# Bank Transactions Batch Data Processing with AWS

1. This project was about implementing a batch data processing pipeline in AWS to process daily bank transactions stored in JSON files. Goal was to set up an AWS data processing pipeline that automatically processes this data as soon as it lands on S3, transforming and storing it for querying.

2. A daily JSON file containing bank transactions was being dropped into an S3 bucket. Once new file uploaded, data processing triggered automatically with AWS Lambda function.

3. Within Lambda, I invoked the Glue job. Following step by step starting, used AWS Glue to read the JSON file from the S3 bucket



4. Implemented transformations such as filtering out any transactions with null values, and deduplicating any repeated transactions based on transaction_ID. Then converted the JSON format into a columnar format which is considered more optimized for querying

5. Further storing the transformed data back into a separate S3 bucket in parquet form. Also enabled the etl job bookmark to get only new or updated data.



6. At last I set up AWS Athena to set up a table and query on the top of transformed data stored in S3.

7. Used aws cloudWatch to monitor the data processing tasks



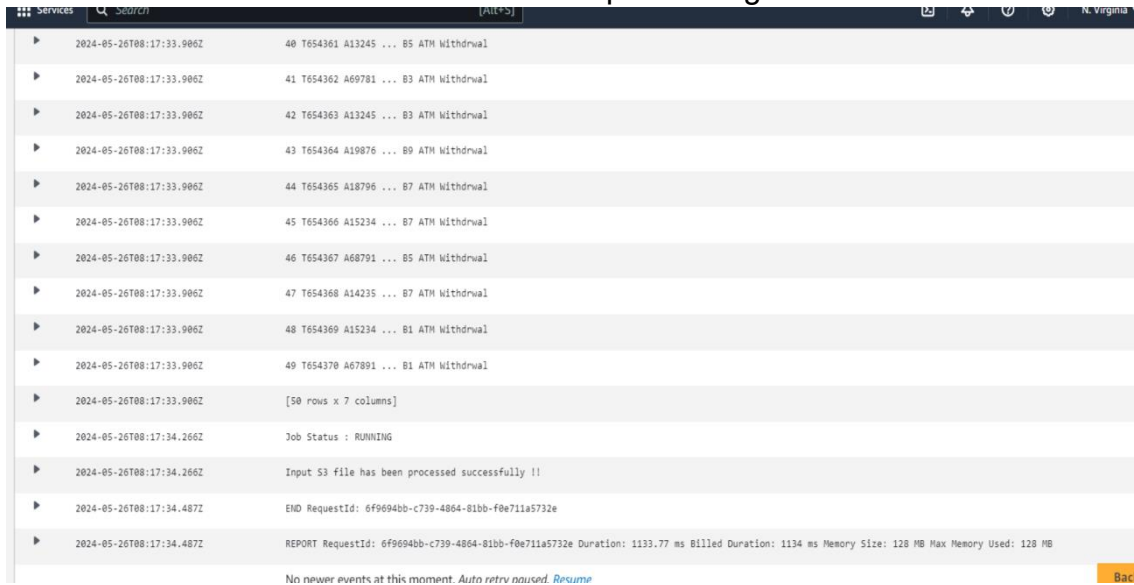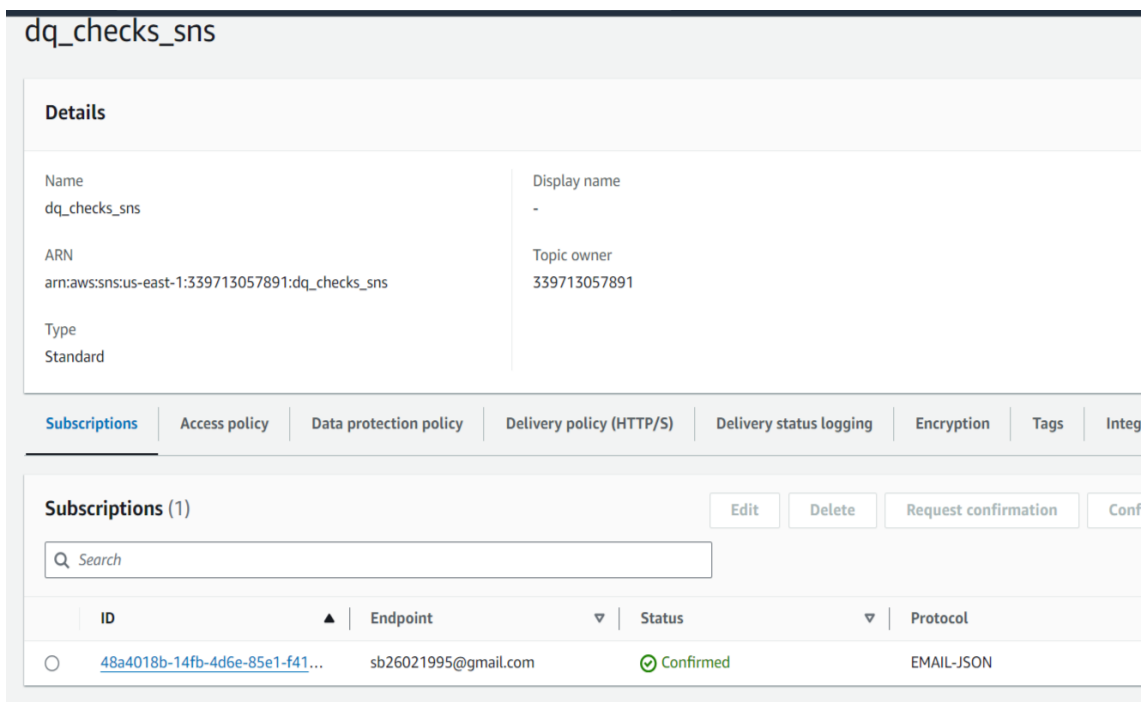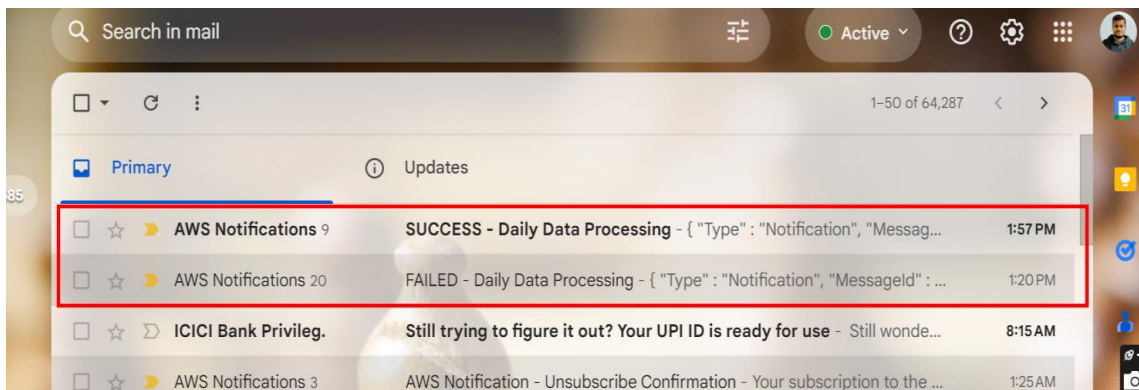| | | |
|---|---|---|
| ▶ | 2024-05-26T08:17:33.906Z | 40 T654361 A13245 ... B5 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 41 T654362 A69781 ... B3 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 42 T654363 A13245 ... B3 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 43 T654364 A19876 ... B9 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 44 T654365 A18796 ... B7 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 45 T654366 A15234 ... B7 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 46 T654367 A68791 ... B5 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 47 T654368 A14235 ... B7 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 48 T654369 A15234 ... B1 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | 49 T654370 A67891 ... B1 ATM Withdrawal |
| ▶ | 2024-05-26T08:17:33.906Z | [50 rows x 7 columns] |
| ▶ | 2024-05-26T08:17:34.266Z | Job Status : RUNNING |
| ▶ | 2024-05-26T08:17:34.266Z | Input S3 file has been processed successfully !! |
| ▶ | 2024-05-26T08:17:34.487Z | END RequestId: 6f9694bb-c739-4864-81bb-f0e711a5732e |
| ▶ | 2024-05-26T08:17:34.487Z | REPORT RequestId: 6f9694bb-c739-4864-81bb-f0e711a5732e Duration: 1133.77 ms Billed Duration: 1134 ms Memory Size: 128 MB Max Memory Used: 128 MB |

