# Content Page

## 1.1 About Report

Project Name: Analysis Data Pipeline for Global Fisheries
Module: Data Engineering in the Cloud
By: Javen Lai Le Yu (2202934B)
Class: P03

## 1.2 Background

I am tasked to create an infrastructure to host fishing data so that Data Analysts can create reports about fishing impacts in the open seas. My solution is a batch Data Pipeline using AWS cloud services that provides analysts with the data they require to perform analytics.
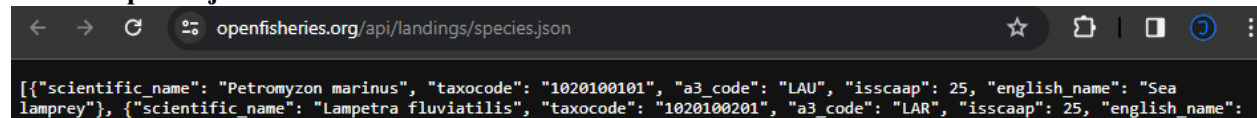
## 1.3 Problem Statement

Data analysts need to analyze catch performance for species and countries, for the years from 1950 to 2018.
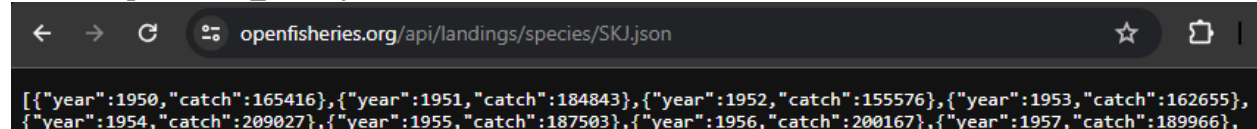
## 1.4 Data Source

The datasets about fishes will be pulled from OpenFisheries.org's publicly available RESTful API, utilizing the following endpoints:

1. **/species.json**



[{"scientific_name": "Petromyzon marinus", "taxocode": "1020100101", "a3_code": "LAU", "isscaap": 25, "english_name": "Sea lamprey"}, {"scientific_name": "Lampetra fluviatilis", "taxocode": "1020100201", "a3_code": "LAR", "isscaap": 25, "english_name":

Provides the list of all species of fish that exist and their respective information.

2. **/species/{a3_code}.json**



[{"year":1950,"catch":165416},{"year":1951,"catch":184843},{"year":1952,"catch":155576},{"year":1953,"catch":162655}, {"year":1954,"catch":209027},{"year":1955,"catch":187503},{"year":1956,"catch":200167},{"year":1957,"catch":189966},
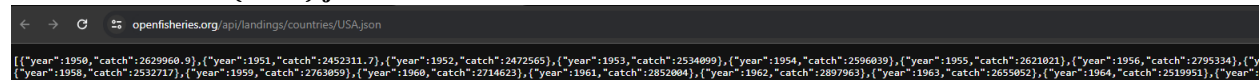
Provides how many tonnes of a fish species were caught each year, accessed using the a3_code of the fish species.

3. **/countries.json**



[{"country": "Afghanistan", "iso3c": "AFG"}, {"country": "Albania", "iso3c": "ALB"}, {"country": "Algeria", "iso3c": "DZA"}, {"country": "American Samoa", "iso3c": "ASM"}, {"country": "Andorra", "iso3c": "AND"}, {"country": "Angola", "iso3c": "AGO"}, {"country": "Anguilla", "iso3c": "AIA"}, {"country": "Antigua and Barbuda", "iso3c": "ATG"}, {"country": "Argentina", "iso3c": "ARG"}, {"country": "Armenia", "iso3c": "ARM"}, {"country": "Aruba", "iso3c": "ABW"}, {"country": "Australia", "iso3c": "AUS"}, {"country": "Austria", "iso3c": "AUT"}, {"country": "Azerbaijan", "iso3c": "AZE"}, {"country": "Bahamas", "iso3c": "BHS"}, {"country": "Bahrain", "iso3c": "BHR"}, {"country": "Bangladesh", "iso3c": "BGD"}, {"country": "Barbados", "iso3c": "BRB"}, {"country": "Belarus", "iso3c": "BLR"},

Provides the list of countries available and their iso3c.

4. **countries/{iso3c}.json**



[{"year":1950,"catch":2629960.9},{"year":1951,"catch":2452311.7},{"year":1952,"catch":2472565},{"year":1953,"catch":2534099},{"year":1954,"catch":2596039},{"year":1955,"catch":2621021},{"year":1956,"catch":2795334},{"y {"year":1958,"catch":2532717},{"year":1959,"catch":2763059},{"year":1960,"catch":2714623},{"year":1961,"catch":2852004},{"year":1962,"catch":2897963},{"year":1963,"catch":2655052},{"year":1964,"catch":2519951},{"year":
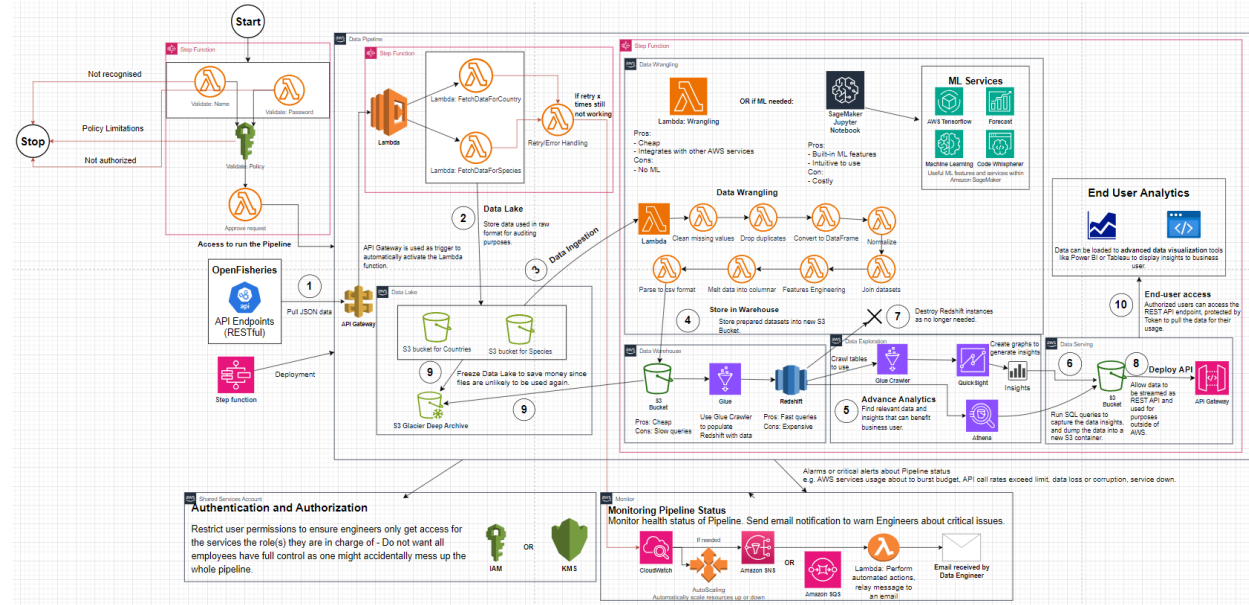
Provides how many tonnes of fish were caught each year for the specified country.

