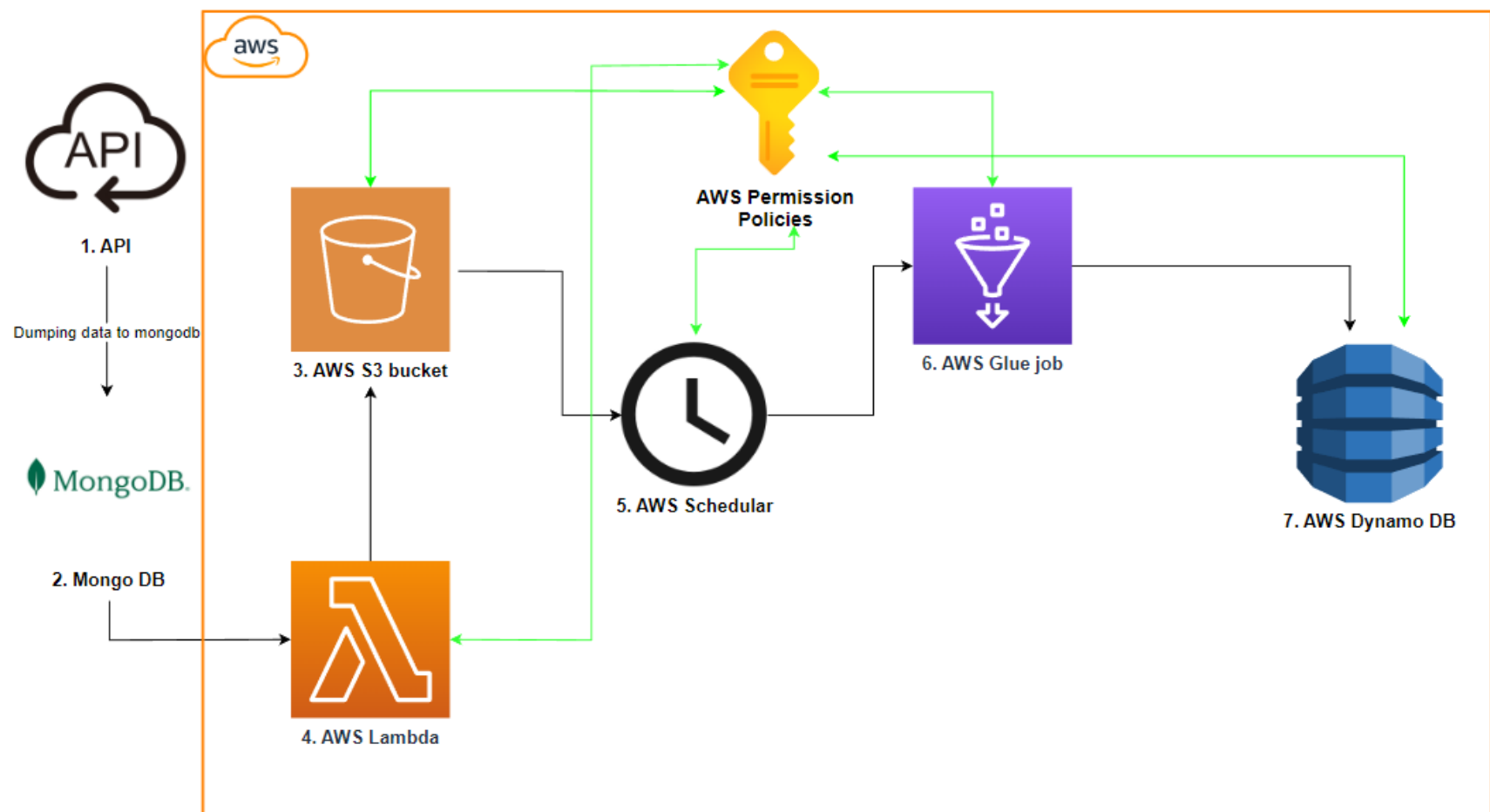
## Introduction

The Aim of the project is to perform ETL Pipeline in AWS using different services of AWS.

## Architecture



## Explanation:

First, we will be extract data from the API (https://www.consumerfinance.gov/data-research/consumer-complaints/search/api/v1/) From some specific interval date and dumping data in MongoDB.

