**CAPSTONE PROJECT**

**ARCHTECTING AND DEPLOYING A DATA LAKE AUTOMATION SOLUTION**



# PROJECT DESCRIPTION

This Data Lake project is designed to help Datasoft Incorporation solve the data needs of its large pool of clients by providing a more streamlined process and platform to enable the company to easily correlate, transform, query, analyze and visualize customers data feeds at every point in time making it easy to generate valuable insights necessary in meeting customers' needs as well as improving the revenue of the company.

# AWS SERVICE TOOLS USED FOR THIS PROJECT

- AWS Lambda
- Amazon Kinesis Data Stream
- Amazon Kinesis Firehose
- Amazon S3
- Amazon Athena
- Amazon SNS TOPIC
- Amazon DynamoDB
- Amazon CloudWatch
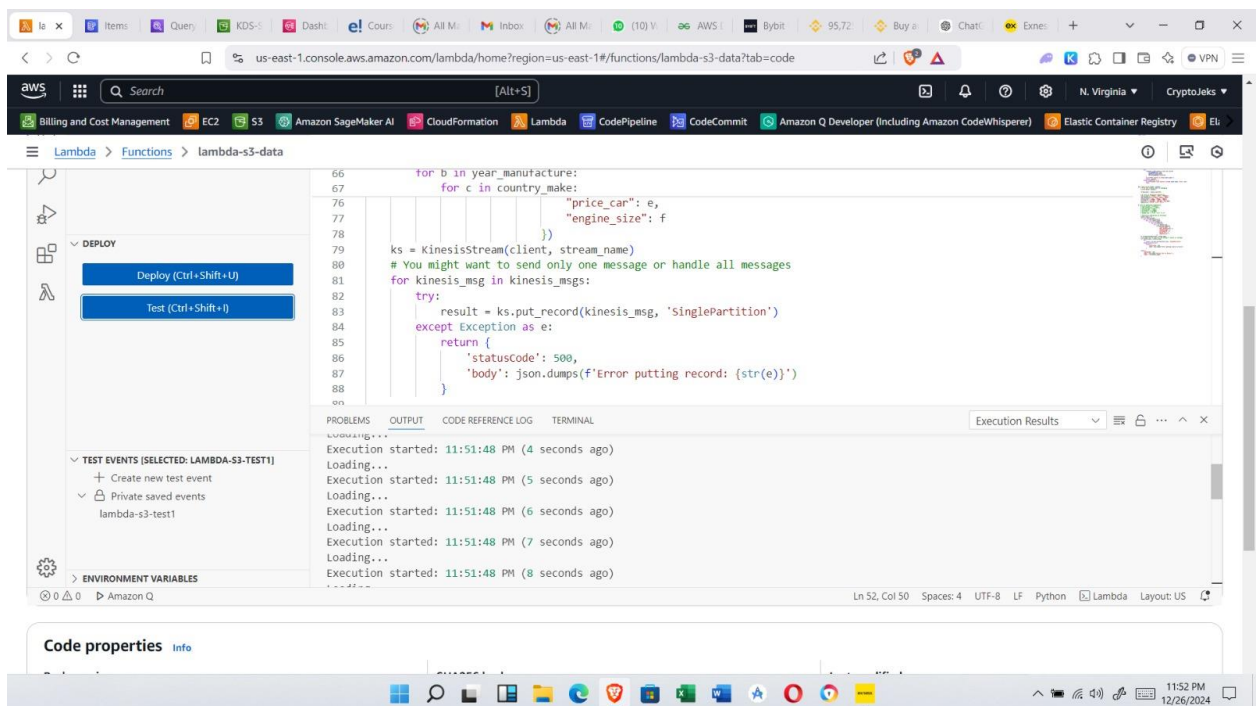- AWS Cognito
- Amazon API Gateway