## 1.5 Objectives of the project

1. Implement a viable prototype of the designed Data Pipeline using AWS services as Proof of Concept (POC).
2. Perform basic Business Intelligence by creating visualizations to answer the problem statement using the prototype to prove the pipeline's usability and judge its feasibility for actual deployment using larger datasets.

## 2. Data Pipeline Architecture

The Data Pipeline architecture has been designed under the Group Project section.



## Limitations of Implementation:

1. Restricted IAM access as I have to work with a controlled AWS environment where I cannot create new roles. Hence, unable to implement authentication and authorization measures.
2. Budget of $100.

## Modifications to Pipeline design:

- Unable to create users or roles to emulate several data engineers due to IAM restrictions.
- I will have a central S3 Bucket as a data lake to store species and countries, separated by folders.
- Use Glue Catalog as the database instead of Redshift due to tight budget.
- Skip Steps 8 and 10. Do not need to release data as API when I'm the only analyst using the data.

## 3. Creating Pipeline

3.1 Preparing S3 Buckets

| | | | | |
|---|---|---|---|---|
| ○ | openfisheries | US East (N. Virginia) us-east-1 | Bucket and objects not public | January 24, 2024, 15:33:34 (UTC+08:00) |
| ○ | openfisheries-lake | US East (N. Virginia) us-east-1 | Bucket and objects not public | February 7, 2024, 10:18:27 (UTC+08:00) |
| ○ | openfisheries-warehouse | US East (N. Virginia) us-east-1 | Bucket and objects not public | February 8, 2024, 21:32:36 (UTC+08:00) |

Created all 3 S3 Buckets that will be used. Default configurations were used. However, for openfisheries, I enabled versioning since this bucket stores the main data files that will be used for analysis, hence keeping historical versions would be good to track deletes or changes.

**Bucket Versioning**

○ Disable

● Enable

## 3.2 Data Lake

Goal: Pull 30 different datasets from OpenFisheries REST API and dump them into S3 Data Lake.

3.2.1 Configurations of Lambda function, fetchAPI.py:



| | | | Runtime settings Info | | | Edit | Edit runtime management configuration |
|---|---|---|---|---|---|---|---|

**Runtime**
Python 3.12

**Handler** Info
fetchAPI.lambda_handler

**Architecture** Info
x86_64

▶ Runtime management configuration

**Layers** Info  [Edit] [Add a layer]

| Merge order | Name | Layer version | Compatible runtimes | Compatible architectures | Version ARN |
|---|---|---|---|---|---|
| 1 | AWSSDKPandas-Python312 | 1 | python3.12 | x86_64 | arn:aws:lambda:us-east-1:336392948345:layer:AWSSDKPandas-Python312:1 |

| 3.2.2 | Modified Timeout to 1 min to ensure sufficient timing is given for code to run before timing out. | |
|---|---|---|
| 3.2.3 | Execute Function under LabRole permissions. | |

Timeout: 1 min 0 sec

Execution role
Choose a role that defines the permissions of your function. To create a custom role, go to the IAM console.
◉ Use an existing role
○ Create a new role from AWS policy templates

### 3.2.4 **Data Ingestion:** Lambda code to pull data from API Endpoints

No retry mechanism because unreachable endpoints do not exist. For error-tolerance, the code moves on to the next country/species in the list. It runs until 1 min timeout OR if 30 (15-15) datasets are loaded.



After running Python code that **pulls the first 15 species and 15 countries** in the lists available, an interesting issue was found: not all species specified on the species list have an accessible endpoint.

### 3.2.5 S3 Bucket results:

| | Name | Type | Last modified | Size | Storage class |
|---|---|---|---|---|---|
| ☐ | countries.json | json | February 8, 2024, 21:49:14 (UTC+08:00) | 10.2 KB | Standard |
| ☐ | countries/ | Folder | - | - | - |
| ☐ | species.json | json | February 8, 2024, 21:49:13 (UTC+08:00) | 1.5 MB | Standard |
| ☐ | species/ | Folder | - | - | - |

### 3.2.6 Proof that data were successfully dumped into Data Lake

Amazon S3 > Buckets > openfisheries-lake > countries/

**countries/**  [Copy S3 URI]

[Objects] [Properties]

**Objects (15)** Info  [↻] [Copy S3 URI] [Copy URL] [Download] [Open] [Delete] [Actions ▾] [Create folder] [Upload]

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

Q Find objects by prefix

| | Name | Type | Last modified | Size | Storage class |
|---|---|---|---|---|---|
| ☐ | ABW.json | json | February 8, 2024, 21:49:44 (UTC+08:00) | 2.0 KB | Standard |
| ☐ | AFG.json | json | February 8, 2024, 21:49:38 (UTC+08:00) | 2.0 KB | Standard |
| ☐ | AGO.json | json | February 8, 2024, 21:49:41 (UTC+08:00) | 2.2 KB | Standard |

## 3.3 Data Wrangling

Extract files from Data Lake into Lambda Function that performs cleaning and transformation to prepare the data the be stored in Data Warehouse schema and usable for analytics. The same configurations for the Lambda function are used (Runtime + Layer to enable pandas libraries to be used).

| Description | Code |
|---|---|
| Added a logger to log errors that occurred in data cleaning.<br><br>This function appends a new column to store the file's name as a value (Species/Country name). Then, it unpivots the data from separate files into a single dataset by stacking them in rows and transforming the data into columnar format (to be used for analysis). | ```python
import boto3
import pandas as pd
from io import import StringIO
import json
import logging

