# Event-Driven
# Data Validation
# and Processing Pipeline
# Using S3,
# Step Functions, and ECS

# Objective

- Trigger: Use Amazon EventBridge to monitor S3 for newly created objects and trigger an AWS Step Function.

- Validate: Perform data validation using Python scripts with Pandas, ensuring data integrity.

- Move: Transfer validated data to a designated S3 location.

- Condition Check: Evaluate the output of the first validation job to determine readiness for further processing.

- Process: Trigger the second task to perform additional data transformations using Python.

- Run Tasks: Execute both Python scripts as containerized tasks on AWS Fargate with images stored in Amazon ECR.

- Track: Record the status and timestamps of the pipeline's execution in Amazon DynamoDB.

# Tools Used

- Amazon S3: As the source and destination for raw and validated data.

- Amazon EventBridge: For detecting object creation events in S3 and triggering the pipeline.

- AWS Step Functions: For orchestrating the pipeline's execution, including validation and condition checks.

- Python: For writing data validation and processing logic using Pandas.

- Amazon ECS (Fargate): For running Python scripts in a serverless, containerized environment.

- Amazon ECR: For storing container images used in ECS tasks.

- Amazon DynamoDB: For logging pipeline status and tracking execution timestamps.

- Pandas: For performing data validation and transformation operations.

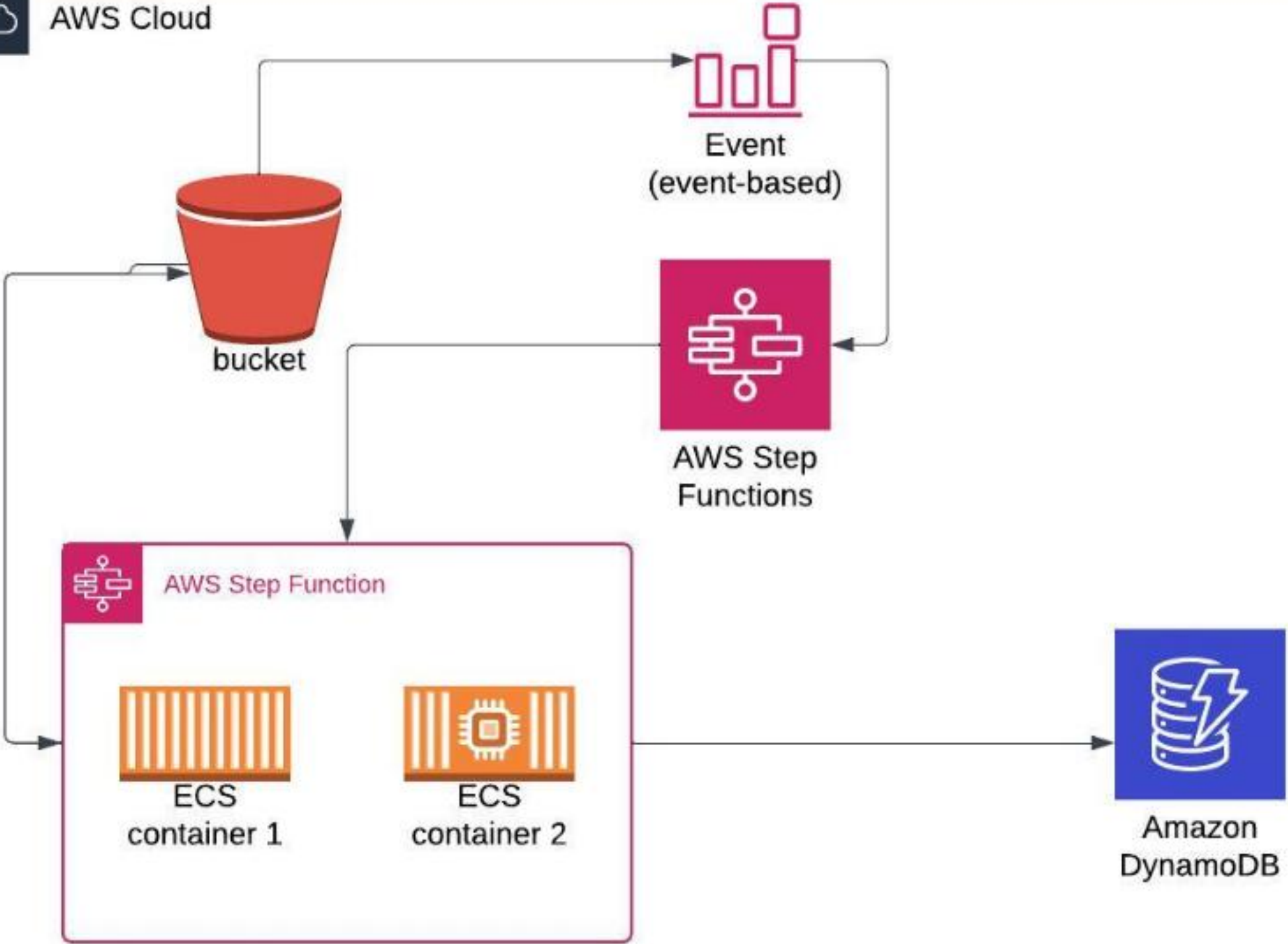- CloudWatch: For monitoring the execution flow and capturing logs.