Before dumping data to s3 we will be attaching policies roles for lambda,s3, event bridge scheduler, glue, dynamo DB
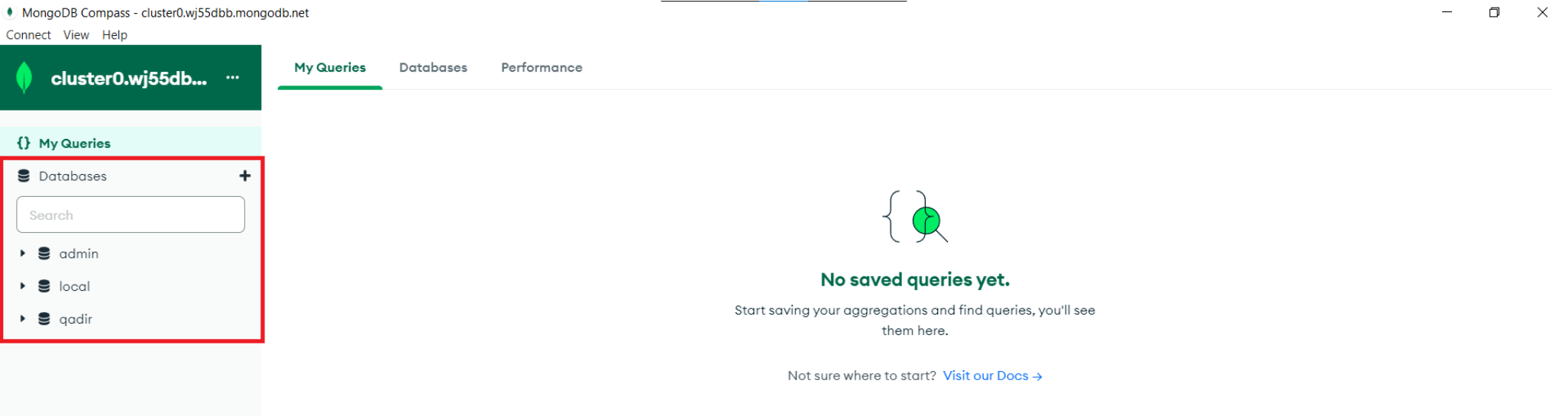
Next from MongoDB dumping data to s3, with help of AWS lambda using python request while dumping data into s3 configured AWS event bridge to dump data for every 2 min with the help of scheduler.

