# Case Study

## Leveraging AWS Lambda, S3, and DynamoDB for Serverless Data Processing

**Name: Prajyot Shinde**

**Roll No: 57   (D15A)**

## Introduction:

In this case study, we explore the benefits, challenges, and practical applications of building a serverless architecture using **AWS Lambda**, **Amazon S3**, and **Amazon DynamoDB**. These services are a part of the AWS ecosystem and allow developers to create scalable, low-maintenance systems that automatically adjust to varying loads and are event-driven by design.

The primary use case for this study is to process files (JSON) uploaded to **S3**, trigger **AWS Lambda** to extract relevant data from the JSON files, and store the data in a **DynamoDB** table. This approach will be compared with traditional methods that use databases such as **MongoDB** and **PostgreSQL**, illustrating the pros, cons, and alternative solutions.

## Background:

Traditionally, applications rely on servers running continuously to handle requests, perform business logic, and manage storage through databases. This can become costly and complex as the infrastructure scales. Moreover, applications that need to process data on demand often require careful capacity planning and load management.

The evolution of cloud computing and serverless technologies has introduced new ways of handling backend workloads. **AWS Lambda** removes the need to manage servers, offering on-demand code execution. Combined with **Amazon S3** for file storage and **DynamoDB** for NoSQL database services, developers can create serverless, event-driven applications that dynamically scale with user demands.

In this case study, we examine an example scenario to demonstrate how Lambda, S3, and DynamoDB are integrated and why this combination is ideal for certain applications.
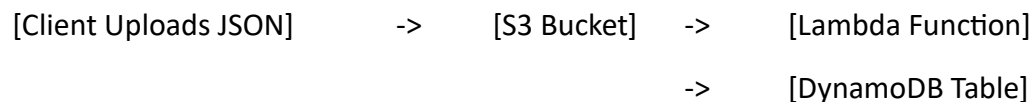
## Scenario Overview

The scenario involves developing an event-driven system where users upload JSON files containing user information (such as userID, timestamp, etc.) to an S3 bucket. Upon the upload, an AWS Lambda function is triggered to:

1. **Read the file** from S3.

2. **Parse the JSON data** to extract specific fields like userID and timestamp.

3. **Store this data** in a DynamoDB table.

We will look at the steps involved, the benefits of using this architecture, and the limitations compared to traditional systems.

## Architecture Diagram:

[Client Uploads JSON]        ->       [S3 Bucket]    ->      [Lambda Function]

                                                              ->      [DynamoDB Table]

## Step-by-Step Breakdown:

**Step 1: Create a DynamoDB Table**

1. **Log in to the AWS Management Console**:

   o Go to the [AWS Management Console](AWS%20Management%20Console).

2. **Navigate to DynamoDB**:

   o In the services menu, search for **DynamoDB** and click on it.

3. **Create a new table**:

   o Click on **Create table**.

   o **Table name**: Choose a name for your table, e.g., UserData.

   o **Partition key**: Define a partition key, e.g., userID (String type).

   o Optionally, you can define a sort key (e.g., timestamp), but for simplicity, we will only use a partition key here.

   o Leave the rest of the settings as default, but feel free to customize them based on your needs.

   o Click on **Create**.

4. **Define Attributes**:

   o Add relevant attributes based on your JSON structure. For this case, you may need attributes like userID, timestamp, and any other fields you plan to store.



**Step 2: Create an S3 Bucket**

1. **Navigate to S3**:

   o In the AWS Management Console, search for **S3** and click on it.

2. **Create a new bucket**:

   o Click on **Create bucket**.

   o **Bucket name**: Choose a unique name (e.g., my-json-bucket).

- o **Region**: Select a region where you want to create the bucket.

- o Leave other options as default or customize as necessary.

- o Click **Create bucket**.





**Steps to Upload Your Lambda Function to S3**

1. **Upload the ZIP File to S3**:

   - o **Open the Amazon S3 Console**: Go to the Amazon S3 console.

   - o **Create a New Bucket (if needed)**: If you don't already have a bucket, you can create one:

- Click on **Create bucket**.

- Provide a unique bucket name and configure settings as needed.

  o **Upload the ZIP File**:

- Click on the bucket you want to use.

- Click on **Upload**.

- Click on **Add files** and select your ZIP file.

- Click on **Upload**.

2. **Get the S3 Object URL**: After uploading, find the object in your S3 bucket, and copy its URL. It should look something like this:

**Step 3: Write the Lambda Function**

1. **Navigate to Lambda**:

   o In the AWS Management Console, search for **Lambda** and click on it.

2. **Create a new Lambda function**:

   o Click on **Create function**.

   o **Function name**: Give it a name (e.g., ProcessJsonFunction).

   o **Runtime**: Select **Node.js 18.x** (or the latest version).

   o **Permissions**: Choose "Create a new role with basic Lambda permissions" or use an existing role if you have one.



3. **Add the necessary permissions**:

   o Go to the **Permissions** tab of your Lambda function.

   o Click on the role linked to your Lambda function to open the IAM role settings.

   o Click on **Add permissions** > **Attach policies**.

   o Search for and attach the following policies:

      ▪ **AmazonDynamoDBFullAccess** (or create a custom policy with limited access to your table).