# PROJECT SECTIONS

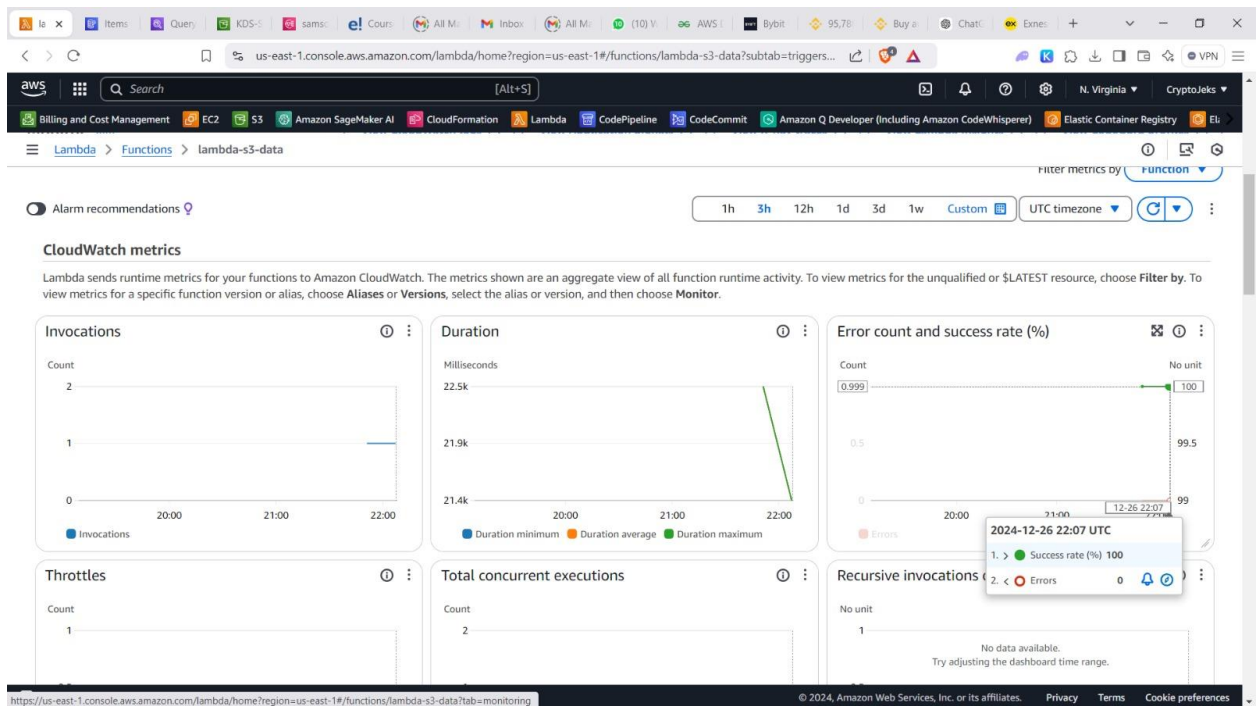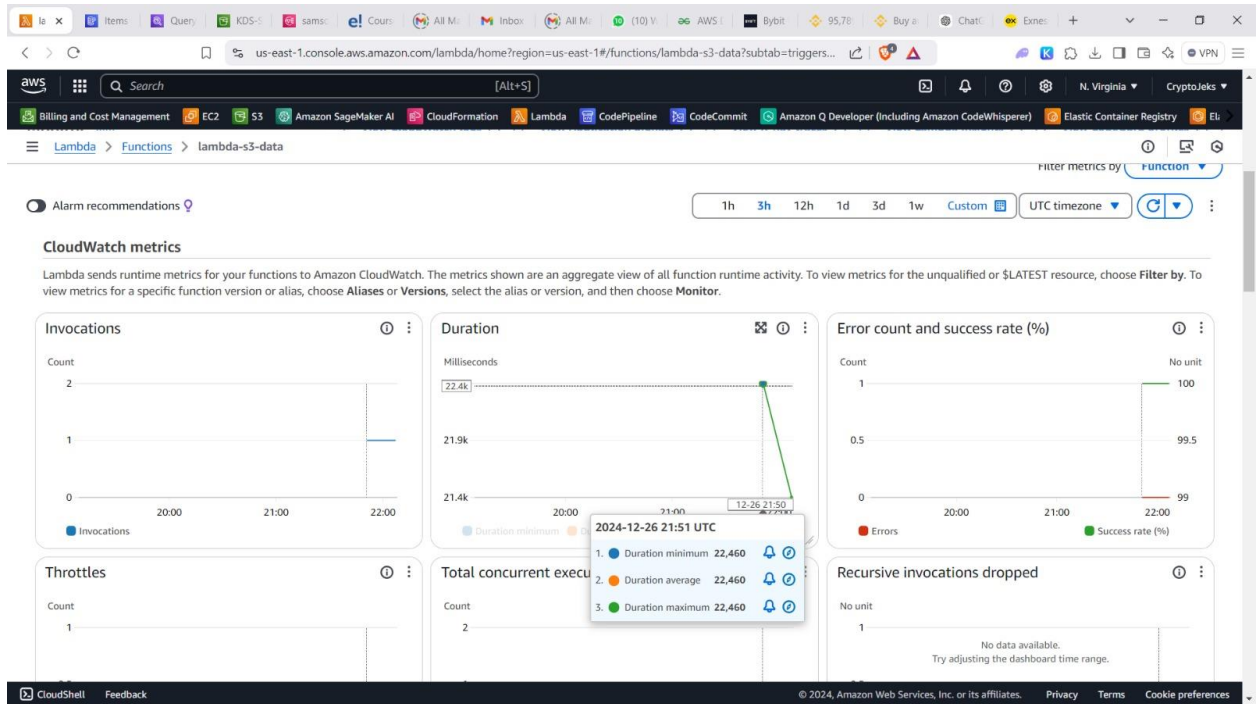This project is divided into three major sections which includes,