AWS Cloud

PC → ECR registry → Amazon ECS Anywhere

Learn

Sandbox

Quick starts

▼ **Buses**

Event buses

**Rules**

Global endpoints

Archives

Replays

▼ **Pipes**

Pipes

▼ **Scheduler**

Schedules

Schedule groups

▼ **Integration**

Partner event sources

API destinations

Connections **Updated**

# 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** (1)　　　　　　　　⟳　Delete　　Enable　　Edit　　CloudFormation Template ▼　**Create rule**

| | | Q Find rules | | Any status ▼ | ‹ 1 › | ⚙ |
|---|---|---|---|---|---|---|

| ☐ | **Name** ▲ | **Status** ▽ | **Type** ▽ | **ARN** ▽ | **Description** ▽ |
|---|---|---|---|---|---|
| ☐ | trigger-ecs-tasks | ⊘ Enabled | Standard | ⬜ arn:aws:events:us-east-1:418272783111: rule/trigger-ecs-tasks | xzzz |

Learn

Sandbox

Quick starts

▼ **Buses**

Event buses

**Rules**

Global endpoints

Archives

Replays

▼ **Pipes**

Pipes

▼ **Scheduler**

Schedules

Schedule groups

▼ **Integration**

Partner event sources

API destinations

Connections **Updated**

-tasks                                                    ult

| **Event pattern** | **Targets** | **Monitoring** | **Tags** |

## Event pattern Info

Edit

```
1  {
2    "source": ["aws.s3"],
3    "detail-type": ["Object Created"],
4    "detail": {
5      "bucket": {
6        "name": ["gd-aws-de-labs"]
7      },
8      "object": {
9        "key": [{
10         "prefix": "ecommerce-data/new/order_items"
11       }]
12     }
13   }
14 }
```

Copy

# new/

Copy S3 URI

**Objects** | Properties

## Objects (3) 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 ↗

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📄 order_items_20240601.csv | csv | December 27, 2024, 14:20:28 (UTC+05:30) | 168.8 KB | Standard |
| ☐ | 📄 orders_20240601.csv | csv | December 27, 2024, 13:58:12 (UTC+05:30) | 94.5 KB | Standard |
| ☐ | 📄 products.csv | csv | December 27, 2024, 13:58:23 (UTC+05:30) | 4.1 MB | Standard |

**Step Functions**

**State machines**

Activities

▼ **Developer resources**

Online learning workshop [↗]

Local Development

Data flow simulator

Feature spotlight

Documentation [↗]

Join our feedback panel [↗]

**Arn**
[□] arn:aws:states:us-east-1:418272783111:stateMachine:testing

**IAM role ARN**
[□] arn:aws:iam::418272783111:role/step-functiontesting [↗]

**Type**
Standard

**Status**
Active

**Creation date**
Dec 27, 2024, 13:37:56 (UTC+05:30)

**X-Ray tracing**
Disabled

| Executions | Monitoring | Logging | Definition | Aliases | Versions | Tags |
|---|---|---|---|---|---|---|

**Executions** (0/2)

Refresh     View details     Stop execution     Redrive     **Start execution**

🔍 Filter executions by property     Filter by status ▼     🗓 Last 15 months     Local timezone ▼     2 matches     < 1 >     ⚙

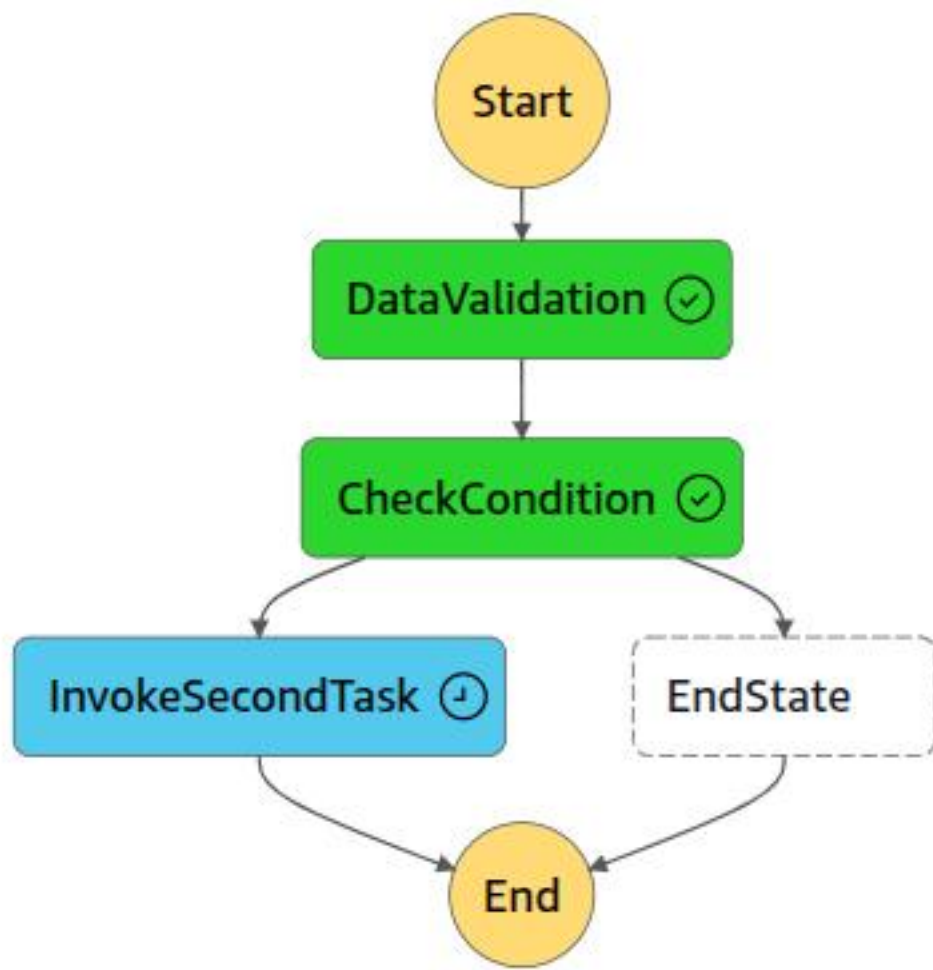| | Name ▽ | Status ▽ | Start Time (local) ▼ | End Time (local) ▽ | Dura |
|---|---|---|---|---|---|
| ☐ | e371b047-2ca5-9879-b99d-73133e1fe1... | ⏱ Running | Dec 27, 2024, 14:20:28 | - | - |
| ☐ | ea4f1378-aa24-4e5e-a78a-7b80eed7bb9a | ⊘ Succeeded | Dec 27, 2024, 13:38:13 | Dec 27, 2024, 13:39:55 | C |

