A PROJECT REPORT

*on*

**"Data Migration and Transformation Tool for Amazon NO SQL Data Bases"**

*Submitted to*

# GUVI GEEK NETWORK

In partial fulfilment of the requirements for the award CCM

**IN**

**MASTER DATA ENGINEERING COURSE**

*By*

**PALUKURI SRINIVAS**

**E5 – Data Engineering batch**

*Under the esteemed guidance of*
**DE mentors,**

**"GUVI"**

## GIVEN PROJECT STATEMENT
## Project 1:

| Project Title | Data Migration and Transformation Tool for Amazon NO SQL Data Bases |
|---|---|
| Technologies | Python/PySpark,Requests,Zipfile,boto3,pandas,sqlalchemy, Amazon S3,Amazon NO SQL Data Bases |

**Problem Statement:**

You have a URL that points to a zip file. The zip file contains multiple JSON files. The JSON files contain multiple documents with various data structures. Your goal is to download the zip file from the URL, extract the data from the JSON files, store it in Amazon S3, and load it into Amazon RDS. You want to use Python or PySpark to perform these tasks. You may use any libraries or tools that are necessary to complete the task.

**Approach:**

To extract the data from a zip file that is available at a URL and load it into Amazon S3 and Amazon RDS (NoSQL), you can follow these steps:

1. Use the requests library to download the zip file from the URL.
2. Use the zipfile module to extract the data from the zip file.
3. Use the boto3 library or PySpark to store the data in Amazon S3.
4. Use the pandas library and sqlalchemy or PySpark to load the data from S3 into Amazon RDS (NoSQL).

**Results:**

The result of following these steps should be that the data from the zip file is extracted and stored in a list of dictionaries (if you are using Python) or a DataFrame (if you are using PySpark). Each dictionary or DataFrame row will represent a document from one of the JSON files in the zip file.

The data in the list or DataFrame will then be stored in Amazon S3 as JSON files. You will be able to access these JSON files using the boto3 library or the Amazon S3 web interface.
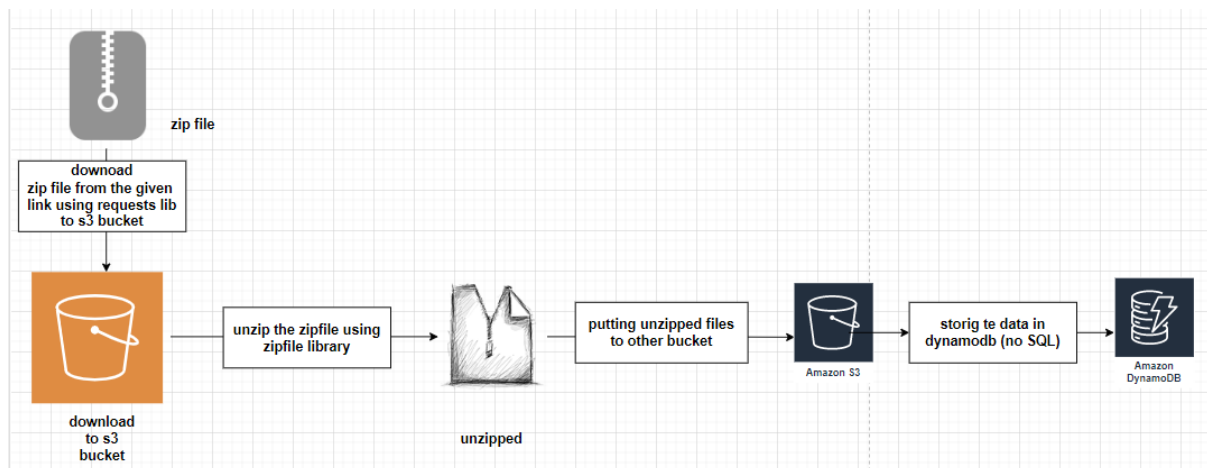
The data from the JSON files will also be loaded into Amazon RDS (NoSQL). You will be able to access the data in RDS using SQL queries. The data will be stored in a table in RDS, and the schema of the table will be determined by the structure of the JSON documents.

I hope this helps. Let me know if you have any further questions or need more assistance.

## Project description:

Data Migration and Transformation Tool for Amazon NOSQL Data Storages.

Project explained in below diagram:



Given link for downloading the zip file : https://www.sec.gov/edgar/sec-api-documentation

By following this link we can get this page.

After reaching this page we can see the option two options for downloading the zip files as shown in the below picture we can download any one for completing this project.