- ➢ Data Ingestion Section
- ➢ Data Processing and Storage Section
- ➢ API Gateway Section
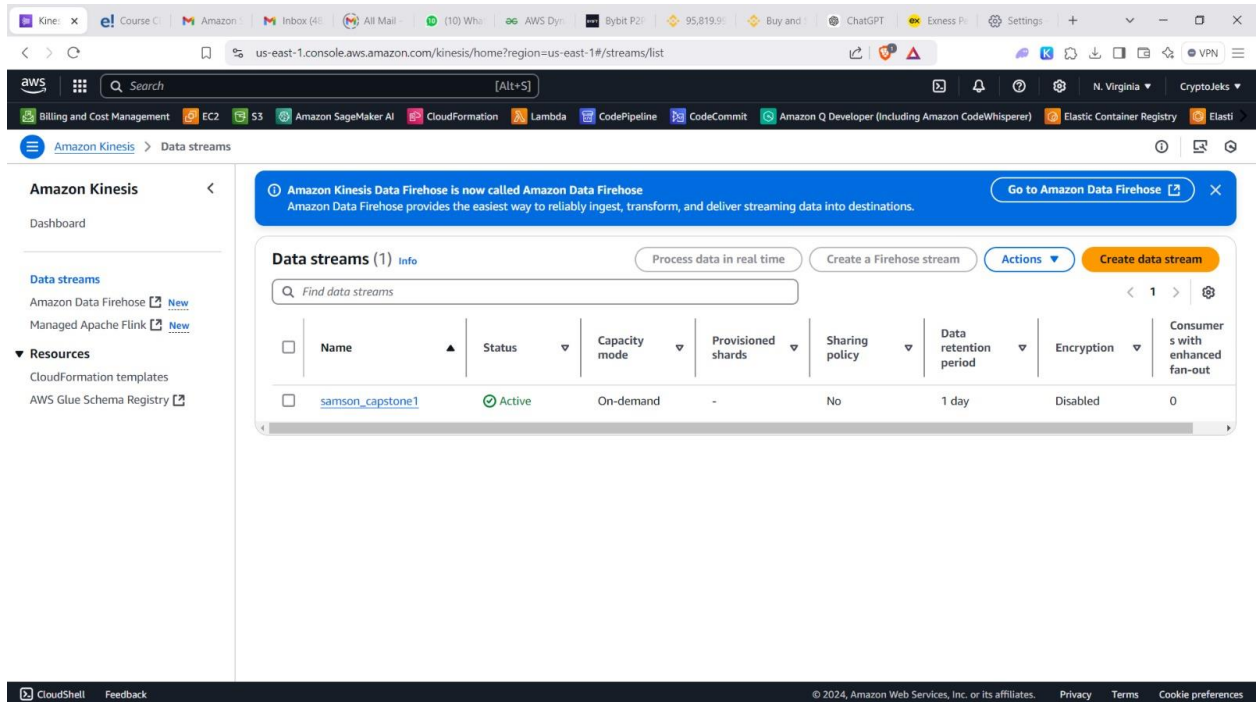
## DATA INJESTION SECTION