# Step Functions

## State machines

## Activities

## ▼ Developer resources

Online learning workshop [↗]

Local Development

Data flow simulator

Feature spotlight

Documentation [↗]

Join our feedback panel [↗]

## Graph view

Actions ▼



Start

DataValidation ⊘

CheckCondition ⊘

InvokeSecondTask ⏲

EndState

End

⏲ In progress | ⊗ Failed | ⚠ Caught error | ⊖ Canceled | ⊘ Succeeded

## Step details

Choose a step to view its details.

**Event view** | **State view**

```json
1  {
2    "Comment": "A Step Function to process ECS tasks and handle output conditions",
3    "StartAt": "DataValidation",
4    "States": {
5      "DataValidation": {
6        "Type": "Task",
7        "Resource": "arn:aws:states:::ecs:runTask.sync",
8        "Parameters": {
9          "Cluster": "ecommerce-data-ingestion-cluster",
10         "TaskDefinition": "arn:aws:ecs:us-east-1:418272783111:task-definition/task-data-validation:4",
11         "LaunchType": "FARGATE",
12         "Overrides": {
13           "ContainerOverrides": [
14             {
15               "Name": "container-data-validation",
16               "Environment": [
17                 {
18                   "Name": "TASK_TOKEN",
19                   "Value.$": "$$.Task.Token"
20                 }
21               ]
22             }
23           ]
24         },
25         "NetworkConfiguration": {
26           "AwsvpcConfiguration": {
27             "Subnets": [
28               "subnet-09fd1652b231f5ca0",
29               "subnet-0cdf12a5d9170a6ff"
30             ],
31             "AssignPublicIp": "ENABLED"
```

```json
  4        "States": {
  5          "DataValidation": {
  8            "Parameters": {
 12              "Overrides": {
 23              ]
 24            },
 25              "NetworkConfiguration": {
 26                "AwsvpcConfiguration": {
 27                  "Subnets": [
 28                    "subnet-09fd1652b231f5ca0",
 29                    "subnet-0cdf12a5d9170a6ff"
 30                  ],
 31                  "AssignPublicIp": "ENABLED"
 32                }
 33              }
 34            },
 35            "Next": "CheckCondition",
 36            "ResultPath": "$.taskOutput"
 37          },
 38          "CheckCondition": {
 39            "Type": "Choice",
 40            "Choices": [
 41              {
 42                "And": [
 43                  {
 44                    "Variable": "$.taskOutput.orders",
 45                    "NumericEquals": 1
 46                  },
 47                  {
 48                    "Variable": "$.taskOutput.order_items",
 49                    "NumericEquals": 1
 50                  },
```

```json
 4        "States": {
38          "CheckCondition": {
40            "Choices": [
41              {
42                "And": [
43                  {
44                    "Variable": "$.taskOutput.orders",
45                    "NumericEquals": 1
46                  },
47                  {
48                    "Variable": "$.taskOutput.order_items",
49                    "NumericEquals": 1
50                  },
51                  {
52                    "Variable": "$.taskOutput.products",
53                    "NumericEquals": 1
54                  }
55                ],
56                "Next": "InvokeSecondTask"
57              }
58            ],
59            "Default": "EndState"
60          },
61          "InvokeSecondTask": {
62            "Type": "Task",
63            "Resource": "arn:aws:states:::ecs:runTask.sync",
64            "Parameters": {
65              "Cluster": "ecommerce-data-ingestion-cluster",
66              "TaskDefinition": "arn:aws:ecs:us-east-1:418272783111:task-definition/task-etl-calculations:3",
67              "LaunchType": "FARGATE",
68              "NetworkConfiguration": {
```

```json
  4        "States": {
 38          "CheckCondition": {
 58            ],
 59            "Default": "EndState"
 60          },
 61          "InvokeSecondTask": {
 62            "Type": "Task",
 63            "Resource": "arn:aws:states:::ecs:runTask.sync",
 64            "Parameters": {
 65              "Cluster": "ecommerce-data-ingestion-cluster",
 66              "TaskDefinition": "arn:aws:ecs:us-east-1:418272783111:task-definition/task-etl-calculations:3",
 67              "LaunchType": "FARGATE",
 68              "NetworkConfiguration": {
 69                "AwsvpcConfiguration": {
 70                  "Subnets": [
 71                    "subnet-09fd1652b231f5ca0",
 72                    "subnet-0cdf12a5d9170a6ff"
 73                  ],
 74                  "AssignPublicIp": "ENABLED"
 75                }
 76              }
 77            },
 78            "End": true
 79          },
 80          "EndState": {
 81            "Type": "Pass",
 82            "End": true
 83          }
 84        }
 85      }
```

```python
import boto3
import pandas as pd
import logging
import json
import os

logging.basicConfig(level=logging.INFO)
logger = logging.getLogger()

s3_bucket = 'gd-aws-de-labs'

folders = {
    'orders': 'ecommerce-data/new/orders',
    'order_items': 'ecommerce-data/new/order_items',
    'products': 'ecommerce-data/new/products'
}

ready_folders = {
    'orders': 'ecommerce-data/orders',
    'order_items': 'ecommerce-data/order_items',
    'products': 'ecommerce-data/products'
}

def list_files(bucket, prefix):
    s3_client = boto3.client('s3')
    paginator = s3_client.get_paginator('list_objects_v2')