I have selected the submissions.zip because it contains the each organization fillings submissions.

## Bulk data

The most efficient means to fetch large amounts of API data is the bulk archive ZIP files, which are recompiled nightly.

- The companyfacts.zip file contains all the data from the XBRL Frame API and the XBRL Company Facts API

  https://www.sec.gov/Archives/edgar/daily-index/xbrl/companyfacts.zip

- The submission.zip file contains the public EDGAR filing history for all filers from the Submissions API

  https://www.sec.gov/Archives/edgar/daily-index/bulkdata/submissions.zip

Selected link for downloading: https://www.sec.gov/Archives/edgar/daily-index/bulkdata/submissions.zip

So now first part of the project:  downloading the ZIPFILE through given link using requests library.

- Importing the required libraries
    1. requests: for downloading the zip file
    2. boto3: a python SDK for AWS
    3. zipfile: for unzipping the zipped file
    4. io : to do input and output operations

```
C: > Users > Palukuri Srinivas > .aws >  migration_proj.py > ...
  1
  2    #importing requred libraries
  3
  4    import requests
  5    import boto3
  6    import zipfile
  7    import io
  8
```

Motioning AWS credentials to connect the AWS resources to python ide (VS code)

```
10     #mentioning aws credentials
11
12     AWS_ACCESS_KEY_ID = 'AKIA              'F'
13     AWS_SECRET_ACCESS_KEY = 'nUpxE         5o+/mKom              v'
14
15
16     #making the clint connecting to bucket
17
18     s3 = boto3.client('s3',
19                       aws_access_key_id=AWS_ACCESS_KEY_ID,
20                       aws_secret_access_key=AWS_SECRET_ACCESS_KEY
21
22                       )
23
```

Observe the above picture we have given the we have mentioned the **AWS_ACCESS_KEY_ID** and **AWS_SECRET_ACCESS_KEY** to connect the AWS resources.

- Created the AWS S3 bucket for downloading the zip file.

```
#creating AWS S3 bucket

response = s3.create_bucket(Bucket="migration-proj",
                            CreateBucketConfiguration={'LocationConstraint': 'ap-northeast-1'})
print("bucket created")



BUCKET_NAME = 'migration-proj'
```

Note: bucket name must be unique in your region so select the unique name to overcome the invalid location constraint error while executing.

- Downloading the zip file into AWS S3 bucket using requests library and using boto3 putting the downloaded file in to S3 bucket.

```
24     #from here downloading the reqired zipfile has started
25
26     url = 'https://www.sec.gov/Archives/edgar/daily-index/xbrl/companyfacts.zip'
27
28     headers= {"user-agent":"Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/110.0.0.0 Safari/537.36"}
29     print("download started")
30     response = requests.get(url,headers = headers, stream= True)
31
32     f = io.BytesIO(response.content)
33     print("download finished")
34
35     print(response)
36
37     #the downloded zipfile put into s3 bucket
38
39     s3.put_object(Bucket=BUCKET_NAME, Key='input/submissions.zip', Body=f)
40
41     print("object succesfully put into s3 bucket")
42
```

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    CODEWHISPERER REFERENCE LOG                                    Python

```
PS C:\Users\Palukuri Srinivas> & C:/python/python.exe "c:/Users/Palukuri Srinivas/.aws/migration_project_2.py"
download started
download finished
<Response [200]>
object succesfully put into s3 bucket
PS C:\Users\Palukuri Srinivas>
```

The mentioned zip file successfully downloaded and put into mentioned S3 bucket: result picture attached below

Unzipped the files using the zipfile module:

Tried to unzip the file using the zip file it took more than 8hrs of time but task has not completed so then I wanted to check what present in that zip file so downloaded to local storage and extracted using extractall
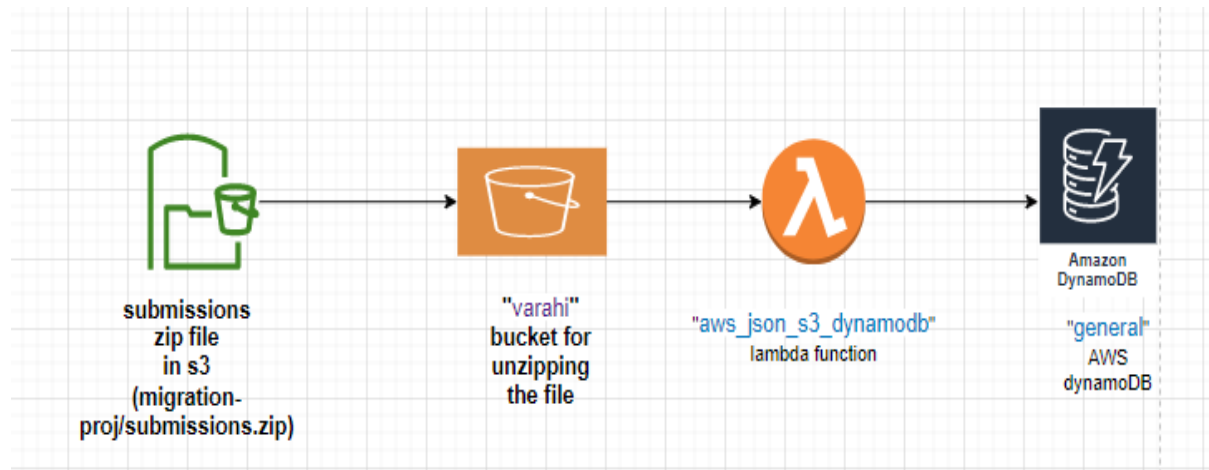
```python
 9    url = 'https://www.sec.gov/Archives/edgar/daily-index/bulkdata/submissions.zip'
10
11    headers= {"user-agent":"Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/110.0.0.0 Safari/537.36"}
12    print("download started")
13    target_loaction = 'D:\GUVI\project\submissions unzip1'
14    response = requests.get(url,headers = headers, stream= True)
15    print(response)
16
17    zip_file = zipfile.ZipFile(io.BytesIO(response.content))
18    zip_file.extrectall(target_loaction)
19
20    print("download finished")
```

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    CODEWHISPERER REFERENCE LOG              Python + ∨ ⬚ 🗑 … ∧