Now, we need to create Dynamo DB without creating any kind of schema, with the help of glue job writing script dynamically we will be dumping data from s3 to dynamo db.
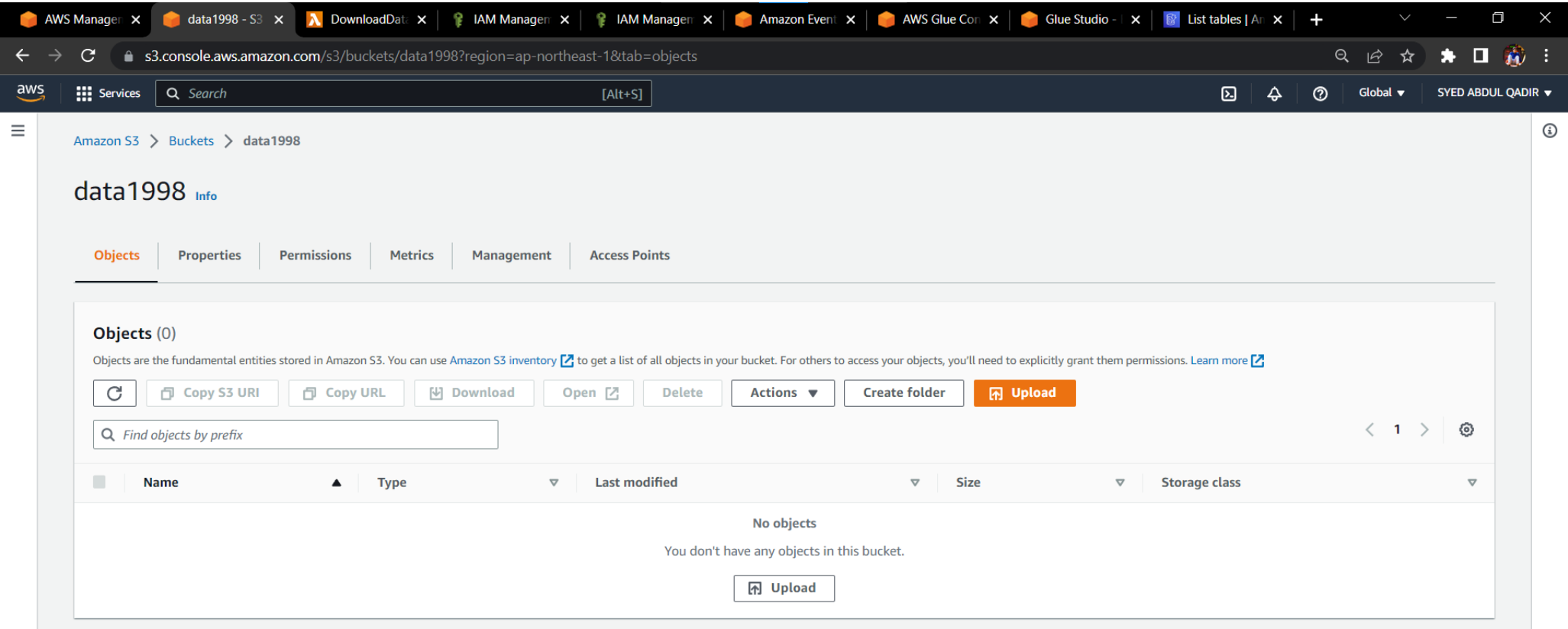
Finally, we create workflow in glue, will we be dumping data in to Dynamo DB.
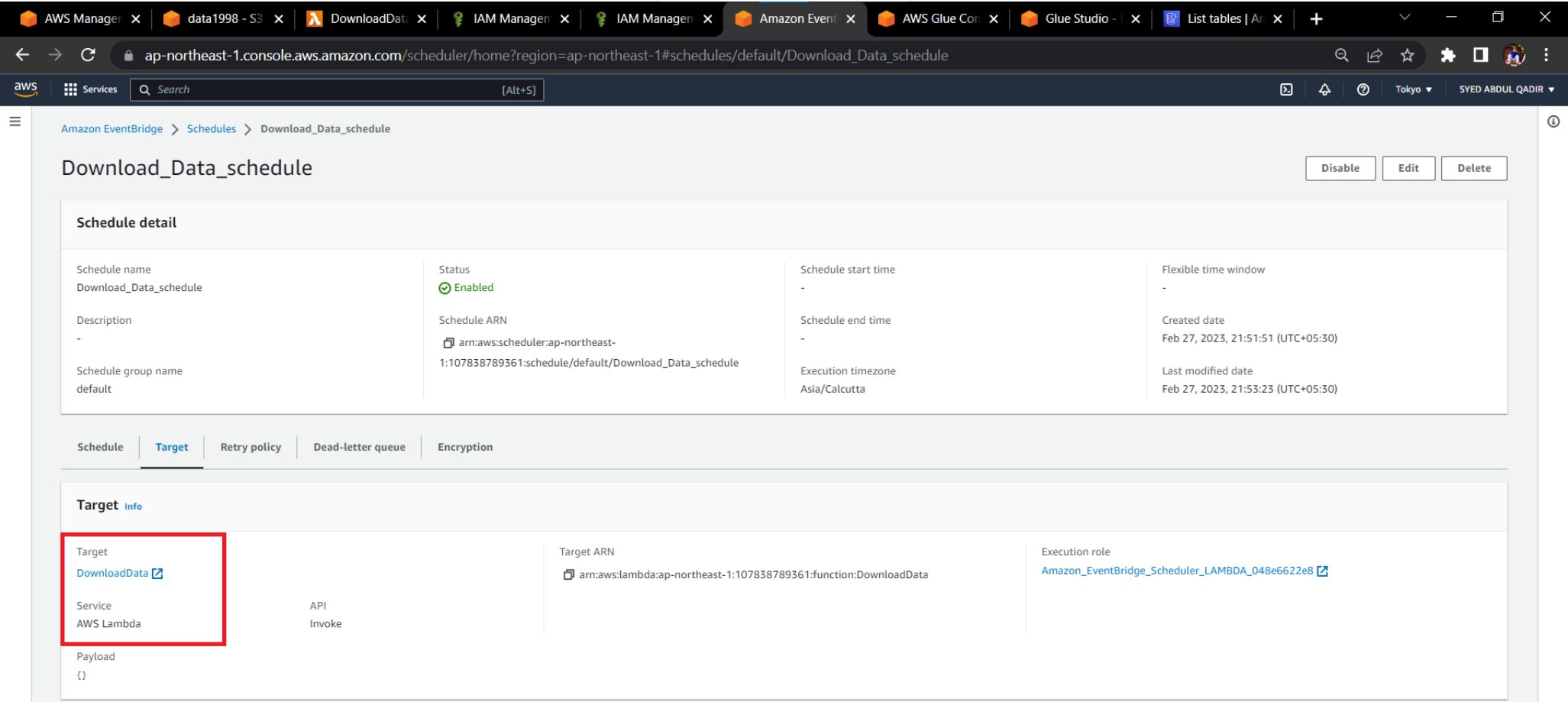
# Implementation:

Mongo db. without dumping any records into it from api



Aws s3 bucket with name **data1998** without dumping any records.



Configuring amazon event bridge schedules, Name it as **Download_Data_schedule** and will be scheduling it for every 2 min and select targeted lambda with name **DowloadData**

Attaching Permission policies with name **S3_lambda_cloud_glue_dynamo_db**



Configuring AWS Lambda with name **DownloadData ,** where we will be dumping api source data to mongo dB and to s3 and we will be filling configuration for environment variables as below.

## Lambda function code:

```python
1.  import json
2.  import pymongo
3.  import certifi
4.  import logging
5.  import os
6.  import boto3
7.  import datetime
8.  import os
9.  import requests
10.
11.     ca = certifi.where()
12.     import os
13.     DATABASE_NAME = os.getenv("DATABASE_NAME")
14.     COLLECTION_NAME = os.getenv("COLLECTION_NAME")
15.     MONGODB_URL = os.getenv("MONGODB_URL")
16.     BUCKET_NAME=os.getenv("BUCKET_NAME")
17.
18.     DATA_SOURCE_URL = f"https://www.consumerfinance.gov/data-research/consumer-complaints/search/api/v1/" \
19.                       f"?date_received_max=<todate>&date_received_min=<fromdate>" \
20.                       f"&field=all&format=json"
21.     client = pymongo.MongoClient(MONGODB_URL, tlsCAFile=ca)
22.
23.     def get_from_date_to_date():
24.         from_date = "2023-01-01"
25.         from_date = datetime.datetime.strptime(from_date, "%Y-%m-%d")
26.
27.         if COLLECTION_NAME in client[DATABASE_NAME].list_collection_names():
28.
29.             res = client[DATABASE_NAME][COLLECTION_NAME].find_one(sort=[("to_date", pymongo.DESCENDING)])
30.             if res is not None:
31.                 from_date = res["to_date"]
32.
33.         to_date = datetime.datetime.now() #current date
34.
35.         response = {
36.             "form_date": from_date.strftime("%Y-%m-%d"),
37.             "to_date": to_date.strftime("%Y-%m-%d"),
38.             "from_date_obj": from_date,
39.             "to_date_obj": to_date
40.         }
41.         logging.info(f"From date and to date {response}")
42.         return response
43.
44.     def save_from_date_to_date(data, status=True):
45.         data.update({"status": status})
46.         logging.info(f"saving from data and to date {data}")
47.         client[DATABASE_NAME][COLLECTION_NAME].insert_one(data)
48.
49.     def lambda_handler(event, context):
50.         print(event,context)
```

```
51.        from_date, to_date, from_date_obj, to_date_obj = get_from_date_to_date().values()
52.        if to_date==from_date:
53.            return {
54.                'statusCode': 200,
55.                'body': json.dumps('Pipeline has already downloaded all data upto yesterday')
56.            }
57.        url = DATA_SOURCE_URL.replace("<todate>", to_date).replace("<fromdate>", from_date)
58.        data = requests.get(url, params={'User-agent': f'your bot '})
59.
60.        finance_complaint_data = list(map(lambda x: x["_source"],
61.                                filter(lambda x: "_source" in x.keys(),
62.                                    json.loads(data.content)))
63.                            )
64.        s3 = boto3.resource('s3')
65.        s3object = s3.Object(BUCKET_NAME, f"inbox/{from_date.replace('-','_')}_{to_date.replace('-','_')}_finance_complaint.json")
66.        s3object.put(
67.            Body=(bytes(json.dumps(finance_complaint_data).encode('UTF-8')))
68.        )
69.
70.        save_from_date_to_date({"from_date": from_date_obj, "to_date": to_date_obj})
71.        return {
72.            'statusCode': 200,
73.            'body': json.dumps('Hello from Lambda!')
74.    }
```

Able to dump data to mongo Db and as well as s3 Bucket.