    response_iterator = paginator.paginate(Bucket=bucket, Prefix=prefix)

    files = []
    for page in response_iterator:
        if 'Contents' in page:
            files += [content['Key'] for content in page['Contents'] if content['Key'].endswith('.csv')]
```

```python
24    def list_files(bucket, prefix):

26        paginator = s3_client.get_paginator('list_objects_v2')
27        response_iterator = paginator.paginate(Bucket=bucket, Prefix=prefix)

29        files = []
30        for page in response_iterator:
31            if 'Contents' in page:
32                files += [content['Key'] for content in page['Contents'] if content['Key'].endswith('.csv')]
33        return files

35    def read_from_s3(bucket, file_key):
36        s3_client = boto3.client('s3')
37        obj = s3_client.get_object(Bucket=bucket, Key=file_key)
38        df = pd.read_csv(obj['Body'])
39        if 'created_at' in df.columns:
40            df['created_at'] = pd.to_datetime(df['created_at'], errors='coerce', utc=True)
41            df['created_at'] = df['created_at'].dt.tz_localize(None)
42        return df

44    def check_format(df, column_formats, file_key):
45        missing_columns = [col for col in column_formats if col not in df.columns]
46        if missing_columns:
47            missing_columns_str = ", ".join(missing_columns)
48            logger.error(f"Missing columns: {missing_columns_str} in {file_key} dataset")
49            return False
50        for column, expected_type in column_formats.items():
51            if column in df.columns:
52                if not pd.api.types.is_dtype_equal(df[column].dtype, expected_type):
53                    logger.error(f"Format check failed for column {column}: expected {expected_type}, found {df[column].dtype} i
54                    return False
55            else:
```

```python
43
44  def check_format(df, column_formats, file_key):
45      missing_columns = [col for col in column_formats if col not in df.columns]
46      if missing_columns:
47          missing_columns_str = ", ".join(missing_columns)
48          logger.error(f"Missing columns: {missing_columns_str} in {file_key} dataset")
49          return False
50      for column, expected_type in column_formats.items():
51          if column in df.columns:
52              if not pd.api.types.is_dtype_equal(df[column].dtype, expected_type):
53                  logger.error(f"Format check failed for column {column}: expected {expected_type}, found {df[column].dtype} i
54                  return False
55          else:
56              logger.error(f"Column {column} is missing from the DataFrame in {file_key} dataset.")
57              return False
58      return True
59
60  def move_file(s3_client, bucket, file_key, ready_prefix):
61      copy_source = {'Bucket': bucket, 'Key': file_key}
62      new_key = ready_prefix + '/' + file_key.split('/')[-1]
63      s3_client.copy_object(CopySource=copy_source, Bucket=bucket, Key=new_key)
64      s3_client.delete_object(Bucket=bucket, Key=file_key)
65
66  def main():
67      task_token = os.getenv('TASK_TOKEN')
68      if not task_token:
69          logger.error("TASK_TOKEN environment variable is missing")
70          return
71
72      results = {"orders": 0, "order_items": 0, "products": 0}
73      s3_client = boto3.client('s3')
74
```

```python
def main():

    if not task_token:
        logger.error("TASK_TOKEN environment variable is missing")
        return

    results = {"orders": 0, "order_items": 0, "products": 0}
    s3_client = boto3.client('s3')

    column_formats = {
        'orders': {
            'order_id': 'int64',
            'user_id': 'int64',
            'status': 'object',
            'created_at': 'datetime64[ns]',
            'num_of_item': 'int64'
        },
        'order_items': {
            'id': 'int64',
            'order_id': 'int64',
            'user_id': 'int64',
            'product_id': 'int64',
            'created_at': 'datetime64[ns]'
        },
        'products': {
            'id': 'int64',
            'sku': 'object',
            'cost': 'float64',
            'category': 'object',
            'name': 'object',
            'retail_price': 'float64',
            'department': 'object'
```

```python
 66    def main():
 91            'id': 'int64',
 92            'sku': 'object',
 93            'cost': 'float64',
 94            'category': 'object',
 95            'name': 'object',
 96            'retail_price': 'float64',
 97            'department': 'object'
 98        }
 99    }

100
101    for folder_key in folders.keys():
102        prefix = folders[folder_key]
103        files = list_files(s3_bucket, prefix)
104
105        for file_name in files:
106            df = read_from_s3(s3_bucket, file_name)
107            format_result = check_format(df, column_formats[folder_key], folder_key)
108            if format_result
109                results[fo    (variable) s3_client: _ | Any
110                move_file(s3_client, s3_bucket, file_name, ready_folders[folder_key])
111
112    logger.info(f"Results: {json.dumps(results)}")
113
114    stepfunctions_client = boto3.client('stepfunctions')
115    stepfunctions_client.send_task_success(
116        taskToken=task_token,
117        output=json.dumps(results)
118    )
119    logger.info("Task output sent to Step Functions successfully")
120
121 if __name__ == '__main__':
```