PS C:\Users\Palukuri Srinivas> & C:/python/python.exe "c:/Users/Palukuri Srinivas/.aws/migration_try.py"
download started
<Response [200]>
download finished
PS C:\Users\Palukuri Srinivas>
```

Files downloaded and extracted to local storage.

| Name | Date modified | Type | Size |
|---|---|---|---|
| CIK0000000003 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000000013 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000000014 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000000017 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000000018 | 20-03-2023 07:49 AM | JSON File | 3 KB |
| CIK0000000020 | 20-03-2023 07:49 AM | JSON File | 62 KB |
| CIK0000000049 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000000051 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000000063 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000001750-submissions-001 | 20-03-2023 07:49 AM | JSON File | 147 KB |
| CIK0000001750 | 20-03-2023 07:49 AM | JSON File | 145 KB |
| CIK0000001761 | 20-03-2023 07:49 AM | JSON File | 6 KB |
| CIK0000001800-submissions-001 | 20-03-2023 07:49 AM | JSON File | 266 KB |
| CIK0000001800-submissions-002 | 20-03-2023 07:49 AM | JSON File | 30 KB |
| CIK0000001800 | 20-03-2023 07:49 AM | JSON File | 140 KB |
| CIK0000001830 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000001841 | 20-03-2023 07:49 AM | JSON File | 5 KB |
| CIK0000001848 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000001853 | 20-03-2023 07:49 AM | JSON File | 22 KB |
| CIK0000001860 | 20-03-2023 07:49 AM | JSON File | 2 KB |
| CIK0000001904 | 20-03-2023 07:49 AM | JSON File | 6 KB |

8,35,133 items

After downloading the file to my local storage then I can know that it containing 835133 (8lakhs of json files)

NOTE: Here we have an issue for unzipping, the zip file contains more than 8 lakh of files, its taking more time to unzip and im using AWS free tier account then I have asked the DE MENTOR during Project doubt sessions as per his suggestion to unzip 100 no of json files are enough to complete the project. But we should have the basic idea to unzip all the files in the zipfile. As per mentor suggestion completed the project by unzipping 150 no of json files.
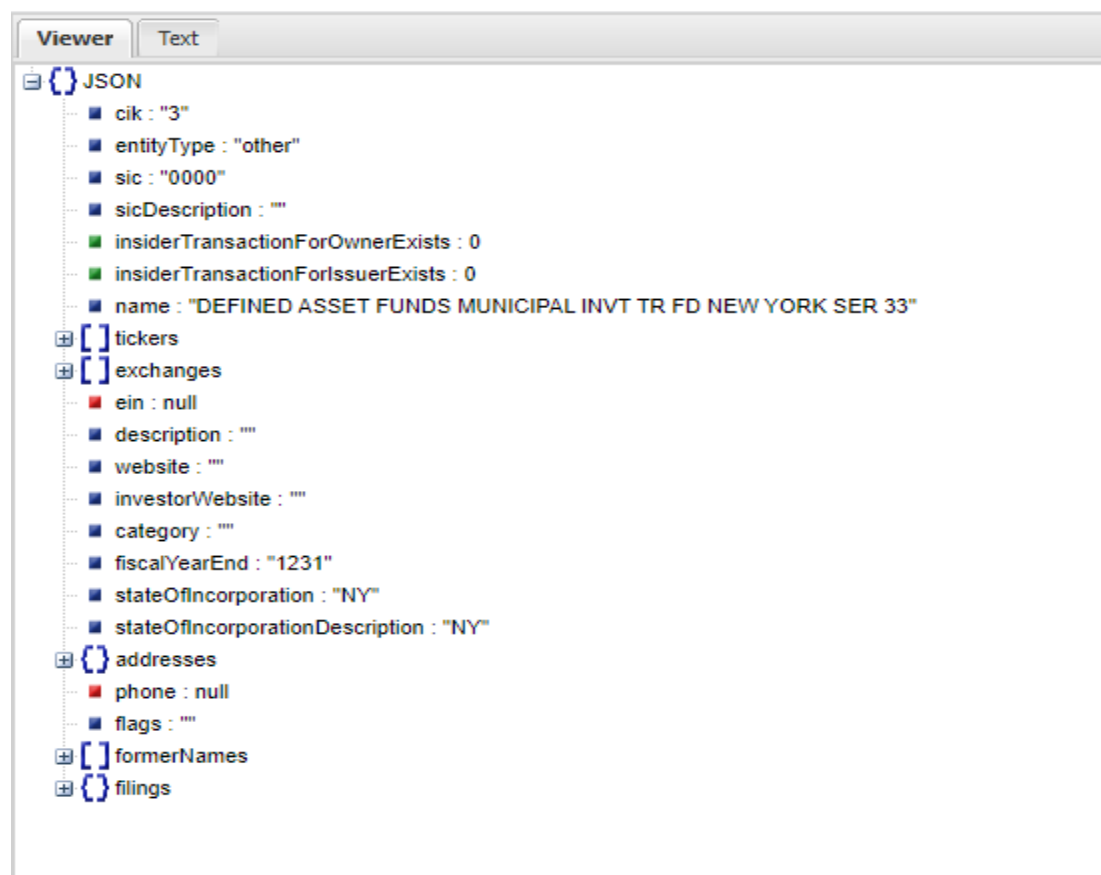
Our strategy to store the json files to dynamoDB from the zip file:

Creating the AWS lambda function which is triggered by S3 "varahi" bucket, then the lambda function will post the data to Dynamo DB "general" as a record.



I have checked the content present in the json in json viewer.

I got the result as shown in below picture and verified many json files then I can find the partition key to create the dynamoDB is "cik" data type is "string" I have observed the many files similar data type and partition key found .



This is that json files containing.

- Creation of dynamo DB for saving the Json data which is present in unzipped files.
- DynamoDB name : general created with partition key as "cik" in "string" data type.

| | Name | | Status | Partition key | Sort key | Indexes | Deletion protection | Read capacity mo |
|---|---|---|---|---|---|---|---|---|
| | general | | ⊘ Active | cik (S) | - | 0 | ⊖ Off | Provisioned with a |

DynamoDB > Tables

**Tables (3)** Info

- Creation of S3 bucket for unzipping, s3 bucket name: "varahi"

```
54   #creating bucket to unzip the file
55
56   print("started")
57
58   response = s3.create_bucket(Bucket="varahi",
59                               CreateBucketConfiguration={'LocationConstraint': 'ap-northeast-1'})
60   print("bucket created")
61
```

- Creation of AWS lambda function with name "aws_json_s3_dynamodb"

Lambda > Functions > aws_json_s3_dynamodb

**aws_json_s3_dynamodb**

▼ **Function overview** Info

aws_json_s3_dynamodb
Layers (0)
S3
+ Add destination
+ Add trigger

Description
-

Last modified
10 days ago

Function ARN
arn:aws:lambda:ap-northeast-1:782483600400:function:aws_json_s3_dynamodb

Added the trigger to lambda function as s3 bucket "varahi" , and Suffix: .json : so that any json file added to "varahi" bucket which have partition key as "cik" in "string" data type can be added to "general" dynamoDB table as an row.

**Triggers (1)** Info

**Trigger**

**S3**: varahi
arn:aws:s3:::varahi

▼ **Details**

Bucket arn: **arn:aws:s3:::varahi**

Event types: **s3:ObjectCreated:Put**

Notification name: **c1cd6c1f-f96d-4dfa-9066-b1e6368d5cad**

Service principal: **s3.amazonaws.com**

Source account: **782483600400**

Statement ID: **lambda-5e26b4ea-e064-4d53-9c7d-e4516fd2f6af**

Suffix: .json

Used lambda code given below:

```
File   Edit   Find   View   Go   Tools   Window        Test  ▼      Deploy

Q      Go to Anything (Ctrl-P)           🗐      lambda_function ×      ⊕

Environment
       ▼ 📁 aws_json_s3_dynar  ⚙▼      1   import json
           </> lambda_function.py         2   import boto3
                                          3
                                          4   def lambda_handler(event, context):
                                          5
                                          6       s3_client = boto3.client('s3')
                                          7       dy_db = boto3.client('dynamodb')
                                          8
                                          9
                                         10       bucket = event['Records'][0]['s3']['bucket']['name']
                                         11       obj_key = event['Records'][0]['s3']['object']['key']
                                         12
                                         13
                                         14       object = s3_client.get_object(Bucket =bucket,Key=obj_key)
                                         15       file = object['Body'].read().decode('utf-8')
                                         16
                                         17       #print(file)
                                         18
                                         19
                                         20       dict = json.loads(file)
                                         21
                                         22       table_name = 'general'
                                         23       table = boto3.resource('dynamodb').Table(table_name)
                                         24       table.put_item(Item=dict)
                                         25
                                         26
                                         27       return {
                                         28           'statusCode': 200,
                                         29           'body': json.dumps('Hello from Lambda!')
                                         30       }
                                         31
```

In the above code we can see :

1. Imported the required libraries.
2. Created S3 client.
3. Created the DynamoDB client.
4. Getting the bucket name and object key by pain event .
5. Retrieving object with bucket name and object key.
6. Reading the content to file variable.
7. Converting the file variable from json to dictionary (as per project statement).
8. Inserting the data to dynamo Db using Put item method.

So to complete our process

1. DynamoDB table created
2. S3 bucket created for unzipping and triggering the lambda function
3. Lambda function also created with required code

Now when any file added to "varahi" we know any Json file have partition key a "cik" with "string"data type it will be added to dynamo db as new record.

- Unzipping the "Submissions.zip" file to "varahi" setting limit as 150 nos, code snippet attached below

s

```python
70    #from here extractting of zip file located in s3 starts unzipping.
71
72    zip_file = zipfile.ZipFile(io.BytesIO(s3_obj['Body'].read()))
73    bucket1='varahi'
74
75    print("extracting the file started to our destination")
76
77
78    i=0
79    for filename in zip_file.namelist():
80
81            unzipped_content = zip_file.read(filename)
82
83            unzipped_key = 'unzipped' + filename
84
85            s3.put_object(Bucket=bucket1, Key= unzipped_key, Body=unzipped_content)
86
87            print("files extracted successsfully to unzipped folder")
88            i=i+1
89            if i==150:
90                break
91    print("process finished")
```

```python
78            i=i+1
79            if i==150:
80                break
81    print("process finished")
```

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    CODEWHISPERER REFERENCE LOG

```
PS C:\Users\Palukuri Srinivas> & C:/python/python.exe "c:/Users/Palukuri Srinivas/.aws/migration_proj.py"
started
bucket created
reading into memory started
extracting the file started to our destination
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
```

Results attached below for executing the above code:

```
PS C:\Users\Palukuri Srinivas> & C:/python/python.exe "c:/Users/Palukuri Srinivas/.aws/migration_proj.py"
started
bucket created
reading into memory started
extracting the file started to our destination
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
files extracted successsfully to unzipped folder
```

Unzipped files to varahi bucket :



We can see the destination as lambda function "aws_json_s3_dynamodb"

Finally as a results the json files updated to dynamoDB as a records (Rows) snippets attached below. We can find the 150 rows updated to "general" dynamoDB.



Cloud watch logs samples attached below as a snippet:

Note: in the above unzipping code we have make limit for 150nos if we remove that each file will be unzipped and uploaded to dynamoDB records.

## CONCLUSION:

- As per project statement we have downloaded the zip file using.
- Zip file is unzipped using zip file python module.
- Using boto3 module uploaded to S3 bucket.
- The json files converted to dictionary and updated to dynamoDB using boto3 SDK

Task completed.