Creating DynamoDB table with name **fc_data** with Partition key **complaint_id** and sort key as **product**





Glue script job as **s3_data_to_dynamodb** code:

```
1.  import sys
2.  from awsglue.transforms import *
3.  from awsglue.utils import getResolvedOptions
4.  from pyspark.context import SparkContext
5.  from pyspark.sql import functions as func
6.  from awsglue.context import GlueContext
7.  from awsglue.dynamicframe import DynamicFrame
8.  from pyspark.sql.types import LongType
9.  from awsglue.job import Job
10.     import os
11.     ## @params: [JOB_NAME]
12.     args = getResolvedOptions(sys.argv, ['JOB_NAME'])
13.
14.     sc = SparkContext()
15.     glueContext = GlueContext(sc)
16.     spark = glueContext.spark_session
17.     job = Job(glueContext)
18.     job.init(args['JOB_NAME'], args)
19.
20.     #declaring constant variables
21.     BUCKET_NAME="data327030"
22.     DYNAMODB_TABLE_NAME="fc_data"
23.     INPUT_FILE_PATH=f"s3://{BUCKET_NAME}/inbox/*json"
24.
25.     #getting logger object to log the progress
26.     logger  = glueContext.get_logger()
27.     logger.info(f"Started reading json file from {INPUT_FILE_PATH}")
28.     df_sparkdf=spark.read.json(INPUT_FILE_PATH)
29.     logger.info(f"Type casting columns of spark dataframe to Long type")
30.     df_sparkdf = df_sparkdf.withColumn("complaint_id",func.col("complaint_id").cast(LongType()))
31.
32.     logger.info(f"Columns in dataframe : {len(df_sparkdf.columns)}--> {df_sparkdf.columns}")
33.     logger.info(f"Number of rows found in file: {df_sparkdf.count()} ")
34.
35.     dyf = glueContext.create_dynamic_frame.from_options(
36.         connection_type="dynamodb",
37.         connection_options={"dynamodb.input.tableName": DYNAMODB_TABLE_NAME,
```
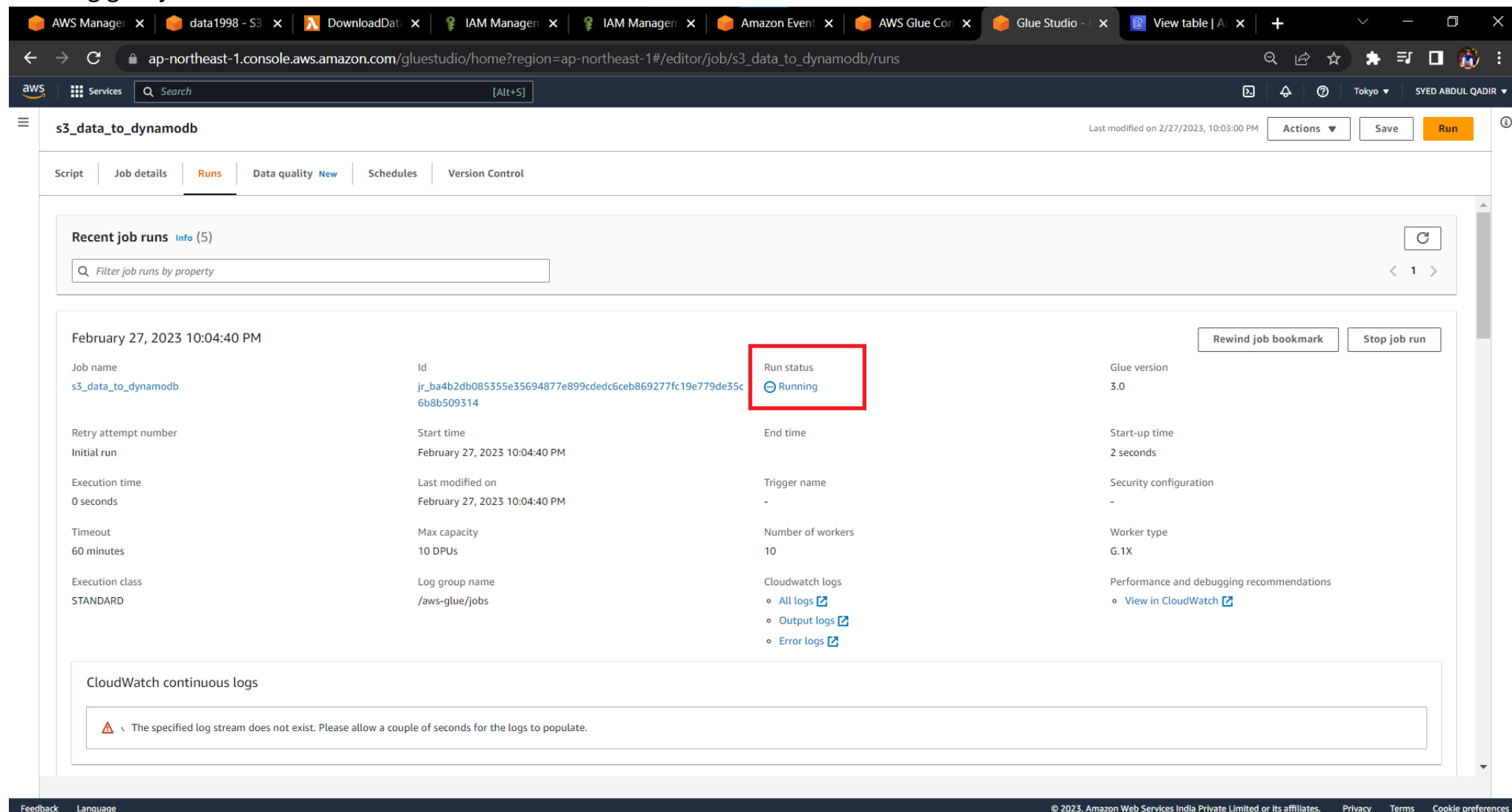
```python
38.                "dynamodb.throughput.read.percent": "1.0",
39.                "dynamodb.splits": "100"
40.             }
41.         )
42.     dyf_sparkdf=dyf.toDF()
43.     new_sparkdf=None
44.     if dyf_sparkdf.count()!=0:
45.         logger.info(f"Columns in dynamodb dataframe : {len(dyf_sparkdf.columns)}--> {dyf_sparkdf.columns}")
46.         logger.info(f"Number of rows found in file: {dyf_sparkdf.count()} ")
47.         logger.info(f"Renaming exiting complaint id column of dynamodb ")
48.         existing_complaint_spark_df =
   dyf_sparkdf.select("complaint_id").withColumnRenamed("complaint_id","existing_complaint_id")
49.         logger.info(f"Applying left join on new dataframe from s3 and dynamo db ")
50.         joined_sparkdf =
   df_sparkdf.join(existing_complaint_spark_df,df_sparkdf.complaint_id==existing_complaint_spark_df.existing_complaint_id,"l
   eft")
51.         logger.info(f"Number of row after left join : {joined_sparkdf.count()}")
52.         new_sparkdf = joined_sparkdf.filter("existing_complaint_id is null")
53.         new_sparkdf.drop("existing_complaint_id")
54.         new_sparkdf=new_sparkdf.coalesce(10)
55.     else:
56.         new_sparkdf=df_sparkdf.coalesce(10)
57.
58.     logger.info(f"Converting spark dataframe to DynamicFrame")
59.     newDynamicFrame= DynamicFrame.fromDF(new_sparkdf, glueContext, "new_sparkdf")
60.     logger.info(f"Started writing new records into dynamo db dataframe.")
61.     logger.info(f"Number of records will be written to dynamodb: {new_sparkdf.count()}")
62.     glueContext.write_dynamic_frame_from_options(
63.         frame=newDynamicFrame,
64.         connection_type="dynamodb",
65.         connection_options={"dynamodb.output.tableName": DYNAMODB_TABLE_NAME,
66.             "dynamodb.throughput.write.percent": "1.0"
67.         }
68.     )
69.
70.     logger.info(f"Data has been dumped into dynamodb ")
71.     logger.info(f"Archiving file from inbox source: s3://{BUCKET_NAME}/inbox  to archive: s3://{BUCKET_NAME}/archive ")
72.     os.system(f"aws s3 sync s3://{BUCKET_NAME}/inbox s3://{BUCKET_NAME}/archive")
73.
74.     logger.info(f"File is successfully archived.")
75.     os.system(f"aws s3 rm s3://{BUCKET_NAME}/inbox/ --recursive")
76.
77. job.commit()
```