# Configure logging
logger = logging.getLogger()
logger.setLevel(logging.INFO)

def process_s3_data(input_bucket, input_folder, output_column):
    s3 = boto3.client('s3')
    df = pd.DataFrame()

    try:
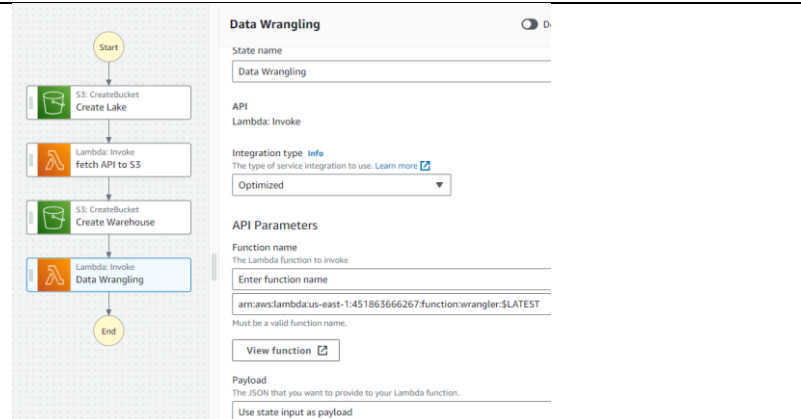        response = s3.list_objects(Bucket=input_bucket, Prefix=f'{input_folder}/')

        for obj in response.get('Contents', []):
            file_key = obj['Key']
            KEY = file_key.split('/')[-1].replace('.json', '')
            response = s3.get_object(Bucket=input_bucket, Key=file_key)
            json_data = response['Body'].read().decode('utf-8')
            temp_df = pd.read_json(StringIO(json_data))
            temp_df[output_column] = KEY
            df = pd.concat([df, temp_df], ignore_index=True)

    except Exception as e:
        logger.error(f"Error processing S3 data: {str(e)}")

    return df
``` |
| Based on OpenFisheries API GitHub documentation, the data in the API is already prepared as it was originally taken from another data source, suggesting that the data quality should be relatively good and clean. | **API data sources**<br><br>**FAO Global fisheries capture landings**<br><br>Data for the OpenFisheries landings API is sourced from the UN Food and Agriculture Organization. FishstatJ Capture dataset is exported to csv in /data/ with the following notes:<br><br>• Country column exported as ISO 3-letter code<br>• Species exported as ASFIS code<br>• Symbols exported as a separate column (the "S" column in the CSV) |
| Basic data cleaning practices still applied despite trust in the API just as precautionary measures.<br><br>This function performs basic data cleaning like dropping missing values and duplicated rows. | ```python
def clean_data(df):
    try:
        df.dropna(inplace=True)
        df.drop_duplicates(inplace=True)
    except Exception as e:
        logger.error(f"Error cleaning data: {str(e)}")

    return df
``` |
| 1st Function stores each file as CSV into a file dir using its name before storing the file because each file dir can only contain 1 file for Glue Crawler to Athena to work.<br><br>2nd Function to parse countries and species lists into CSV before storing them in S3. | ```python
def to_warehouse(df, folder_name, output_bucket):
    try:
        s3 = boto3.client('s3')
        csv_data = df.to_csv(index=False)
        csv_key = f'{folder_name}.csv'
        s3.put_object(Bucket=output_bucket, Key=f"{csv_key}/{csv_key}", Body=csv_data)
        logger.info(f'CSV file created successfully in bucket: {output_bucket}')
    except Exception as e:
        logger.error(f"Error uploading CSV to warehouse: {str(e)}")

def parse_csv(bucket, key, destination_bucket):
    s3 = boto3.client('s3')
    try:
        csv_obj = s3.get_object(Bucket=bucket, Key=f'{key}.json')
        json_data = csv_obj['Body'].read().decode('utf-8')
        df = pd.read_json(StringIO(json_data), orient='records')
        csv_data = df.to_csv(index=False)
        s3.put_object(Bucket=destination_bucket, Key=f'dim_{key}/dim_{key}.csv', Body=csv_data)

    except Exception as e:
        logger.error(f"Error parsing and uploading CSV: {str(e)}")
``` |
| Handler is the function that gets run when this Lambda function is run. Hence, this acts as a main() that activates the data wrangling process.<br><br>The countries and species list doesn't need cleaning as it's used as a look-up table (dictionary to find endpoint) and not actual data rows. | ```python
def lambda_handler(event, context):
    LAKE = 'openfisheries-lake'
    WAREHOUSE = 'openfisheries-warehouse'

    parse_csv(LAKE, 'species', WAREHOUSE)
    parse_csv(LAKE, 'countries', WAREHOUSE)

    try:
        # SPECIES
        species = process_s3_data(LAKE, 'species', 'a3_code')
        species = clean_data(species)
        to_warehouse(species, 'fact_species', WAREHOUSE)

        # COUNTRY
        country = process_s3_data(LAKE, 'countries', 'iso3c')
        country = clean_data(country)
        to_warehouse(country, 'fact_countries', WAREHOUSE)

        return {
            'statusCode': 200,
            'body': json.dumps('Processing completed successfully!')
        }
    except Exception as e:
        logger.error(f"Error in lambda_handler: {str(e)}")
        return {
            'statusCode': 500,
            'body': json.dumps('Error during processing. Check logs for details.')
        }
``` |

## 3.4  Automating using Step Functions

Goal: Create a function that runs all pre-processing work with 1 click to ensure easy Pipeline reusability.
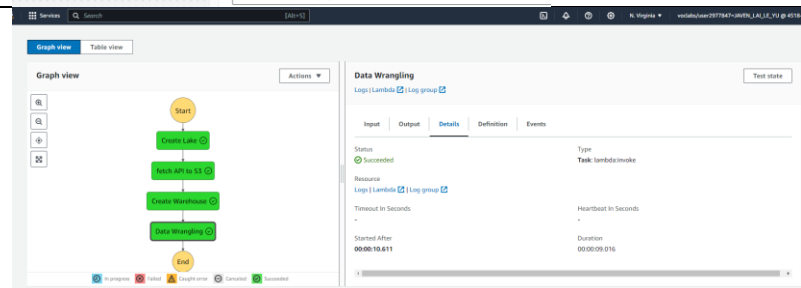
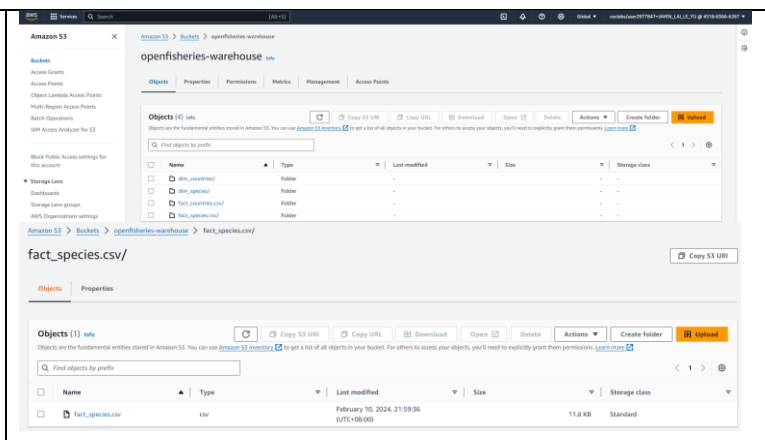| | |
|---|---|
| Lake = openfisheries-lake,<br>Fetch API = fetchAPI,<br>Warehouse = openfisheries-warehouse,<br>Data Wrangling = wrangler<br><br>This workflow was manually created using drag and drop. |  |
| This workflow was able to successfully run without encountering any errors.<br><br>This is validated as I cleared the S3 Buckets, ran this Step function, and the S3 buckets were populated with data. |  |
| This is the list of events triggered by the Step function.<br><br>As seen, all the steps were successfully run and the functions were executed accordingly. |  |