ECS-STEP-PYTHON

- > data
- ∨ docker-data-validity
  - 🐍 app.py
  - 🐳 Dockerfile
  - ☰ requirements.txt
- ∨ docker-dynamo-etl-wrangler
  - 🐍 app.py
  - 🐳 Dockerfile
  - ☰ requirements.txt
- ∨ step-functions
  - {} step-function.json
  - {} step-functions-iam-execution-pol...
- $ docker-commands.sh
- {} event-pattern.json
- 📘 test.ipynb

```dockerfile
FROM python:3.10

WORKDIR /usr/src/app

COPY . .

RUN pip install --no-cache-dir -r requirements.txt

CMD ["python3", "app.py"]
```

> OUTLINE

> TIMELINE

```
docker-data-validity  >  ☰ requirements.txt

     1    pandas
     2    boto3
```

```python
import pandas as pd
import awswrangler as wr
import boto3
from decimal import Decimal, Inexact, Rounded
import logging

logging.basicConfig(level=logging.INFO, format='%(asctime)s - %(levelname)s - %(message)s')
logger = logging.getLogger()

base_path = "s3://gd-aws-de-labs/ecommerce-data/"
orders_path = f"{base_path}orders/"
order_items_path = f"{base_path}order_items/"
products_path = f"{base_path}products/"

orders_archive_path = f"{base_path}archived/"
order_items_archive_path = f"{base_path}archived/"
products_archive_path = f"{base_path}archived/"

try:
    logger.info("Reading datasets from S3")
    orders = wr.s3.read_csv(path=orders_path)
    order_items = wr.s3.read_csv(path=order_items_path)
    products = wr.s3.read_csv(path=products_path)
except Exception as e:
    logger.error(f"Error reading datasets: {e}")
    raise
try:
    logger.info("Joining datasets")
    df = orders.merge(order_items, on='order_id', how='inner', suffixes=('_order', '_item'))
    df = df.merge(products, left_on="product_id", right_on="id", how='inner', suffixes=('', '_product'))

    df.rename(columns={'created_at_order': 'order_date', 'status_order': 'order_status'}, inplace=True)
```

```python
27     try:
28         logger.info("Joining datasets")
29         df = orders.merge(order_items, on='order_id', how='inner', suffixes=('_order', '_item'))
30         df = df.merge(products, left_on="product_id", right_on="id", how='inner', suffixes=('', '_product'))
31
32         df.rename(columns={'created_at_order': 'order_date', 'status_order': 'order_status'}, inplace=True)
33
34         df['order_date'] = pd.to_datetime(df['order_date'].str.replace(' UTC', ''), format='%Y-%m-%d %H:%M:%S', errors='coerce')
35
36         grouped = df.groupby([df['order_date'].dt.date, 'category'])
37
38     except Exception as e:
39         logging.error(f"Error processing data: {e}")
40         raise
41
42     def safe_convert_decimal(val):
43         try:
44             return Decimal(str(val))
45         except:
46             return Decimal(str(round(val, 2)))
47
48     def calculate_kpis(group):
49         total_revenue = (group['sale_price'] * group['num_of_item']).sum()
50         avg_order_value = group['sale_price'].mean()
51         avg_return_rate = group[group['order_status'] == 'Returned']['order_id'].nunique() / group['order_id'].nunique()
52         return pd.Series({
53             'daily_revenue': safe_convert_decimal(total_revenue),
54             'avg_order_value': safe_convert_decimal(avg_order_value),
55             'avg_return_rate': safe_convert_decimal(avg_return_rate)
56         })
57
58     try:
```

```python
46            return Decimal(str(round(val, 2)))
47
48    def calculate_kpis(group):
49        total_revenue = (group['sale_price'] * group['num_of_item']).sum()
50        avg_order_value = group['sale_price'].mean()
51        avg_return_rate = group[group['order_status'] == 'Returned']['order_id'].nunique() / group['order_id'].nunique()
52        return pd.Series({
53            'daily_revenue': safe_convert_decimal(total_revenue),
54            'avg_order_value': safe_convert_decimal(avg_order_value),
55            'avg_return_rate': safe_convert_decimal(avg_return_rate)
56        })
57
58    try:
59        df_category_wise_summary = grouped.apply(calculate_kpis).reset_index()
60    except Exception as e:
61        logger.error(f"Error calculating category wise KPIS: {e}")
62        raise
63
64    def calculate_order_kpis(group):
65        total_orders = group['order_id'].nunique()
66        total_revenue = (group['sale_price'] * group['num_of_item']).sum()
67        total_items_sold = group['num_of_item'].sum()
68        return_rate = group[group['order_status'] == 'Returned']['order_id'].nunique() / total_orders
69        unique_customers = group['user_id_order'].nunique()
70
71        return pd.Series({
72            'total_orders': total_orders,
73            'total_revenue': safe_convert_decimal(total_revenue),
74            'total_items_sold': total_items_sold,
75            'return_rate': safe_convert_decimal(return_rate),
76            'unique_customers': unique_customers
```

```python
64      def calculate_order_kpis(group):

71          return pd.Series({
72              'total_orders': total_orders,
73              'total_revenue': safe_convert_decimal(total_revenue),
74              'total_items_sold': total_items_sold,
75              'return_rate': safe_convert_decimal(return_rate),
76              'unique_customers': unique_customers
77          })
78
79      try:
80          df_daily_order_summary = df.groupby(df['order_date'].dt.date).apply(calculate_order_kpis).reset_index()
81      except Exception as e:
82          logger.error(f"Error calculating daily order KPIs: {e}")
83          raise
84
85      dynamodb = boto3.resource('dynamodb')
86      kpi_table = dynamodb.Table('category_wise_summary')
87      order_kpi_table = dynamodb.Table('daily_order_summary')
88
89      def upsert_dynamodb(table, item):
90          try:
91              table.put_item(Item=item)
92          except Exception as e:
93              logging.error(f"Error upserting data into DynamoDB: {e}")
94
95      try:
96          logger.info("Inserting category-wise KPI data into DynamoDB")