This section of the project comprises of AWS services such as AWS lambda, Amazon Kinesis Data Stream, Amazon Kinesis Firehose, and Amazon S3 Bucket. In this section, I ingested a Kinesis Producer Python Code which contains a kinesis message about vehicle information using an AWS lambda function. Kindly click on this link to access my code on Github, https://github.com/CryptoJeks/capstone-code. Using IAM, I assigned a role for the lambda function with full permission to access Amazon Kinesis Data Stream, Kinesis Firehose, and Amazon S3 bucket. This made it possible for my Kinesis Producer Python Code to send the kinesis message to the s3 Bucket via Kinesis Data Stream and Kinesis Firehose. Find below some screenshots of my code and also the data ingestion process.



**The Kinesis Producer Python Code**

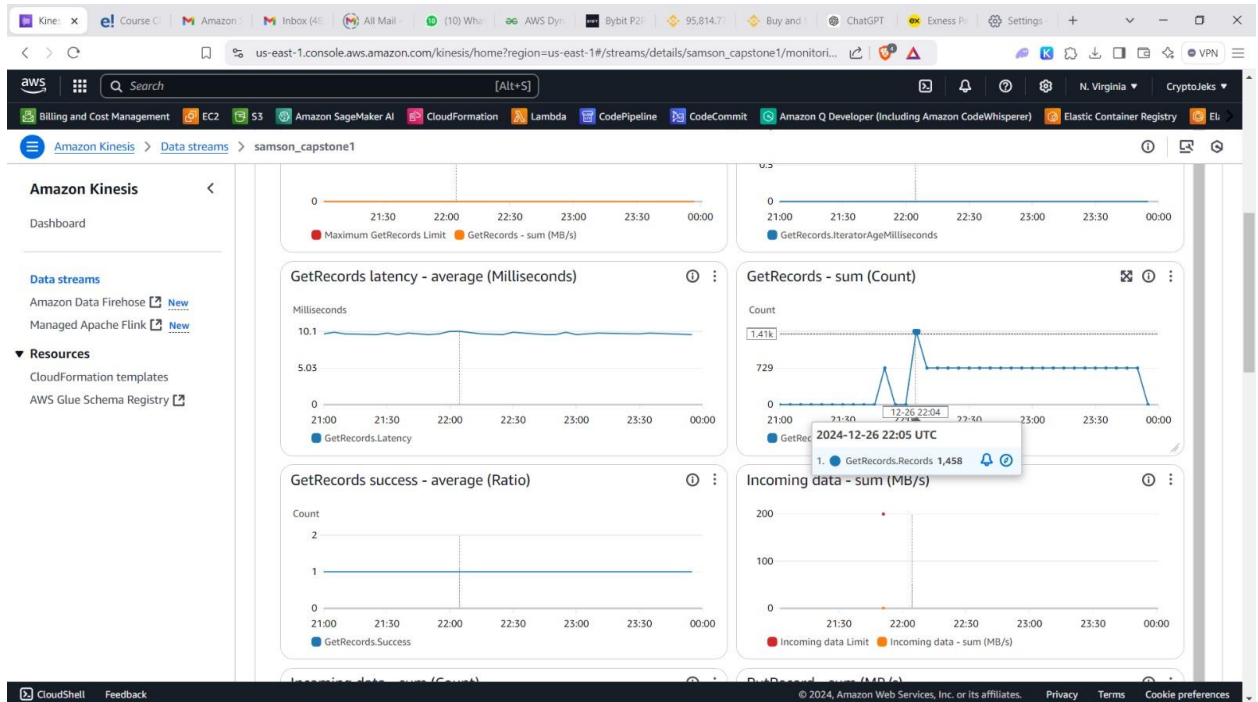**Cloud Watch showing the activities of Kinesis Producer**

**My Amazon Kinesis Data Stream service**

**Real-Time Streaming of Kinesis Message via Kinesis Data Stream**



**My Amazon Kinesis Firehose Stream service**

**Live Streaming of Kinesis Message via Amazon Kinesis Firehose**

**Kinesis Message Successfully pushed to my S3 bucket using a lambda function via kinesis data stream and Firehose.**

## DATA PROCESSING AND STORTAGE SECTION

The data processing section comprises of AWS services such as the S3 Bucket, SNS Topic, AWS Lambda, Amazon DynamoDB and Amazon Athena. In this section, the kinesis message stored in the S3 bucket in JASON format triggers an SNS Topic as soon as there is a PUT into the S3 bucket. SNS Topic then publishes the JASON Kinesis Message into DynamoDB in a visible structural form with the help of a lambda function. Once the kinesis message is published on DynamoDB, Amazon Athena is used to analyze and generate insight from the published Data on DynamoDB. Below is the screenshot of my published kinesis message on my S3 Bucket, SNS Topic and the kinesis message on my DynamoDB.