      ▪ **AmazonS3ReadOnlyAccess** (or create a custom policy to read from your specific bucket).

**Go to the Function Code Section**:

- In the **Code** section, look for the option that allows you to upload from S3.

**Enter the S3 URL**:

- Paste the S3 object URL you copied earlier into the field for the S3 location.

**Save or Update the Function**: After entering the S3 URL, click on **Save** or **Deploy**.

**Step 4: Configure the S3 Trigger**

1.  **Set up the trigger**:

    o   Go to the **Configuration** tab of your Lambda function.

    o   Click on **Triggers** and then **Add trigger**.

    o   Select **S3** from the dropdown.

    o   Choose the bucket you created (my-json-bucket).

    o   For the event type, select **All object create events**.

    o   Click on **Add**

**Step 5: Upload a Sample JSON File to S3**

1. **Navigate to your S3 bucket**:

   o Go back to your S3 bucket (my-json-bucket).

2. **Upload a sample JSON file**:

   o Click on **Upload**.

   o Choose **Add files**, and select a sample JSON file with the following structure:

```
3. [
4.   {
5.     "userID": "user123",
6.     "timestamp": "2024-10-18T10:30:00Z",
7.     "name": "Prajyot Shinde",
8.     "email": "prajyot.shinde@example.com",
9.     "age": 22
10.     },
11.     {
12.       "userID": "user456",
13.       "timestamp": "2024-10-18T11:00:00Z",
14.       "name": "Jane Smith",
15.       "email": "jane.smith@example.com",
16.       "age": 25
17.     },
18.     {
19.       "userID": "user789",
20.       "timestamp": "2024-10-18T11:30:00Z",
21.       "name": "Mike Johnson",
22.       "email": "mike.johnson@example.com",
23.       "age": 40
24.     }
25.   ]
26.
```

**Step 6: Verify Data in DynamoDB**

1. **Navigate to DynamoDB**:

   o In the AWS Management Console, go back to DynamoDB.

2. **View the items in your table**:

   o Click on the table you created (UserData-prajyot).

   o Go to the **Items** tab to see if the data from the JSON file has been successfully written to your DynamoDB table.

## Benefits of AWS Lambda, S3, and DynamoDB Architecture

**AWS Lambda:**

- **Serverless**: There is no need to manage or provision any infrastructure. AWS handles all the scaling, maintenance, and updates.

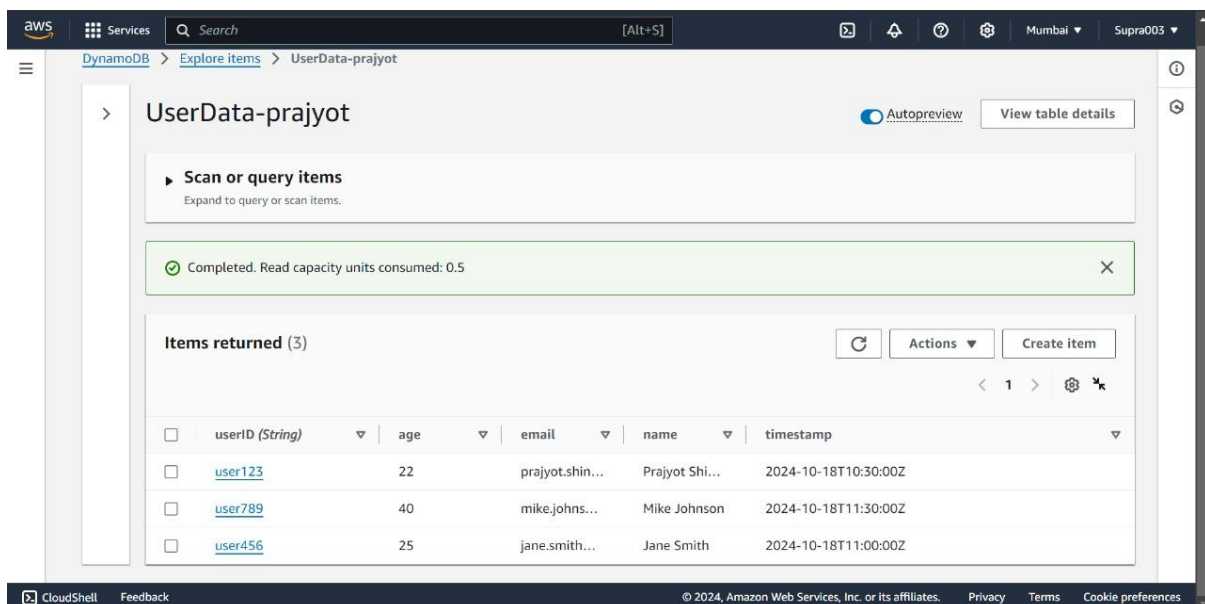- **Event-Driven**: Lambda functions can automatically trigger based on predefined events, such as an S3 file upload or a DynamoDB stream.