97          for index, row in df_category_wise_summary.iterrows():
98              kpi_data = {
99                  'category': row['category'],
100                 'order_date': str(row['order_date']),
```

```python
         logging.error(f"Error upserting data into DynamoDB: {e}")

 try:
     logger.info("Inserting category-wise KPI data into DynamoDB")
     for index, row in df_category_wise_summary.iterrows():
         kpi_data = {
             'category': row['category'],
             'order_date': str(row['order_date']),
             'daily_revenue': safe_convert_decimal(row['daily_revenue']),
             'avg_order_value': safe_convert_decimal(row['avg_order_value']),
             'avg_return_rate': safe_convert_decimal(row['avg_return_rate'])
         }
         upsert_dynamodb(kpi_table, kpi_data)
 except Exception as e:
     logger.error(f"Error inserting category-wise KPI data: {e}")
     raise

 try:
     logger.info("Inserting daily order KPI data into DynamoDB")
     for index, row in df_daily_order_summary.iterrows():
         order_kpi_data = {
             'order_date': str(row['order_date']),
             'total_orders': row['total_orders'],
             'total_revenue': safe_convert_decimal(row['total_revenue']),
             'total_items_sold': row['total_items_sold'],
             'return_rate': safe_convert_decimal(row['return_rate']),
             'unique_customers': row['unique_customers']
         }
         upsert_dynamodb(order_kpi_table, order_kpi_data)
 except Exception as e:
     logger.error(f"Error inserting daily order KPI data: {e}")
     raise
```

```python
119                unique_customers": row["unique_customers"]
120            }
121            upsert_dynamodb(order_kpi_table, order_kpi_data)
122    except Exception as e:
123        logger.error(f"Error inserting daily order KPI data: {e}")
124        raise
125
126    logger.info("KPIs calculated and saved successfully to DynamoDB tables category_wise_summary and daily_order_summary")
127
128    s3_client = boto3.client('s3')
129
130    def move_s3_files(source_path, archive_path):
131        s3_client = boto3.client('s3')
132        bucket = 'gd-aws-de-labs'
133        source_prefix = source_path.replace("s3://gd-aws-de-labs/", "")
134        archive_prefix = archive_path.replace("s3://gd-aws-de-labs/", "")
135
136        try:
137            response = s3_client.list_objects_v2(Bucket=bucket, Prefix=source_prefix)
138            if 'Contents' in resp          (variable) obj: Any
139                for obj in respon
140                    source_key = obj['Key']
141                    archive_key = source_key.replace(source_prefix, archive_prefix)
142                    if not archive_key.startswith(archive_prefix):
143                        archive_key = archive_prefix + source_key.split('/')[-1]
144
145                    s3_client.copy_object(
146                        CopySource={'Bucket': bucket, 'Key': source_key},
147                        Bucket=bucket,
148                        Key=archive_key
149                    )
150                    s3_client.delete_object(Bucket=bucket, Key=source_key)
```

```python
130    def move_s3_files(source_path, archive_path):
131        s3_client = boto3.client('s3')
132        bucket = 'gd-aws-de-labs'
133        source_prefix = source_path.replace("s3://gd-aws-de-labs/", "")
134        archive_prefix = archive_path.replace("s3://gd-aws-de-labs/", "")
135
136        try:
137            response = s3_client.list_objects_v2(Bucket=bucket, Prefix=source_prefix)
138            if 'Contents' in response:
139                for obj in response['Contents']:
140                    source_key = obj['Key']
141                    archive_key = source_key.replace(source_prefix, archive_prefix)
142                    if not archive_key.startswith(archive_prefix):
143                        archive_key = archive_prefix + source_key.split('/')[-1]
144
145                    s3_c    (function) CopySource: Any
146                        CopySource={'Bucket': bucket, 'Key': source_key},
147                        Bucket=bucket,
148                        Key=archive_key
149                    )
150                    s3_client.delete_object(Bucket=bucket, Key=source_key)
151            logger.info(f"Moved files from {source_path} to {archive_path} successfully")
152        except Exception as e:
153            logging.error(f"Error files {source_path} to {archive_path} successfully")
154            raise Exception(f"Error files {source_path} to {archive_path} successfully")
155
156    try:
157        move_s3_files(orders_path, orders_archive_path)
158        move_s3_files(order_items_path, order_items_archive_path)
159        move_s3_files(products_path, products_archive_path)
160    except Exception as e:
161        logger.error(f"Error moving files to archived folders: {e}")
```

```python
130     def move_s3_files(source_path, archive_path):
144
145                     s3_client.copy_object(
146                         CopySource={'Bucket': bucket, 'Key': source_key},
147                         Bucket=bucket,
148                         Key=archive_key
149                     )
150                     s3_client.delete_object(Bucket=bucket, Key=source_key)
151             logger.info(f"Moved files from {source_path} to {archive_path} successfully")
152         except Exception as e:
153             logging.error(f"Error files {source_path} to {archive_path} successfully")
154             raise Exception(f"Error files {source_path} to {archive_path} successfully")
155
156     try:
157         move_s3_files(orders_path, orders_archive_path)
158         move_s3_files(order_items_path, order_items_archive_path)
159         move_s3_files(products_path, products_archive_path)
160     except Exception as e:
161         logger.error(f"Error moving files to archived folders: {e}")
162         raise
```

