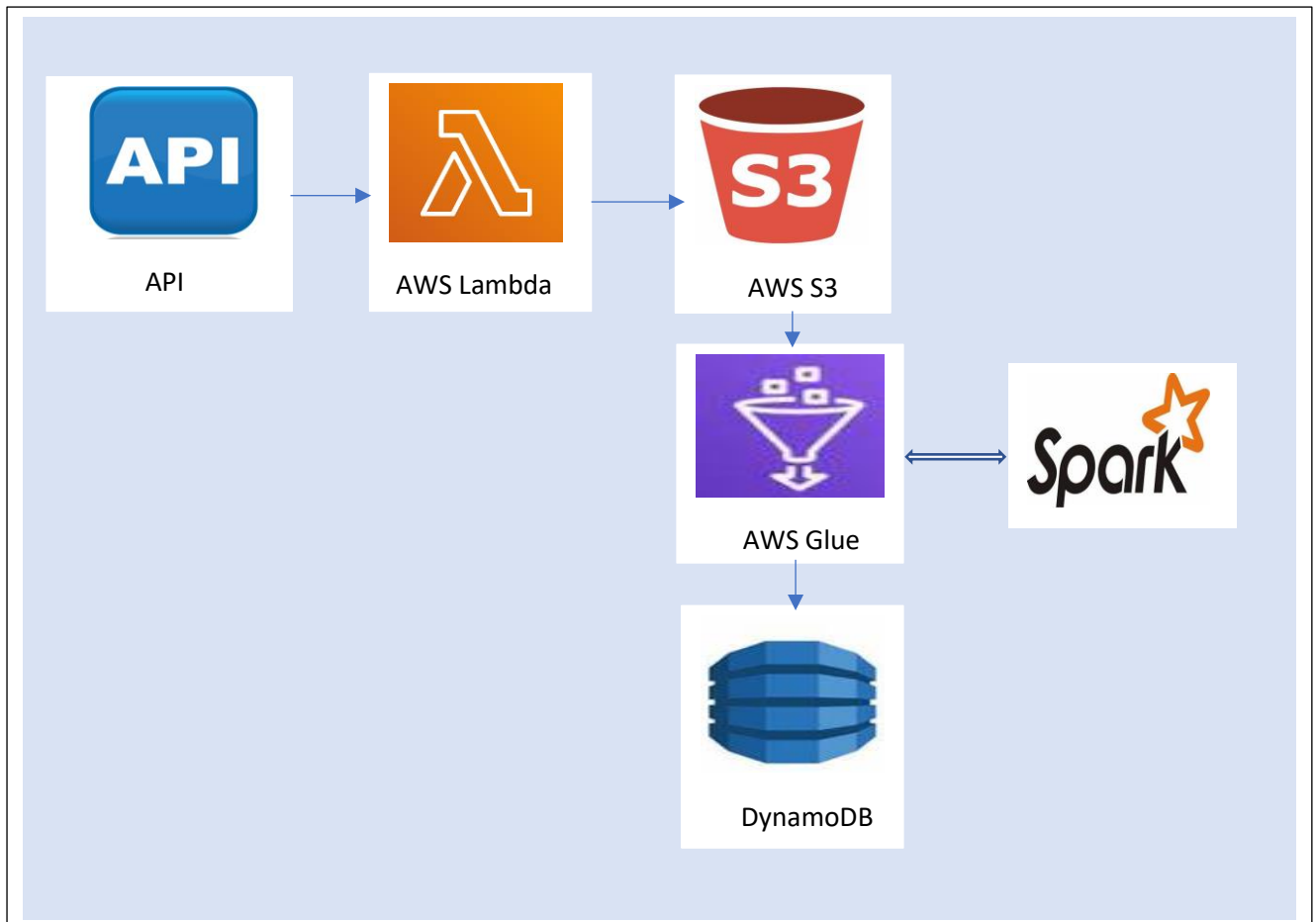# AWS PROJECT Version 2

_____

## Project Architecture:



## AWS Lambda:

Web API -----> JSON Data FROM_DATE – TO_DATE -----> S3 Bucket

This Lambda function retrieves financial complaint data from the Consumer Finance Protection Bureau (CFPB) API between a specified date range. The date range is determined by checking if there is any previously downloaded data in a MongoDB database and using that to set the from_date, otherwise it defaults to a hardcoded from_date. The data is then filtered and saved in a JSON format in an S3 bucket. Finally, the from_date and to_date are saved in the MongoDB database for future reference.