## 3.5  Data Warehouse

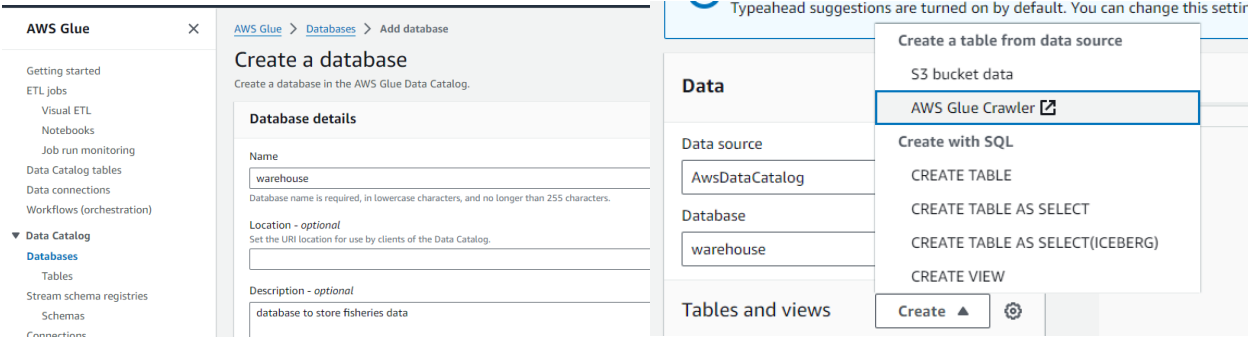| | |
|---|---|
| All the prepared datasets for the data warehouse are stored in a S3 Bucket.<br><br>The intention of doing so is to ensure dataset files are easily and conveniently accessible when needed, as the relational database that stores these data can be destroyed to cut costs.<br><br>Due to a tight budget, Glue Catalog is used as a database to store the data to be queryable using Athena. |  |

### 3.5.1 Creating Database at Glue Data Catalog



### 3.5.2 Glue Crawler

To effortlessly load the datasets from the S3 warehouse into the database, a Crawler is used to infer and automate the creation of tables based on the dataset files and load the data into the tables.

**Some additional configurations:**
Firstly, I declared that each S3 file directory path has its own schema.

Next, I set the Table level to 2 to ensure that the files nested within the respective sub-folders get crawled:

| | Name | ▲ | Type |
|---|---|---|---|
| ☐ | 📄 countries.json | | json |
| ☐ | 📄 countries/ | | Folder |
| ☐ | 📄 species.json | | json |
| ☐ | 📄 species/ | | Folder |



### 3.6 Exploratory Data Analysis (EDA) + Preparing Data for Analytics

#### 3.6.0 Configurations for Athena

Configured Athena to dump results of all queries run into the openfisheries bucket meant for storing the insights from SQL queries.

When an SQL query is run from Editor console, the results will be logged and sent to 'openfisheries' S3 Bucket with versioning.



#### 3.6.1 Species:

To prepare the data to be used for EDA, I created a View which essentially acts as a table that only has the relevant columns I need to use.

```
1  CREATE VIEW species AS (
2  SELECT d.scientific_name, d.english_name, d.isscaap, f.*
3  FROM dim_species d, fact_species_csv f
4  WHERE d.a3_code = f.a3_code
5  )
```

**SELECT \* FROM species** to check the data inside the View table. There are 1035 observations from 2018 to 1950 for the 15 selected species.

Found that the iso3c for certain countries was wrong. Then, I noticed that the country name was duplicated for these records. Another issue with the countries list dataset was that the column headers were not configured properly.

However, these issues should not be a concern as they do not affect the main data and can be easily fixed with SQL.

### 3.6.2 Countries:

The 15 files I sampled can NEVER contain a country with faulty iso3c as an endpoint with wrong iso3c does not exist.

A possible fix to this issue would be to source online for a flawless country and iso3c list with no data issues instead of using the OpenFisheries countries list if all countries' fishery data needs to be extracted and used for analysis.



Similar to species, I created a View for countries that only contains the relevant columns needed for analysis. Additionally, column headers are fixed to be an accurate representation of what the values that the column holds.



**Code used to create View:**

```
CREATE VIEW countries AS (
SELECT d.col0 AS name, f.year, f.catch
FROM dim_countries d, fact_countries_csv f
WHERE d.col1 = f.iso3c
)
```

These 2 main datasets are the final product of this Pipeline, and are stored in the openfisheries S3 Bucket for indefinite usage (until Openfisheries updates again).



**How openfisheries S3 bucket looks:**



The CSV files are the main products.

To show that the version control for 'openfisheries' S3 bucket works, the photo shows that all the previous files that were deleted are still recorded and traceable in the bucket when 'Show Versions' is checked.
To delete the trace of the file permanently:





Deleted files were unnecessary queries
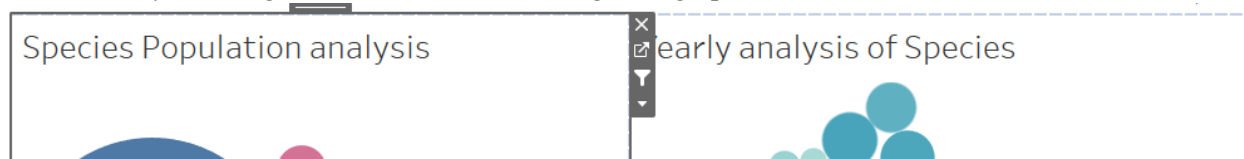
## 3.7 Advanced Analytics using Tableau

The files were downloaded from the insights S3 Bucket and loaded onto Tableau to conduct professional Data Analytics to report any insights or findings regarding fishes… and to answer the problem statement.

### 3.7.1 Some steps taken to build visualization

1. Filter Catch=0 *for the KPI graph* as 0 is almost the same as having no fish caught that year since the sum is less than 1 tonne. These rows are retained in data because a country that has nearly 0 catches for a year is still considered an insight.