docker-dynamo-etl-wrangler > Dockerfile

- data
- docker-data-validity
  - app.py
  - Dockerfile
  - requirements.txt
- docker-dynamo-etl-wrangler
  - app.py
  - Dockerfile
  - requirements.txt
- step-functions
  - {} step-function.json
  - {} step-functions-iam-execution-pol...
- $ docker-commands.sh
- {} event-pattern.json
- test.ipynb

```dockerfile
1  FROM python:3.10
2
3  WORKDIR /usr/src/app
4
5  COPY . .
6
7  RUN pip install --no-cache-dir -r requirements.txt
8
9  CMD ["python3", "app.py"]
10
```

OUTLINE

TIMELINE

(i) Do you want to install the recommended 'Docker' extension from Microsoft for the Docker language?

Install    Show Recommendations

ECS-STEP-PYTHON

> data
∨ docker-data-validity
  🐍 app.py
  🐳 Dockerfile
  ≡ requirements.txt
∨ docker-dynamo-etl-wrangler
  🐍 app.py
  🐳 Dockerfile
  ≡ requirements.txt
∨ step-functions
  {} step-function.json
  {} step-functions-iam-execution-pol...
$ docker-commands.sh
{} event-pattern.json
📓 test.ipynb

```
1  pandas
2  awswrangler
```

ⓘ Do you want to install the recommended 'Docker' extension from Microsoft for the Docker language?

Install    Show Recommendations

OUTLINE

TIMELINE

```sh
$ docker-commands.sh
1   # Docker Authentication with ECR
2   aws ecr get-login-password \
3           --region us-east-1 | docker login \
4           --username AWS 418272783111\
5           --password-stdin .dkr.ecr.us-east-1.amazonaws.com
6
7   # Commands for data-validation task
8   docker build -t ecom_data_validation .
9   docker run -d -v ~/.aws:/root/.aws ecom_data_validation
10  docker tag ecom_data_validation:latest 418272783111.dkr.ecr.us-east-1.amazonaws.com/ecommerce-pipelines:ecom_data_validation
11  docker push 418272783111.dkr.ecr.us-east-1.amazonaws.com/ecommerce-pipelines:ecom_data_validation
12
13  # Newly tagged image
14  docker tag ecom_data_validation:latest {aws-account-number}.dkr.ecr.us-east-1.amazonaws.com/ecommerce-pipelines:ecom_data_va
15  docker push {aws-account-number}.dkr.ecr.us-east-1.amazonaws.com/ecommerce-pipelines:ecom_data_validation_v1
16
17
18  # Commands for ETL Job
19  docker build -t etl_aggregations .
20  docker run -d -v ~/.aws:/root/.aws etl_aggregations
21  docker tag etl_aggregations:latest 418272783111.dkr.ecr.us-east-1.amazonaws.com/ecommerce-pipelines:etl_aggregations
22  docker push 418272783111.dkr.ecr.us-east-1.amazonaws.com/ecommerce-pipelines:etl_aggregations
23
24
25  # Get task definition once created
26  aws ecs describe-task-definition \
27                  --task-definition redshift-ingestion \
28                  --query taskDefinition "Angle Bracket" task-definition.json
29
30
31  aws ecs describe-task-definition --task-definition redshift-ingestion
```

# ecommerce-data/

Copy S3 URI

**Objects** | **Properties**

## Objects (2) 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 ↗

🔍 Find objects by prefix                                          < 1 >    ⚙

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📁 archived/ | Folder | - | - | - |
| ☐ | 📁 new/ | Folder | - | - | - |

| | Name | Event types | Filters | Destination type | Destination |
|---|---|---|---|---|---|

### No event notifications

Choose **Create event notification** to be notified when a specific event occurs.

Create event notification

## Amazon EventBridge

Edit

For additional capabilities, use Amazon EventBridge to build event-driven applications at scale using S3 event notifications. Learn more ☐ or see EventBridge pricing ☐

**Send notifications to Amazon EventBridge for all events in this bucket**
On

## Transfer acceleration

Edit

Use an accelerated endpoint for faster data transfers. Learn more ☐

**Transfer acceleration**
Disabled

## Object Lock

Edit

Store objects using a write-once-read-many (WORM) model to help you prevent objects from being deleted or overwritten for a fixed amount of time or indefinitely. Object Lock works only in versioned buckets. Learn more ☐
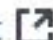
# Amazon Elastic Container Service

Clusters **Updated**

Namespaces

**Task definitions**

Account settings **Updated**

Install AWS Copilot [↗]

Amazon ECR [↗]

Repositories [↗]

AWS Batch [↗]

Documentation [↗]

Discover products [↗]

Subscriptions [↗]

## Task definitions (2) Info

Last updated
27 December 2024 at 14:30 (UTC+5:30)

Deploy ▼    Create new revision ▼    **Create new task definition** ▼

**Filter by status**

🔍 Filter task definitions

Active ▼    ‹ 1 ›    ⚙

| Task definition ▽ | Status of last revision |
|---|---|
| ◯ task-data-validation | ⊘ ACTIVE |
| ◯ task-etl-calculations | ⊘ ACTIVE |