Running glue job below

**Screenshot 1**

s3_data_to_dynamodb

Last modified on 2/27/2023, 10:03:00 PM    Actions ▼    Save    Run

Script | Job details | Runs | Data quality New | Schedules | Version Control

**Recent job runs** Info (5)

Filter job runs by property

< 1 >

**February 27, 2023 10:04:40 PM**

Rewind job bookmark    Stop job run

| | | | |
|---|---|---|---|
| Job name | Id | Run status | Glue version |
| s3_data_to_dynamodb | jr_ba4b2db085355e35694877e899cdedc6ceb869277fc19e779de35c6b8b509314 | ⊘ Running | 3.0 |
| Retry attempt number | Start time | End time | Start-up time |
| Initial run | February 27, 2023 10:04:40 PM | - | 7 seconds |
| Execution time | Last modified on | Trigger name | Security configuration |
| 1 minute 13 seconds | February 27, 2023 10:04:43 PM | - | - |
| Timeout | Max capacity | Number of workers | Worker type |
| 60 minutes | 10 DPUs | 10 | G.1X |
| Execution class | Log group name | Cloudwatch logs | Performance and debugging recommendations |
| STANDARD | /aws-glue/jobs | ○ All logs 🗗<br>○ Output logs 🗗<br>○ Error logs 🗗 | ○ View in CloudWatch 🗗 |

**CloudWatch continuous logs**

⌄ Driver logs

Driver and executor log streams 🗗

```
23/02/27 16:36:08 INFO DAGScheduler: Job 0 is finished. Cancelling potential speculative or zombie tasks for this job
23/02/27 16:36:02 INFO MultipartUploadOutputStream: close closed:false s3://aws-glue-assets-107838789361-ap-northeast-1/sparkHistoryLogs/spark-application-1677515741306.inprogress
```

---

**Screenshot 2**

**CloudWatch continuous logs**

⌄ Driver logs

Driver and executor log streams 🗗

```
23/02/27 16:36:18 INFO DAGScheduler: Job 3 is finished. Cancelling potential speculative or zombie tasks for this job
23/02/27 16:36:15 INFO DAGScheduler: Got job 3 (fromRDD at DynamicFrame.scala:322) with 10 output partitions
23/02/27 16:36:15 INFO GlueContext: The DataSource in action : com.amazonaws.services.glue.ConnectionDataSource
23/02/27 16:36:15 INFO GlueContext: Glue secret manager integration: secretId is not provided.
23/02/27 16:36:15 INFO GlueLogger: Number of rows found in file: 156849
23/02/27 16:36:15 INFO DAGScheduler: Job 2 finished: count at NativeMethodAccessorImpl.java:0, took 0.522618 s
23/02/27 16:36:15 INFO DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job
23/02/27 16:36:14 INFO DAGScheduler: Got job 2 (count at NativeMethodAccessorImpl.java:0) with 1 output partitions
23/02/27 16:36:08 INFO GlueLogger: Columns in dataframe : 18--> ['company', 'company_public_response', 'company_response', 'complaint_id', 'complaint_what_happened', 'consumer_consent_provi
23/02/27 16:36:08 INFO GlueLogger: Type casting columns of spark dataframe to Long type
23/02/27 16:36:08 INFO DAGScheduler: Job 0 finished: json at NativeMethodAccessorImpl.java:0, took 18.492476 s
23/02/27 16:36:08 INFO DAGScheduler: Job 0 is finished. Cancelling potential speculative or zombie tasks for this job
23/02/27 16:36:02 INFO MultipartUploadOutputStream: close closed:false s3://aws-glue-assets-107838789361-ap-northeast-1/sparkHistoryLogs/spark-application-1677515741306.inprogress
23/02/27 16:35:50 INFO DAGScheduler: Got job 0 (json at NativeMethodAccessorImpl.java:0) with 29 output partitions
23/02/27 16:35:44 INFO GlueLogger: Started reading json file from s3://data1998/inbox/*json
23/02/27 16:35:44 INFO GlueContext: GlueMetrics configured and enabled
23/02/27 16:35:40 INFO Utils: Successfully started service 'sparkDriver' on port 36637.
```

Finally able to dump data form s3 to Dynamo DB.