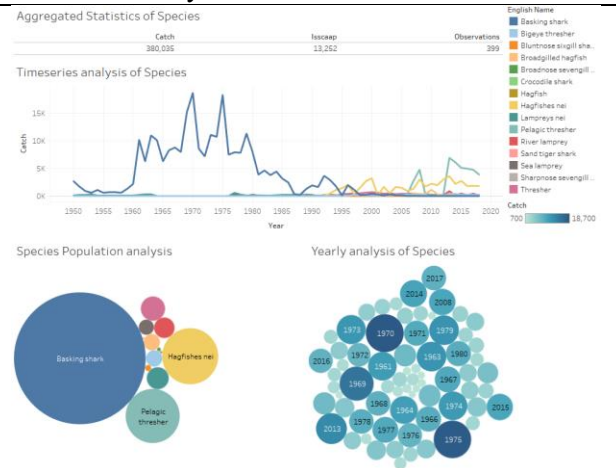| Filters | species (2) | | | Catch |
|---|---|---|---|---|
| Catch | | | | 1        18700 |
| Measure Names | Catch | Isscaap | Observations | |
| | 380,035 | 13,252 | 399 | |

2. Created a Dashboard for both Species and Countries to conduct analysis. Designed the dashboard to be interactive by allowing the data to be filtered using each graph in the dashboard.
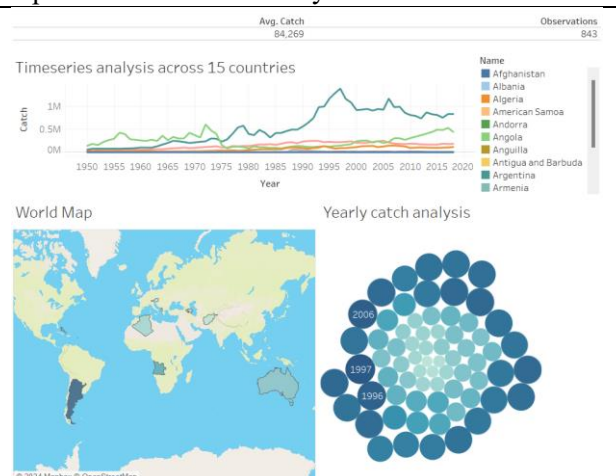
Species Population analysis   early analysis of Species

### 3.7.2 Species Dashboard to analyze Species performance across the year.

1. The basking shark was an extremely popular fish during 1960-1980, but the catch quantity drastically fell after. This suggests that this species is either going extinct as it was overfished or that people lost interest in catching this species.

2. The most popular fish in the modern 21$^{st}$ century is the Bigeye thresher, which especially gained popularity from 2012 onwards.

3. The Hagfish Nei was only caught in 1993, suggesting that its existence was found in 1993.

4. Most of the fish were caught during 1975, 1970, and 1969.

### 3.7.3 Countries Dashboard to analyze countries' catch performance across the year.

1. Argentina is the heaviest fishing country out of the 15 countries ever since 1975 onwards.

2. Angola's fishing intensity dropped sharply after 1975, but they picked up their pace again from 2010 onwards.

3. Countries surrounding the open ocean waters tend to have higher catches in tonnes.

4. There seem to be no peak years for catching fish for these 15 countries since there are ample years with almost the same size and color. However, the top 3 years are 2006, 1997, and 1996.

Link to video explanation of how Dashboard was interpreted:
https://youtu.be/o_mECUFbNJg

### 3.8 Archiving S3 Buckets (End of Pipeline procedure)

Since Data Lake and Warehouse S3 Buckets will no longer be used by Pipeline, they will be archived.

Configured my Lake and Warehouse S3 buckets to transition the files into Glacier Deep Archive S3 Buckets after 30 days, since these files will not be used anymore but is good to keep in case they are needed for future auditing and checking purposes.

Using Glacier instead of S3 helps optimize cloud expenses by cutting down on costs as Glacier is cheaper than S3.

Deep Archive is used as it's the cheapest option.

### 3.9 Security

Since I am restricted from IAM, I decided to configure KMS to ensure the pipeline is secured and protected.

For each S3 bucket, enable default server-side encryption using SSE-KMS that protects S3 objects when no other key is defined.

KMS key to protect SNS data when no other key is defined.

### 4 Monitoring (CloudWatch)

To monitor the health status of the Pipeline, I created an Alarm to alert potential issues that could arise from Glue usage and Errors that occur when trying to fetch data from API.

The alarm will trigger when the blue line goes beyond the threshold within 5 minutes (to alert engineers about any sudden sharp rise in spending).

To alert issues with fetching data from OpenFisheries API, any error that occurs will be flagged as an error occurring that would cause the error count to exceed the threshold of 0 error.

To monitor Glue usage, the Billing for Glue is used, and the bill for Glue cannot exceed the maximum threshold set in USD. The threshold is $2. If Glue usage deviates by $2 from normal expected usage, this anomalous activity will be notified to the engineers.

**Unresolved error when fetching data from API:**

**Glue's Billing to ensure no exceed budget:**

## 4.1 Critical error (unable to be resolved) in fetching data from Openfisheries API

| | |
|---|---|
| When an error occurs, the blue line WILL exceed the 0 threshold, hence engineers will receive the notification of the failure.<br><br>The alarm will send a notification to the fetching_API topic, and the email endpoint is set to the engineer's email address. | **Conditions**<br><br>Threshold type<br>● Static — Use a value as a threshold / ○ Anomaly detection — Use a band as a threshold<br><br>Whenever Errors is... Define the alarm condition.<br>● Greater > threshold / ○ Greater/Equal >= threshold / ○ Lower/Equal <= threshold / ○ Lower < threshold<br><br>than... Define the threshold value.<br>0<br>Must be a number<br><br>Step 1 Specify metric and conditions<br>Step 2 Configure actions<br>Step 3 Add name and description<br>Step 4 Preview and create<br><br>**Configure actions**<br>**Notification**<br>Alarm state trigger — Define the alarm state that will trigger this action. [Remove]<br>● In alarm — The metric or expression is outside of the defined threshold. / ○ OK — The metric or expression is within the defined threshold. / ○ Insufficient data — The alarm has just started or not enough data is available.<br><br>Send a notification to the following SNS topic — Define the SNS (Simple Notification Service) topic that will receive the notification.<br>● Select an existing SNS topic<br>○ Create new topic<br>○ Use topic ARN to notify other accounts<br><br>Send a notification to...<br>🔍 Fetching_API ✕<br>Only topics belonging to this account are listed here. All persons and applications subscribed to the selected topic will receive notifications.<br><br>Email (endpoints)<br>2202934b@student.tp.edu.sg - View in SNS Console 🗗 |
| This is the message that will appear in the email message if the alarm is triggered.<br><br>The engineer can recognize the email using the distinct alarm name, to know that the notification alarm is regarding an error with fetching the REST API from OpenFisheries. | **Name and description**<br><br>Alarm name<br>error_fetching<br><br>Alarm description - *optional*  View formatting guidelines<br>Edit    Preview<br><br>ALERT: **ERROR** occurred while fetching data from OpenFisheries API<br><br>Up to 1024 characters (119/1024) |