- **Cost-Efficient**: You only pay for the compute time when the Lambda function is executed, making it highly cost-effective for applications with sporadic usage.

**Amazon S3:**

- **Scalability**: S3 scales effortlessly as the amount of data increases, without any intervention required from the user.

- **Durability**: S3 offers 99.999999999% durability (11 nines), ensuring that data is protected from data loss.

- **Integration with Lambda**: S3 natively integrates with Lambda, which allows files to trigger backend processing with ease.

**Amazon DynamoDB:**

- **Low Latency**: DynamoDB provides single-digit millisecond response times, which makes it suitable for applications requiring fast access to data.

- **Scalability**: Like S3, DynamoDB automatically scales with the load, so it can handle a high volume of data without manual intervention.

- **Flexible NoSQL**: DynamoDB is schema-less, which allows you to store and query JSON-like documents flexibly.

## Limitations of AWS Lambda, S3, and DynamoDB

**AWS Lambda:**

- **Cold Starts**: When a Lambda function has not been invoked for some time, there may be a delay (cold start) when it's triggered again.

- **Execution Limits**: Lambda has a maximum execution time of 15 minutes, so long-running processes are not ideal.

- **Stateless**: Lambda functions are stateless, meaning you must use external services like S3 or DynamoDB to store state between executions.

**Amazon S3:**

- **Eventual Consistency**: S3 follows an eventual consistency model for some operations, such as overwriting an object or listing buckets.

- **Not a File System**: While S3 excels in object storage, it's not suited for high-frequency, small data transactions like a traditional file system.

**Amazon DynamoDB:**

- **Complex Queries**: DynamoDB is optimized for simple key-value lookups and basic queries. Complex queries with joins or aggregates are not supported natively.

- **Pricing Model**: Costs can increase for applications with high write and read throughput if not carefully optimized.

## Use Cases of AWS Lambda, S3, and DynamoDB

- **Serverless Web Applications**: Frontend web applications can store files in S3, process them with Lambda, and store metadata in DynamoDB, creating a completely serverless architecture.

- **IoT Systems**: IoT devices can send data to S3 or DynamoDB through Lambda, providing scalable data processing pipelines.

- **Real-time Analytics**: Event-driven systems like Lambda can be used to process and analyze data in real-time, triggered by uploads or changes in data streams.

## Alternatives to AWS Lambda, S3, and DynamoDB

**Serverless Compute Alternatives:**

- **Google Cloud Functions**: Google Cloud's serverless offering similar to AWS Lambda.

- **Azure Functions**: Microsoft's serverless compute option with seamless integration into the Azure ecosystem.

**Object Storage Alternatives:**

- **Google Cloud Storage**: Similar to Amazon S3 but on Google Cloud, with comparable scalability and durability.

- **Azure Blob Storage**: Microsoft's cloud object storage service.

**NoSQL Alternatives:**

- **MongoDB Atlas**: A managed version of MongoDB offering rich query capabilities.

- **Cassandra**: An open-source NoSQL database designed for distributed data storage.

## Problems Faced and Solutions

1. **Error: "Cannot find module 'index'"**

   o **Problem**: This error occurred due to AWS Lambda not recognizing the index.mjs file when attempting to load the handler function.

   o **Solution**: The issue was resolved by ensuring that the **handler function** was correctly set in the AWS Lambda configuration and the **filename** was properly named index.mjs in both the ZIP file and AWS Lambda's settings.

2. **Packaging Node Modules**

   o **Problem**: I initially faced issues related to missing the node_modules directory and other dependencies when deploying the Lambda function.

   o **Solution**: By ensuring the **aws-sdk** was installed and packaging only the necessary files (index.mjs and package.json) into the ZIP file, I overcame this issue. The **aws-sdk** is pre-installed in AWS Lambda, so no need to include it in the package.

3. **Permission Errors in S3**

   o **Problem**: While trying to read the JSON file from S3, permission-related errors occurred due to missing access rights.

   o **Solution**: I resolved this by assigning the **proper permissions** to the Lambda function through an **IAM role** that allowed it to access the S3 bucket and DynamoDB.

4. **Data Processing Logic**

   o **Problem**: Initially, I was able to process only basic user data fields (userID and timestamp), but more fields and multiple users were needed.

   o **Solution**: I enhanced the Lambda function logic to handle **additional fields** and multiple users by updating the JSON schema, dynamically parsing the incoming data, and writing each user's data to DynamoDB.

## Conclusion

AWS Lambda, S3, and DynamoDB together provide a powerful, serverless infrastructure that is highly scalable, cost-efficient, and ideal for event-driven architectures. They allow for a high degree of automation and are well-suited for modern applications that need to scale dynamically. However, their limitations in execution time, cold starts, and query complexity make them less suitable for stateful applications or those with highly complex queries.

For applications requiring fast, scalable, event-driven processing, this architecture is a robust choice. However, if your application demands advanced querying, transactions, or stateful workflows, traditional databases such as **MongoDB** and **PostgreSQL** or other alternatives might be more appropriate.