# Amazon Elastic Container Service

Clusters **Updated**

Namespaces

**Task definitions**

Account settings **Updated**

Install AWS Copilot [↗]

Amazon ECR [↗]

Repositories [↗]

AWS Batch [↗]

Documentation [↗]

Discover products [↗]

Subscriptions [↗]

## task-data-validation:4

Deploy ▼   Actions ▼   **Create new revision ▼**

### Overview  Info

**ARN**
[⧉] arn:aws:ecs:us-east-1:418
272783111:task-definition/ta
sk-data-validation:4

**Status**
⊘ ACTIVE

**Time created**
27 December 2024 at
13:22 (UTC+5:30)

**App environment**
Fargate

**Task role**
custom-ecs-task-execution-
role [↗]

**Task execution role**
custom-ecs-task-
execution-role [↗]

**Operating
system/Architecture**
Linux/X86_64

**Network mode**
awsvpc

**Fault injection**
⊖ Turned off

---

**Containers**  |  JSON  |  Task placement  |  Volumes (0)  |  Requires attributes  |  Tags

### Task size

**Task CPU**
2,048 units (2 vCPU)

**Task memory**
8,192 MiB (8 GB)

**Amazon Elastic Container Service**  ‹

Clusters  **Updated**

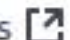Namespaces

**Task definitions**

Account settings  **Updated**

Install AWS Copilot ↗

Amazon ECR ↗

Repositories ↗

AWS Batch ↗

Documentation ↗

Discover products ↗

Subscriptions ↗

---

**Containers**   JSON   **Task placement**   **Volumes (0)**   **Requires attributes**   **Tags**

## Task size

**Task CPU**
2,048 units (2 vCPU)

**Task CPU maximum allocation for containers**

CPU (unit) [ bar chart extending from 0 to ~1900 ]

0   200   400   600   800   1000   1400   1800

☐ container-data-validation   ■ Shared task CPU

**Task memory**
8,192 MiB (8 GB)

**Task memory maximum allocation for container memory reservation**

Memory (MiB) [ bar chart extending from 0 to ~8000 ]

0   1000   2000   3000   4000   5000   6000   7000   8000

☐ container-data-validation   ■ Shared task memory

## Containers  Info

| Container name | Image | Private reg... | Essential | CPU | Memory ha... | GPU |
|---|---|---|---|---|---|---|
| container-dat... | 🗗 41827... | - | Yes | 0 | -/- | - |

# Amazon Elastic Container Service

Clusters **Updated**

Namespaces

**Task definitions**

Account settings **Updated**

Install AWS Copilot [↗]

Amazon ECR [↗]

Repositories [↗]

AWS Batch [↗]

Documentation [↗]

Discover products [↗]

Subscriptions [↗]

## task-etl-calculations:3

Deploy ▼   Actions ▼   **Create new revision** ▼

### Overview Info

**ARN**
[⧉] arn:aws:ecs:us-east-1:418
272783111:task-definition/ta
sk-etl-calculations:3

**Status**
⊘ ACTIVE

**Time created**
27 December 2024 at
12:43 (UTC+5:30)

**App environment**
Fargate

**Task role**
custom-ecs-task-execution-
role [↗]

**Task execution role**
custom-ecs-task-
execution-role [↗]

**Operating
system/Architecture**
Linux/X86_64

**Network mode**
awsvpc

**Fault injection**
⊖ Turned off

---

**Containers** | JSON | Task placement | Volumes (0) | Requires attributes | Tags

---

### Task size

**Task CPU**
2,048 units (2 vCPU)

**Task memory**
10,240 MiB (10 GB)

**Amazon Elastic Container Service**

Clusters **Updated**

Namespaces

**Task definitions**

Account settings **Updated**

Install AWS Copilot ↗

Amazon ECR ↗

Repositories ↗

AWS Batch ↗

Documentation ↗

Discover products ↗

Subscriptions ↗

| Containers | JSON | Task placement | Volumes (0) | Requires attributes | Tags |

## Task size

**Task CPU**
2,048 units (2 vCPU)

**Task CPU maximum allocation for containers**

CPU (unit)

0    200   400   600   800  1000        1400        1800

⬛ container-etl-aggregations    ⬛ Shared task CPU

**Task memory**
10,240 MiB (10 GB)

**Task memory maximum allocation for container memory reservation**

Memory (MiB)

0        2000      4000      6000      8000

⬛ container-etl-aggregations    ⬛ Shared task memory

## Containers  Info

| Container name | Image | Private reg... | Essential | CPU | Memory ha... | GPU |
|---|---|---|---|---|---|---|
| container-etl-... | ⧉ 41827... | - | Yes | 0 | -/- | - |

# DynamoDB

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

# PartiQL editor

Operations performed using the PartiQL editor might incur charges. Learn more [↗]

## Tables (2)

Find tables

< 1 >  ⚙

▶ category_wise_summary  •••

▼ daily_order_summary  •••

order_date  •••
Partition key

### Query 1  +

```
1 select * from daily_order_summary;
```

Run    Clear

# DynamoDB

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

Run    Clear

**Table view**    **JSON view**

⊘ Completed

**Started on 27/12/2024, 16:18:49**

**Elapsed time 268ms**

## Items returned (2)

Download results to CSV

🔍 Find items

< **1** >    ⚙

| order_date ▽ | unique_customers ▽ | return_rate ▽ | total_revenue |
|---|---|---|---|
| 2024-06-01 | 387 | 0.106965174... | 56816.6200499... |
| 2024-06-02 | 366 | 0.109947643... | 63619.1099960... |