## 4.2 Exceeding Glue Budget

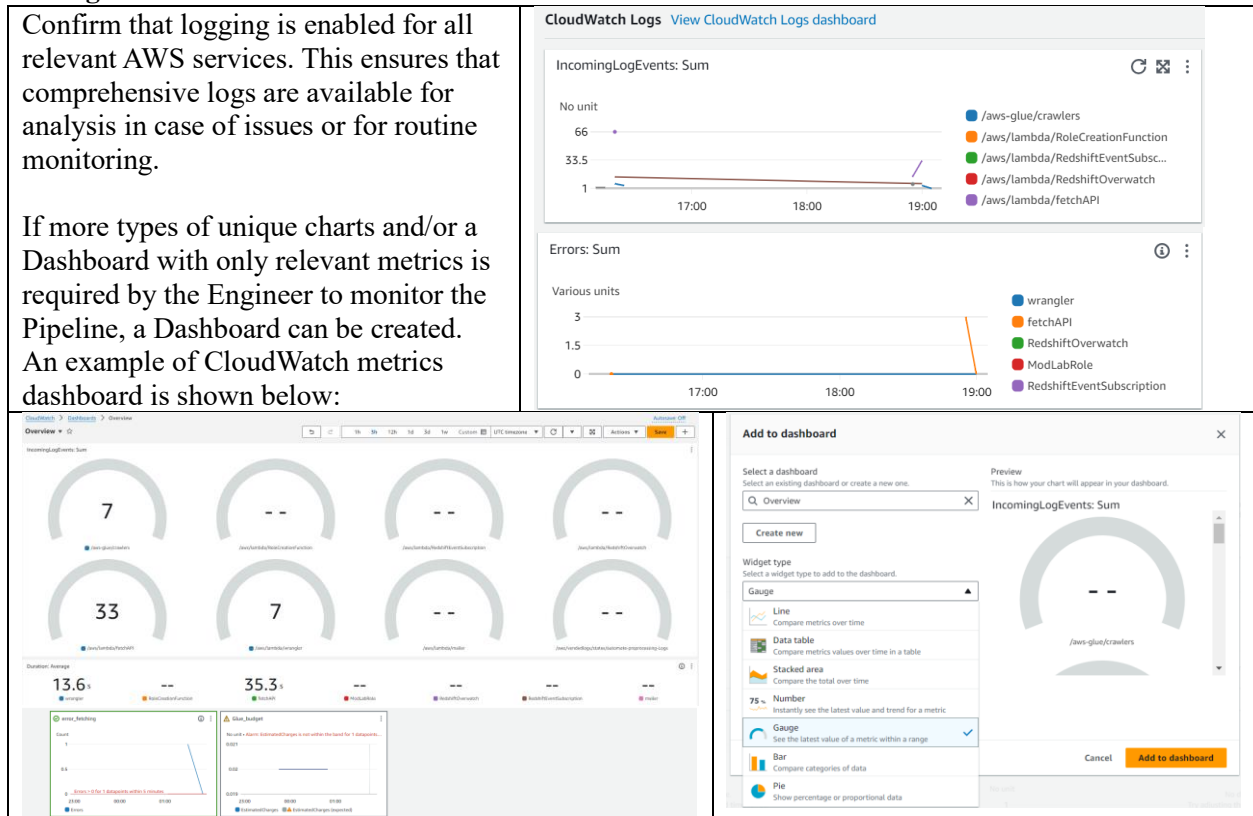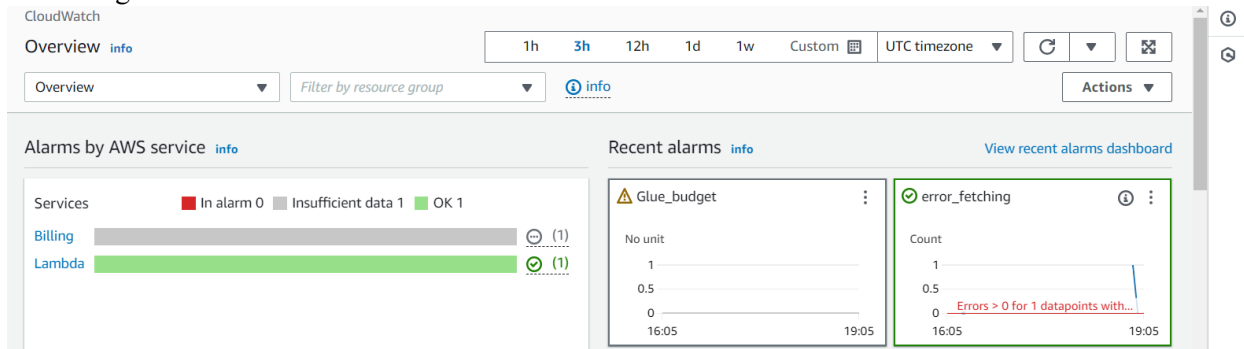| | |
|---|---|
| If Glue usage exceeds the designated budget threshold, the Engineer's email will receive a message to warn him/her.<br><br>Mail will alert the engineer to assess the situation and disable service (if needed) before usage spirals out of control.<br><br>Do not auto-disable service as approval from an engineer is crucial before any major change that disrupts the pipeline. | **Configure actions** — **Add name and description**<br><br>**Notification**<br>Alarm state trigger — Define the alarm state that will trigger this action.<br>● In alarm — The metric or expression is outside of the defined threshold. / ○ OK — The metric or expression is within the defined threshold.<br><br>Send a notification to the following SNS topic — Define the SNS (Simple Notification Service) topic that will receive the notification.<br>● Select an existing SNS topic<br>○ Create new topic<br>○ Use topic ARN to notify other accounts<br><br>Send a notification to...<br>🔍 Glue_budget ✕<br>Only topics belonging to this account are listed here. All persons and applications subscribed to the selected topic will receive notifications.<br><br>Email (endpoints)<br>2202934b@student.tp.edu.sg - View in SNS Console 🗗<br>[Add notification]<br><br>**Name and description**<br>Alarm name<br>glue_budget<br><br>Alarm description - *optional*  View formatting guidelines<br>Edit    Preview<br><br>WARNING: Glue usage has exceeded budget!<br><br>Up to 1024 characters (40/1024)<br><br>ⓘ Markdown formatting is only applied when viewing your alarm in the console. The des... plain text in the alarm notifications. |
| Since exceeding budget is a SERIOUS concern that demands immediate action as the consequences could be dire, a subscription was also added to send an SMS to the engineer's number (because he/she may not check email in time to assess the problem) whereas he should have a phone with him at all times. | Amazon SNS > Subscriptions > Create subscription<br>**Create subscription**<br><br>**Details**<br><br>Topic ARN<br>🔍 arn:aws:sns:us-east-1:451863666267:Glue_budget ✕<br><br>Protocol — The type of endpoint to subscribe<br>SMS ▼<br><br>Endpoint — A mobile number that can receive notifications from Amazon SNS. |

## 4.3 Logs for all services are enabled:

Confirm that logging is enabled for all relevant AWS services. This ensures that comprehensive logs are available for analysis in case of issues or for routine monitoring.