Implementation for AWS Lambda:

```python
lambda_function.py  ×
lambda_function.py > ...
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 = "2023-02-08"

PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL                              Pip package upd

Loading configuration....
Done loading configuration
```

```
File   Edit   Selection   View   Go   Run   Terminal   Help

≡ requirements.txt  ×

boto3 > ≡ requirements.txt
1    pymongo[srv]
2    boto3
3    requests
4
```

conda create -p venv python==3.9 -y --> to create a new virtual environment using conda

conda activate venv/                    --> to activate the virtual environment

pip install -r requirements.txt         --> to install the required packages

pip install --platform manylinux2014_x86_64 --target=lambda_function_code --implementation cp

Select all and create a .zip file:

Big Data  >  AWS Project Version 2  >  lambda_function_code

| Name | Date modified | Type | Size |
|---|---|---|---|
| __pycache__ | 10-02-2023 15:43 | File folder | |
| bin | 10-02-2023 15:43 | File folder | |
| boto3 | 10-02-2023 15:44 | File folder | |
| boto3-1.26.68.dist-info | 10-02-2023 15:44 | File folder | |
| botocore | 10-02-2023 15:44 | File folder | |
| botocore-1.29.68.dist-info | 10-02-2023 15:44 | File folder | |
| bson | 10-02-2023 15:43 | File folder | |
| certifi | 10-02-2023 15:43 | File folder | |
| certifi-2022.12.7.dist-info | 10-02-2023 15:43 | File folder | |
| charset_normalizer | 10-02-2023 15:43 | File folder | |
| charset_normalizer-3.0.1.dist-info | 10-02-2023 15:43 | File folder | |
| dateutil | 10-02-2023 15:43 | File folder | |
| dns | 10-02-2023 15:43 | File folder | |
| dnspython-2.3.0.dist-info | 10-02-2023 15:43 | File folder | |
| gridfs | 10-02-2023 15:43 | File folder | |
| idna | 10-02-2023 15:43 | File folder | |
| idna-3.4.dist-info | 10-02-2023 15:43 | File folder | |
| jmespath | 10-02-2023 15:43 | File folder | |
| jmespath-1.0.1.dist-info | 10-02-2023 15:43 | File folder | |
| pymongo | 10-02-2023 15:43 | File folder | |
| pymongo-4.3.3.dist-info | 10-02-2023 15:44 | File folder | |
| python_dateutil-2.8.2.dist-info | 10-02-2023 15:43 | File folder | |
| requests | 10-02-2023 15:43 | File folder | |
| requests-2.28.2.dist-info | 10-02-2023 15:43 | File folder | |
| s3transfer | 10-02-2023 15:44 | File folder | |
| s3transfer-0.6.0.dist-info | 10-02-2023 15:44 | File folder | |
| six-1.16.0.dist-info | 10-02-2023 15:43 | File folder | |

| Name | Date modified |
|---|---|
| botocore | 10-02-2023 15:44 |
| botocore-1.29.68.dist-info | 10-02-2023 15:44 |
| bson | 10-02-2023 15:43 |
| certifi | 10-02-2023 15:43 |
| certifi-2022.12.7.dist-info | 10-02-2023 15:43 |
| charset_normalizer | 10-02-2023 15:43 |
| charset_normalizer-3.0.1.dist-info | 10-02-2023 15:43 |
| dateutil | 10-02-2023 15:43 |
| dns | 10-02-2023 15:43 |
| dnspython-2.3.0.dist-info | 10-02-2023 15:43 |
| gridfs | 10-02-2023 15:43 |
| idna | 10-02-2023 15:43 |
| idna-3.4.dist-info | 10-02-2023 15:43 |
| jmespath | 10-02-2023 15:43 |
| jmespath-1.0.1.dist-info | 10-02-2023 15:43 |
| pymongo | 10-02-2023 15:43 |
| pymongo-4.3.3.dist-info | 10-02-2023 15:44 |
| python_dateutil-2.8.2.dist-info | 10-02-2023 15:43 |
| requests | 10-02-2023 15:43 |
| requests-2.28.2.dist-info | 10-02-2023 15:43 |
| s3transfer | 10-02-2023 15:44 |
| s3transfer-0.6.0.dist-info | 10-02-2023 15:44 |
| six-1.16.0.dist-info | 10-02-2023 15:43 |
| urllib3 | 10-02-2023 15:43 |
| urllib3-1.26.14.dist-info | 10-02-2023 15:43 |
| aws_application | 10-02-2023 15:53 |
| lambda_function | 10-02-2023 15:14 |
| six | 10-02-2023 15:43 |

We will upload this .zip file to AWS Lambda Function later.

# AWS S3:

Create S3 Bucket -> finance-complaint-data:

Create two folder inbox and archive inside bucket finance-complaint-data:

Inbox -> to store the downloaded data.

Archive -> to store processed data.