**Kinesis Message pushed to S3 via AWS Lambda**

# Amazon Event Bridge Creation

Amazon Event Bridge is a serverless service that uses events to connect application components together, making it easier for developers to build scalable event-driven applications. In this case, the Event Bridge is connecting my S3 bucket to DynamoDB by publishing my kinesis message PUT on S3 to DynamoDB using an SNS Topic. Find below some of my live screenshots explaining the process.

Event type
The type of events as the source of the matching
pattern

Amazon S3 Event Notification ▼

ⓘ S3 Event Notifications will only
match your rules if you have
configured your S3 bucket(s) to
publish event notifications to
EventBridge. Learn more.

Event Type Specification 1
◯ Any event
◉ Specific event(s)

Specific event(s)

▼

Object Created ✕   Object Deleted ✕

Object ACL Updated ✕

Event Type Specification 2
◉ Any bucket
◯ Specific bucket(s) by name

Cancel    Previous    Next

---

Step 2
Build event pattern

Step 3
Select target(s)

Step 4 - optional
Configure tags

Step 5
Review and create

# Select target(s)

ⓘ Permissions
Note: When using the EventBridge console, EventBridge will automatically configure the proper permissions for
the selected targets. If you're using the AWS CLI, SDK, or CloudFormation, you'll need to configure the proper
permissions.

## Target 1

Target types
Select an EventBridge event bus, EventBridge API destination (SaaS partner), or another AWS service as a target.
◯ EventBridge event bus
◯ EventBridge API destination
◉ AWS service

Select a target   Info
Select target(s) to invoke when an event matches your event pattern or when schedule is triggered (limit of 5 targets per rule)

SNS topic ▼

Topic

samson_cap ▼   ⟳

▶ Additional settings

Add another target    Cancel    Skip to Review and create    Previous    Next

Amazon EventBridge > Rules > Create rule

# Review and create

## Step 1: Define rule detail
Edit

### Define rule detail

| Rule name | Status | Event bus |
| --- | --- | --- |
| samson_event | ⊘ Enabled | default |
| **Description** | **Rule type** | |
| event from S3 to SNS Topic | Standard rule | |

## Step 2: Build event pattern
Edit

### Event pattern  Info

```
1 {
2   "source": ["aws.s3"],
3   "detail-type": ["Object Created", "Object Deleted", "Object ACL Updated", "Object Restore Initiated"]
4 }
```