No newer events at this moment. Auto retry paused. Resume

8. set up SNS notifications for any failures in the pipeline or if any data quality checks fail



## dq_checks_sns

### Details

| | |
|---|---|
| **Name** | **Display name** |
| dq_checks_sns | - |
| **ARN** | **Topic owner** |
| arn:aws:sns:us-east-1:339713057891:dq_checks_sns | 339713057891 |
| **Type** | |
| Standard | |

| Subscriptions | Access policy | Data protection policy | Delivery policy (HTTP/S) | Delivery status logging | Encryption | Tags | Integ |

### Subscriptions (1)

Edit | Delete | Request confirmation | Conf

Q Search

| | ID | ▲ | Endpoint | ▽ | Status | ▽ | Protocol |
|---|---|---|---|---|---|---|---|
| ○ | 48a4018b-14fb-4d6e-85e1-f41... | | sb26021995@gmail.com | | ⊘ Confirmed | | EMAIL-JSON |



Q Search in mail                                            ⊙ Active ⌄   ?  ⚙  ⋮⋮⋮

☐ ▾  C  ⋮                                            1–50 of 64,287  <  >

📥 **Primary**              ⓘ  Updates

| ☐ ☆ ⟩ | AWS Notifications 9 | **SUCCESS - Daily Data Processing** - { "Type" : "Notification", "Messag... | 1:57 PM |
| ☐ ☆ ⟩ | AWS Notifications 20 | FAILED - Daily Data Processing - { "Type" : "Notification", "MessageId" : ... | 1:20 PM |
| ☐ ☆ ⅀ | ICICI Bank Privileg. | **Still trying to figure it out? Your UPI ID is ready for use** - Still wonde... | 8:15 AM |
| ☐ ☆ ⟩ | AWS Notifications 3 | AWS Notification - Unsubscribe Confirmation - Your subscription to the ... | 1:25 AM |