# AWS Lambda:

Create lambda function download_data:

**Update IAM Role:** Go to configuration-> open role -> attach policy-> AWSS3FullAccess



**Update memory, timeout and other basic settings according to your requirement:**

Add code through .zip file:

Add environment variables:



MONGODB_URL to connect to the mongodb cluster:

# AWS EventBridge:

Add trigger to schedule lambda function using EventBridge:

Add target -> in this case lambda function download_data is the target.

Amazon EventBridge 〉 Rules

# Rules

A rule watches for specific types of events. When a matching event occurs, the event is routed to the targets associated with the rule. A rule can be associated with one or more targets.

## Select event bus

Event bus
Select or enter event bus name

| default |
| --- |

### Rules (1/1)

| | Delete | Enable | Edit | CloudFormation Template ▼ |

🔍 Find rules      Any status ▼

| | Name ▲ | Status ▽ | Type ▽ | Description |
| --- | --- | --- | --- | --- |
| ☐ | DownloadDataWeekly | ⊘ Enabled | Scheduled Standard | |

Trigger added.

Lambda 〉 Functions 〉 download_data

# download_data

Throttle  |  🗐 Copy ARN  |  Actions ▼

▼ Function overview   Info



| download_data |
| Layers (0) |

EventBridge (CloudWatch Events)

+ Add trigger

+ Add destination

Description
-

Last modified
30 seconds ago

Function ARN
🗐 arn:aws:lambda:ap-south-1:788659805261:function:download_data

Function URL   Info
-

Code | Test | Monitor | Configuration | Aliases | Versions

### Code source   Info

Upload from ▼

ⓘ The deployment package of your Lambda function "download_data" is too large to enable inline code editing. However, you can still invoke your function.

As soon as the trigger arrives, lambda function will download the data from API and store it in S3 bucket.

Amazon S3 > Buckets > finance-complaint-data > inbox/

## inbox/

Objects | Properties

### Objects (1)

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

Copy S3 URI | Copy URL | Download | Open | Delete | Actions ▼ | Create folder | Upload

Find objects by prefix

| | Name | Type | Last modified | Size |
|---|---|---|---|---|
| | 2023_02_08_2023_02_11_finance_complaint.json | json | February 11, 2023, 14:10:10 (UTC+05:30) | |

This is the data downloaded:

{} 2023_02_08_2023_02_11_finance_complaint.json ×

C: > Users > namrt > Downloads > {} 2023_02_08_2023_02_11_finance_complaint.json > ...
1  [{"product": "Credit reporting, credit repair services, or other personal consumer reports", "complaint_what_happened": "",
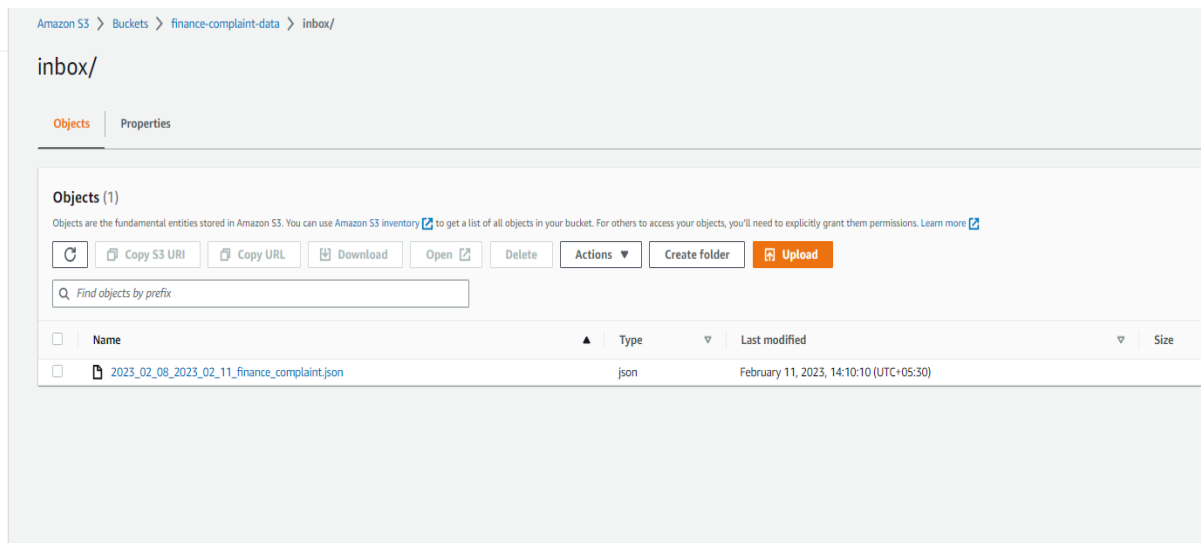   "date_sent_to_company": "2023-02-08T12:00:00-05:00", "issue": "Improper use of your report", "sub_product": "Credit reporting",
   "zip_code": "32780", "tags": null, "complaint_id": "6548160", "timely": "Yes", "consumer_consent_provided": null, "company_response": "In
   progress", "submitted_via": "Web", "company": "EQUIFAX, INC.", "date_received": "2023-02-08T12:00:00-05:00", "state": "FL",
   "consumer_disputed": "N/A", "company_public_response": null, "sub_issue": "Credit inquiries on your report that you don't recognize"},
   {"product": "Credit reporting, credit repair services, or other personal consumer reports", "complaint_what_happened": "",
   "date_sent_to_company": "2023-02-08T12:00:00-05:00", "issue": "Problem with a credit reporting company's investigation into an existing
   problem", "sub_product": "Credit reporting", "zip_code": "27209", "tags": null, "complaint_id": "6548416", "timely": "Yes",
   "consumer_consent_provided": null, "company_response": "In progress", "submitted_via": "Web", "company": "EQUIFAX, INC.",
   "date_received": "2023-02-08T12:00:00-05:00", "state": "NC", "consumer_disputed": "N/A", "company_public_response": null, "sub_issue":
   "Their investigation did not fix an error on your report"}, {"product": "Credit reporting, credit repair services, or other personal
   consumer reports", "complaint_what_happened": "", "date_sent_to_company": "2023-02-08T12:00:00-05:00", "issue": "Problem with a credit
   reporting company's investigation into an existing problem", "sub_product": "Credit reporting", "zip_code": "07011", "tags": null,
   "complaint_id": "6540060", "timely": "Yes", "consumer_consent_provided": null, "company_response": "In progress", "submitted_via": "Web",
   "company": "TRANSUNION INTERMEDIATE HOLDINGS, INC.", "date_received": "2023-02-08T12:00:00-05:00", "state": "NJ", "consumer_disputed": "N/
   A", "company_public_response": null, "sub_issue": "Their investigation did not fix an error on your report"}, {"product": "Credit
   reporting, credit repair services, or other personal consumer reports", "complaint_what_happened": "", "date_sent_to_company":
   "2023-02-08T12:00:00-05:00", "issue": "Incorrect information on your report", "sub_product": "Credit reporting", "zip_code": "33054",
   "tags": null, "complaint_id": "6540512", "timely": "Yes", "consumer_consent_provided": null, "company_response": "In progress",
   "submitted_via": "Web", "company": "EQUIFAX, INC.", "date_received": "2023-02-08T12:00:00-05:00", "state": "FL", "consumer_disputed": "N/
   A", "company_public_response": null, "sub_issue": "Information belongs to someone else"}, {"product": "Credit reporting, credit repair
   services, or other personal consumer reports", "complaint_what_happened": "", "date_sent_to_company": "2023-02-08T12:00:00-05:00",
   "issue": "Incorrect information on your report", "sub_product": "Credit reporting", "zip_code": "30092", "tags": null, "complaint_id":
   "6540441", "timely": "Yes", "consumer_consent_provided": null, "company_response": "In progress", "submitted_via": "Web", "company":

PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL                                              Pip package updater  ∨
Loading configuration....
Done loading configuration

from_date and to_date are saved in the MongoDB database for future reference.

# DynamoDB:

Create DynamoDB table to store the final result.

# AWS Glue:

Will use spark script to write our Glue job

## move_S3_data_to_DynamoDB

| Script | Job details | Runs | Schedules | Version Control |

### Script Info

```python
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="finance-complaint-data"
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)
```

Python    Ln 56, Col 39    ⊗ Errors: 0    ⚠ Warnings: 0

## Glue Job:

The script first reads the JSON file from the specified S3 location and converts it into a Spark DataFrame. Then it creates a DynamicFrame from the existing data in the DynamoDB table and converts it to spark dataframe. If the DynamoDB table has any data, the script performs a left join of the Spark DataFrame from S3 and the spark dataframe from DynamoDB on the "complaint_id" column, filters out the rows where the "existing_complaint_id" is null, and creates a new Spark DataFrame. If the DynamoDB table is empty, the Spark DataFrame from S3 is directly used as the new Spark DataFrame. The new Spark DataFrame is then converted back to a DynamicFrame and written to the DynamoDB table. The script also archives the original JSON file in S3 by copying it to another S3 location. The progress of the script is logged using a logger object.

Add all these policies to IAM Role so that AWS services can interact with one another:



Run the Glue Job:

Can se the live item count in DynamoDB:

Query the DynamoDB Table:



Output:

We can also schedule our Glue job to run on specific frequency.

Here we are selecting weekly run.

As we are downloading our data on Saturday so I am scheduling my job to run on Sunday.