Copy

---



# Amazon EventBridge ✕

Dashboard New

▼ Developer resources
  Learn
  Sandbox
  Quick starts

▼ Buses
  Event buses
  **Rules**
  Global endpoints
  Archives
  Replays

▼ Pipes
  Pipes

▼ Scheduler
  Schedules
  Schedule groups

▼ Integration

# 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 (2)

Delete  Enable  Edit  CloudFormation Template ▼  **Create rule**

Find rules | Any status ▼ | ‹ 1 › ⚙

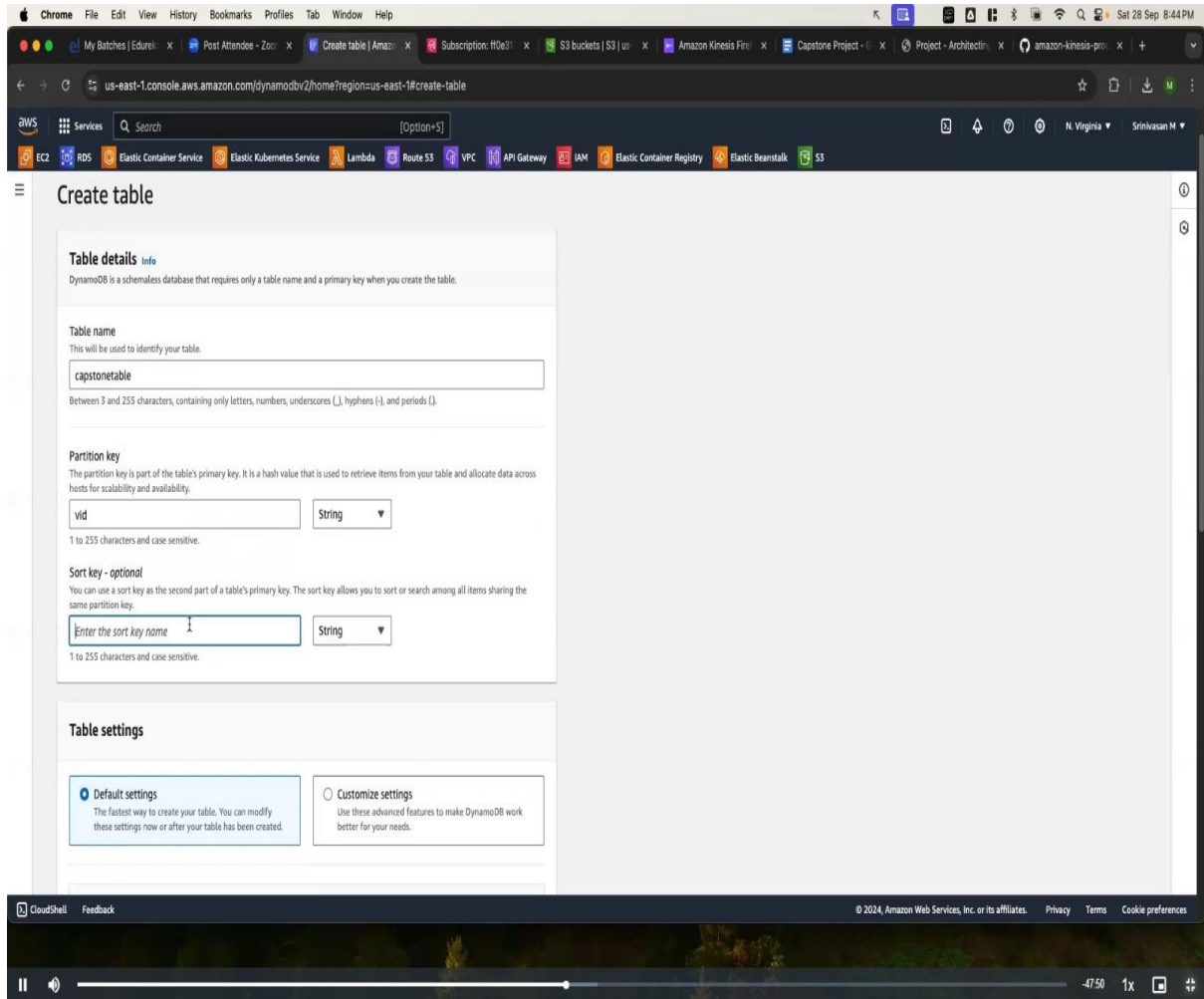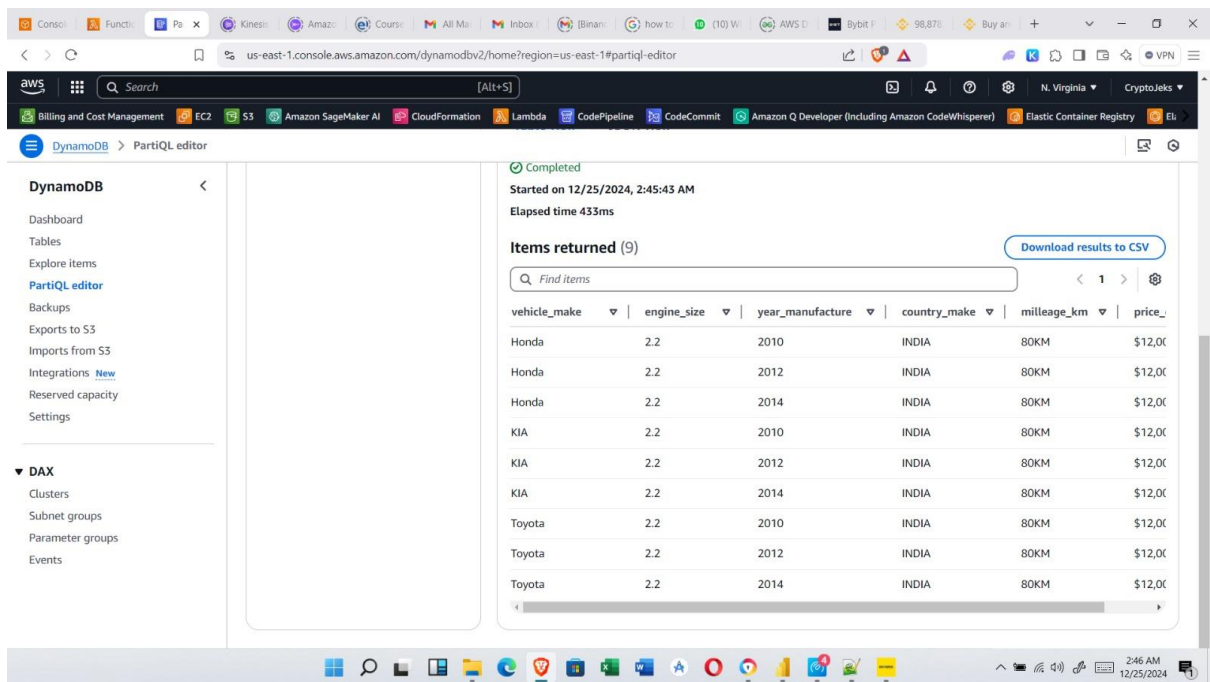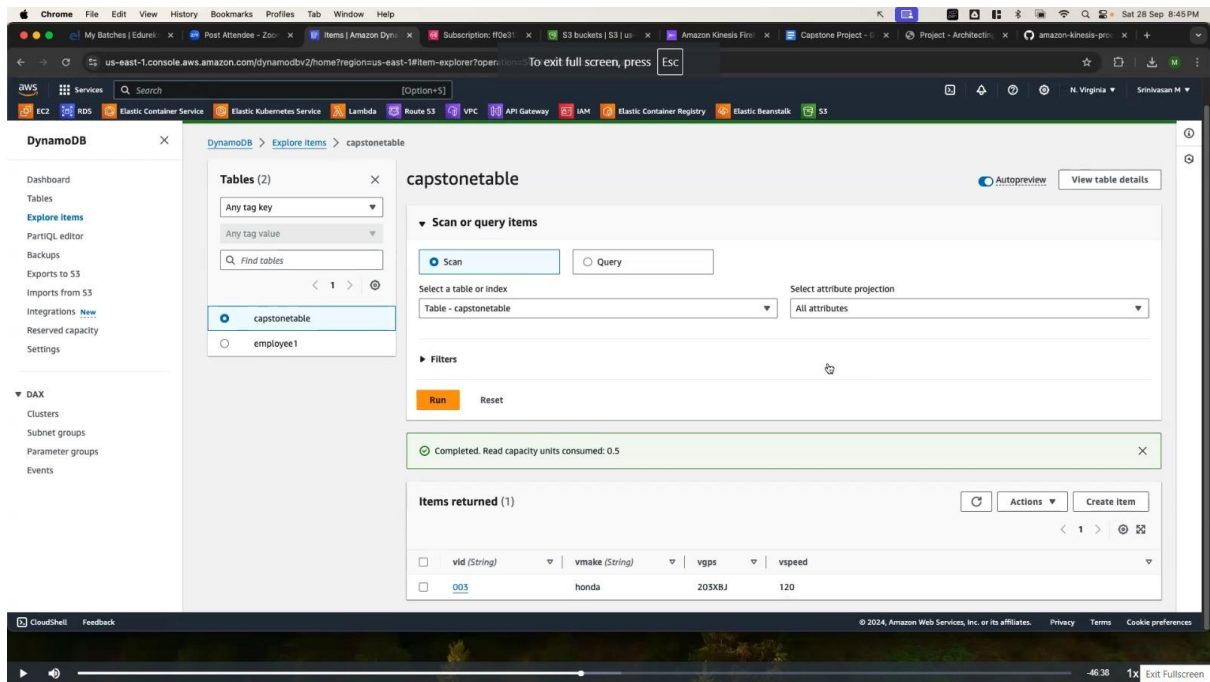| | Name ▲ | Status ▽ | Type ▽ | ARN ▽ | Description ▽ |
| --- | --- | --- | --- | --- | --- |
| ☐ | samson_event | ⊘ Enabled | Standard | arn:aws:events:us-east-1:905418056715:rule/samson_event | event from S3 to SNS Topic |
| ☐ | samson-capstone1 | ⊘ Enabled | Standard | arn:aws:events:us-east-1:905418056715:rule/samson-capstone1 | my_capstone1 |