If more types of unique charts and/or a Dashboard with only relevant metrics is required by the Engineer to monitor the Pipeline, a Dashboard can be created.
An example of CloudWatch metrics dashboard is shown below:



## 4.4 Alarms enabled:

When the issue is fixed (after re-running the Step function and successful run), the Service alarm status will change back to OK.



## 4.5 SNS Subscription:

Confirmed and activated, ready to send message to engineer when needed.



[*To see an SNS Mail in action, refer to section 6*]

## 5. Deployment (Automating Pipeline)

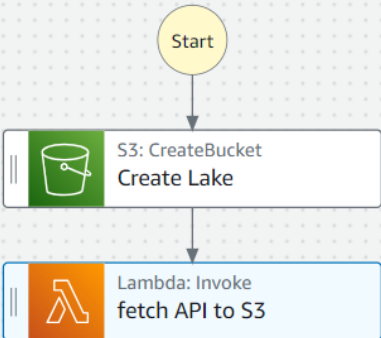1. Added feature to dynamically declare how many datasets to be used.

Modified fetchAPI.py Lambda function to take LIMIT_ROWS from Step input to allow analysts to choose how many datasets they want to use. By default 15 datasets for countries and species are used to align with the project purpose of analyzing 30 datasets.

```
54   def lambda_handler(event, context):
55       LIMIT_ROWS = event.get("LIMIT_ROWS")  # Extract the LIMIT_ROWS value from the input event
56       data_lake(LIMIT_ROWS=LIMIT_ROWS)
```

Analysts using the pipeline can conveniently use this Step function and can choose how many species and countries they want for their analysis by stating it in the input of the Step function.

| Details | Execution input and output | Definition |
|---|---|---|

Input

```
1▾  {
2       "LIMIT_ROWS": 15
```
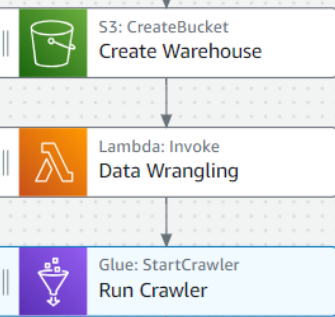
2. Added a timeout to stop fetching API if it takes too long.



+ Add new catcher

**Timeouts**

TimeoutSeconds - *optional*
Fail the state if it runs longer than the specified seconds.

Enter TimeoutSeconds ▾

60    seconds

Must be greater than zero.

3. Run the crawler to populate (or update the existing) the Data Catalog database with data.
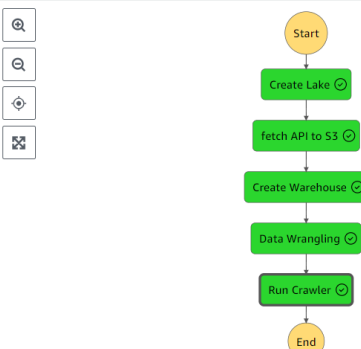


**API Parameters**

JSON object containing the parameters to pass into the {{apiName}} API.

```
1  {
2    "Name": "crawl_data"
3  }
```

4. Prove of that Step workflow is fully functional and works:



**Run Crawler**                                    Test state
aws-sdk:glue:startCrawler

| Input | Output | Details | Definition | Events |
|---|---|---|---|---|

Status                          Type
⊘ Succeeded                      Task: aws-sdk:glue:startCrawler

Resource
aws-sdk:glue:startCrawler

Timeout In Seconds              Heartbeat In Seconds
-                               -

Started After                   Duration
00:00:46.490                    00:00:02.469

This Step function fully prepares the data for analyst's use. Analyst only has to use Athena > Query editor to query the data they need to build their visualizations and report their findings and insights.

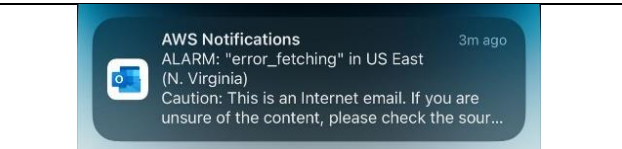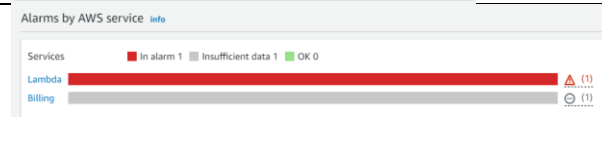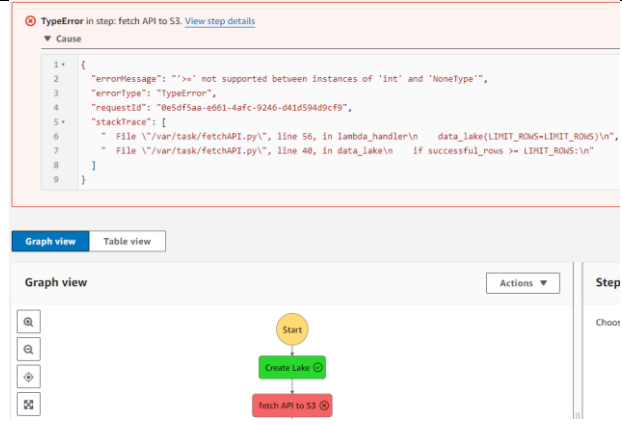## 6.  Evidence that Services work as intended

**Footage of configuring Pipeline**:
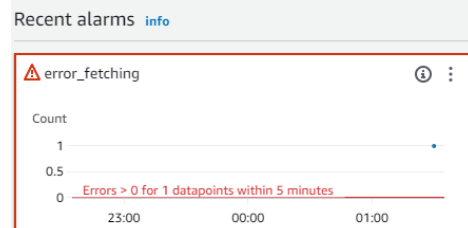
https://youtu.be/qmR2lkdPQWs (Had to redo as my Main Lab got terminated)

**Prove that CloudWatch Alert powered by SNS works:**

To trigger the Error alert, I broke fetchAPI.py script such that it will encounter an error:

```
def lambda_handler(event, context):
    LIMIT_ROWS = event.get("8")  # Ext
    data_lake(LIMIT_ROWS=LIMIT_ROWS)
```

After Step execution failed, this log turned red:









As shown above, when Lambda code has unresolved issues when run, an alarm is triggered to send an email to the engineer to alert them about the issue. The email sent notifies the issue that went wrong.