# DynamoDB Creation

DynamoDB is a fully managed, key-value, and document database that delivers single-digit-millisecond performance at any scale. In this case, my event, stored in my S3 bucket, is pushed on DynamoDB using an SNS Topic via a Lambda function. Below is a screenshot of my DynamoDB creation.

To exit full screen, press `Esc`

**DynamoDB** ✕

DynamoDB > Explore items > capstonetable

## capstonetable

Autopreview   **View table details**

Dashboard
Tables
**Explore items**
PartiQL editor
Backups
Exports to S3
Imports from S3
Integrations  *New*
Reserved capacity
Settings

▼ DAX
Clusters
Subnet groups
Parameter groups
Events

**Tables** (2)   ✕

Any tag key ▼

Any tag value ▼

🔍 Find tables

‹ 1 ›   ⚙

◉ capstonetable
○ employee1

▼ **Scan or query items**

◉ Scan       ○ Query

Select a table or index
Table - capstonetable ▼

Select attribute projection
All attributes ▼

▶ Filters

**Run**   Reset

✓ Completed. Read capacity units consumed: 0.5   ✕

**Items returned** (1)   ↻   Actions ▼   **Create item**

‹ 1 ›   ⚙ ⤢

| | vid *(String)* ▽ | vmake *(String)* ▽ | vgps ▽ | vspeed ▽ |
|---|---|---|---|---|
| ☐ | 003 | honda | 203XBJ | 120 |

CloudShell   Feedback      © 2024, Amazon Web Services, Inc. or its affiliates.   Privacy   Terms   Cookie preferences

▶ 🔊 ———————————————————————————————   -46:38   1x   Exit Fullscreen

---

**DynamoDB** > PartiQL editor

✓ Completed
Started on 12/25/2024, 2:45:43 AM
Elapsed time 433ms

Dashboard
Tables
Explore items
**PartiQL editor**
Backups
Exports to S3
Imports from S3
Integrations  *New*
Reserved capacity
Settings

▼ DAX
Clusters
Subnet groups
Parameter groups
Events

**Items returned** (9)   **Download results to CSV**

🔍 Find items   ‹ 1 ›   ⚙

| vehicle_make ▽ | engine_size ▽ | year_manufacture ▽ | country_make | milleage_km ▽ | price_ |
|---|---|---|---|---|---|
| Honda | 2.2 | 2010 | INDIA | 80KM | $12,00 |
| Honda | 2.2 | 2012 | INDIA | 80KM | $12,00 |
| Honda | 2.2 | 2014 | INDIA | 80KM | $12,00 |
| KIA | 2.2 | 2010 | INDIA | 80KM | $12,00 |
| KIA | 2.2 | 2012 | INDIA | 80KM | $12,00 |
| KIA | 2.2 | 2014 | INDIA | 80KM | $12,00 |
| Toyota | 2.2 | 2010 | INDIA | 80KM | $12,00 |
| Toyota | 2.2 | 2012 | INDIA | 80KM | $12,00 |
| Toyota | 2.2 | 2014 | INDIA | 80KM | $12,00 |

2:46 AM
12/25/2024

My S3 Event published on my DynamoDB in a Tabular Form

## Amazon Athena

Amazon Athena is an interactive query service that makes it easy to analyze data in Amazon S3 and other federated data sources using standard SQL. In this case, I used Amazon Athena to analyze the above data published on my DynamoDB.

## API Gateway Section

Amazon API Gateway helps developers to create and manage APIs to back-end systems running on Amazon EC2, AWS Lambda, or any publicly addressable web service. With Amazon API Gateway, you can generate custom client SDKs for your APIs, to connect your back-end systems to mobile, web, and server applications or services. In this case, I configured Amazon API Gateway to enable external users to have access to the real time data published on DynamoDB via Lambda function. Users can access this data by first, validating their authentication using Amazon Cognito after which, they can connect to the API Gateway either by using Data Lake Console or Data Lake CLI.