After testing that CloudWatch Alarm + SNS works properly. The broken fetchAPI function is then reverted to its original working version.



**Proof that CloudWatch Logs for services are enabled:**

## Checklist (configurations of services)
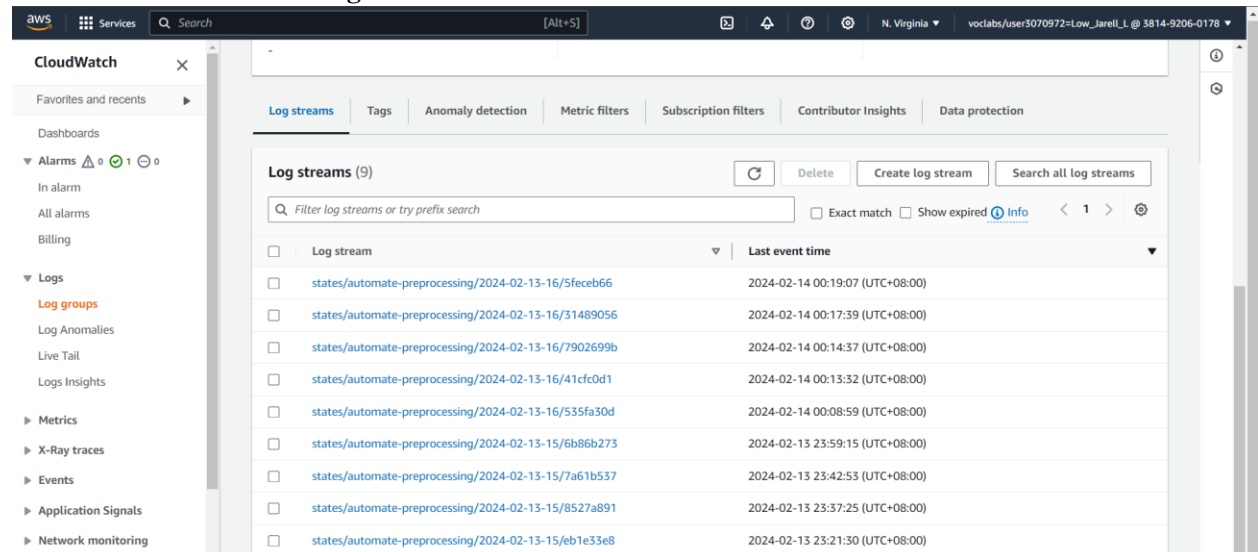### Logging and Monitoring Mechanism for Compliance
Lambda Functions:
- Enable CloudWatch Logs for Lambda functions. (Y)
- Implemented structured logging of unsuccessful operations compliance and auditing. (Y)
- Configure CloudWatch Alarms for critical Lambda metrics: Errors + Exceeding Budget. (Y)

S3 Buckets:
- Enable CloudWatch Metrics for S3 buckets. (Y)
- AWS CloudTrail for S3 bucket access monitoring. Not needed as API not implemented. (N)
- Configure S3 bucket access logging for compliance. (N)
- Enable version control to track any changes made. Only applied to insights bucket (Y)

AWS Glue:
- Enable CloudWatch Logs for AWS Glue ETL jobs and CloudWatch Alarms for Glue job metrics. (Y)

Athena:
- Enable CloudWatch Metrics for Athena. (Y)
- Monitor query execution times and errors. (Y)
- Set up CloudWatch Alarms for Athena query metrics. (Y)

Amazon QuickSight – Visualization tool replaced with Tableau due to AWS budget. (N)
Access Control Management (IAM, Users, Roles, and Policies) – Restrict access from school (N)

### Ideal Configurations for Each Component
S3 Bucket:
- Block public access to the S3 bucket. (Y)
- Implement versioning and logging for S3 bucket. – Track all activity in insights bucket (Y)
- Configure lifecycle policies for data retention. – Archive Lake and Warehouse after usage (Y)

AWS Glue:
- Set up AWS Glue crawlers to automatically infer the structure of the data. (Y)
- Utilize Glue triggers for job scheduling. – Used Step function to trigger Crawler (-)
- Optimize Glue job configurations for resource utilization. – Not applicable (N)

Amazon Athena:
- Organized data into Views for usage for analytics (Y)
- Implement query execution time limits for resource optimization. – No need since dataset is small (N)
- Define query result location for efficient retrieval. – Created Views to store in insights S3 (Y)

Tableau:

### Plan-Perform-Monitor-Reflect (PPMR)
https://tasks.office.com/tp.edu.sg/en-US/Home/Planner/#/plantaskboard?groupId=cb0614be-0efc-4c92-a96a-975d80323615&planId=dM_PKRaeNkqXQGMZPLXNssgACuh1

### Additional Info
My original Lab, Javen_Lai was terminated and Low_Jarell is used in the second half where I had to redo everything as Javen_Lai progress was completely gone.

I skipped API Gateway as a trigger for Pipeline considering how OpenFisheries was last updated in 2018, it wouldn't be efficient to have a live hook to detect changes in API and update accordingly. Instead, I made the Pipeline into a Step function, so analysts can manually update the batch data if needed.

I skipped QuickSight since Tableau is used for visualizations, hence redundant to do visualizations twice which is why I only did advanced analytics to gather all insights in 1 round of EDA.

**References:**

- https://tplms.polite.edu.sg/d2l/le/enhancedSequenceViewer/372754?url=https%3A%2F%2F425c3724-b655-4262-9c13-4699f01f20f3.sequences.api.brightspace.com%2F372754%2Factivity%2F7215919%3FfilterOnDatesAndDepth%3D1 [week 14 workshop to pull data from API]
- https://tplms.polite.edu.sg/d2l/le/enhancedSequenceViewer/372754?url=https%3A%2F%2F425c3724-b655-4262-9c13-4699f01f20f3.sequences.api.brightspace.com%2F372754%2Factivity%2F7839269%3FfilterOnDatesAndDepth%3D1 [week 18 workshop to enable SNS mail]
- https://youtu.be/qmR2lkdPQWs [Speed running AWS Setup]
- https://youtu.be/o_mECUFbNJg [Showcase of Dashboard + How to use and interpret]