

Data Pipeline Design Document

Table of Contents

1. [Introduction](#)

- Purpose of the document
- Scope

2. [System Architecture](#)

- High-Level Architecture Diagram
- Components and Technologies used
 - Google Cloud Storage
 - BigQuery
 - Google Cloud Composer (Airflow)
 - Other tools and services

3. [Pipeline Details \[pax, sales, loading\]](#)

- Source Data
- Destination
- Steps and transformations

4. [Common Pipeline Components](#)

- Data Extraction using storage client
- Loading Data into BigQuery
- Data cleaning and transformations

5. [Micro Batching using Cloud functions](#)

- Deploying a cloud function
- Trigger based on storage event

6. [Infrastructure setup and configuration](#)

- IAM considerations
- External package dependencies on composer

7. [Leveraging Snowflake for Machine Learning Workloads](#) (optional)

8. [Summary and Appendices](#)



Orchestrating ETL data pipelines on Google Cloud Platform using Airflow and Python

1. Introduction

In this document, we'll dive into the design and implementation of our data pipelines for zeroG Airline. These pipelines handle the critical task of processing and transforming various data sources—like passenger information (pax), in-flight sales data, and goods loading information—into a structured and accessible format for our data scientists and analysts.

We'll explore each component, the workflow, and the technology stack we use to ensure efficient and reliable data processing.

Purpose:

The purpose of this document is to provide a comprehensive overview of our data pipeline architecture. It's meant to guide you through the steps we've taken to build and maintain these pipelines, ensuring you understand the underlying design principles and can effectively troubleshoot or expand the system as needed.

Whether you're new to data processing workflows on Google Cloud Platform, or looking to set-up batch pipelines using Airflow, this document will serve as your go-to reference.

Scope:

This document covers the end-to-end design and implementation of our data pipelines, specifically focusing on the ingestion, transformation, and loading of pax, sales, and loading data.

We'll discuss the use of Google Cloud Storage (GCS) | Google Cloud Composer (Airflow) | BigQuery | Cloud Functions | and Snowflake.

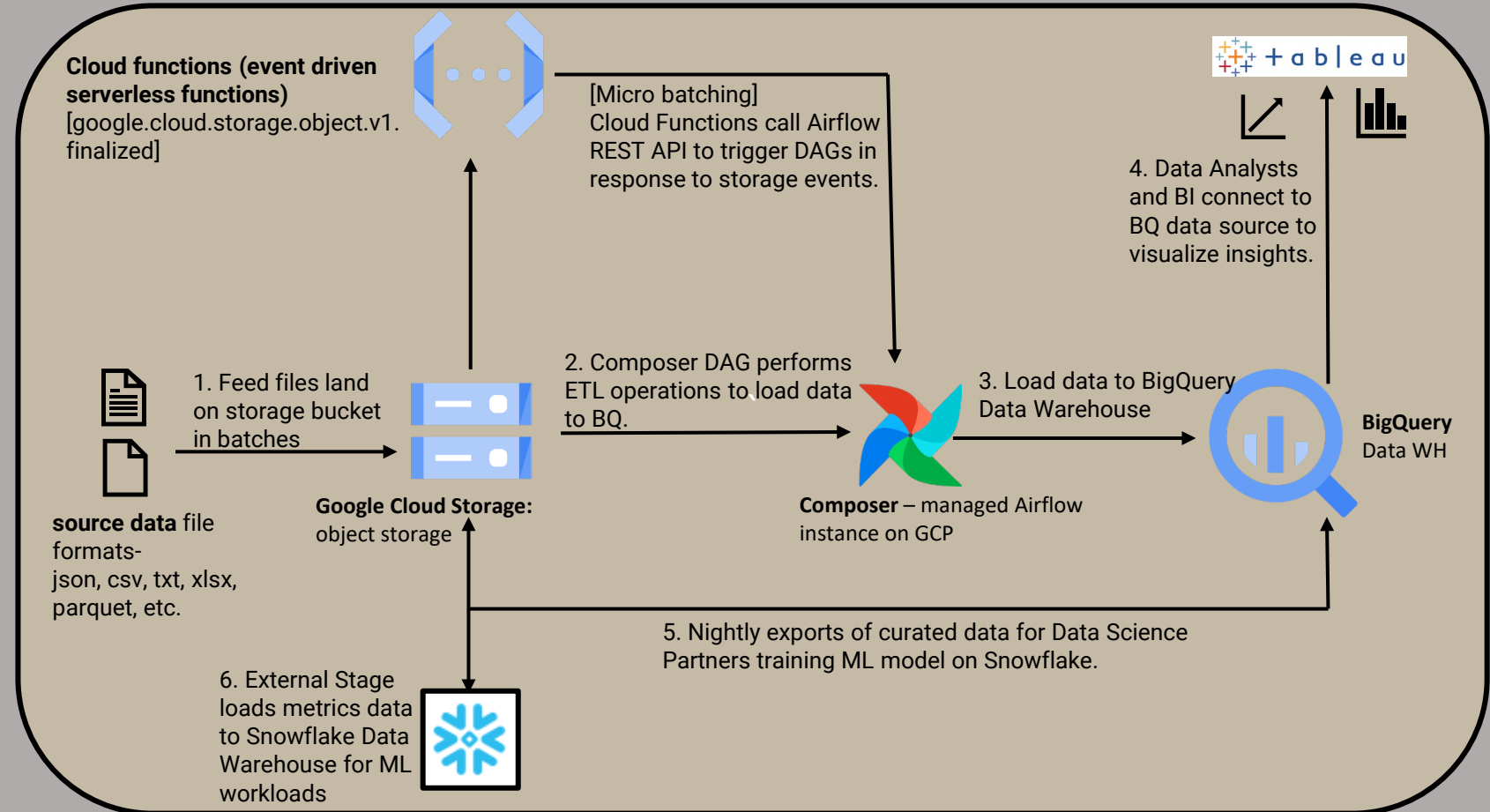
- The scope includes the daily and event-driven workflows.
- The data transformation processes using SQL queries.
- And the export mechanisms that integrate with Snowflake.

By the end of this document, you should have a clear understanding of how these components work together to deliver clean, reliable data to our stakeholders.

2. System Architecture

Process Flow:

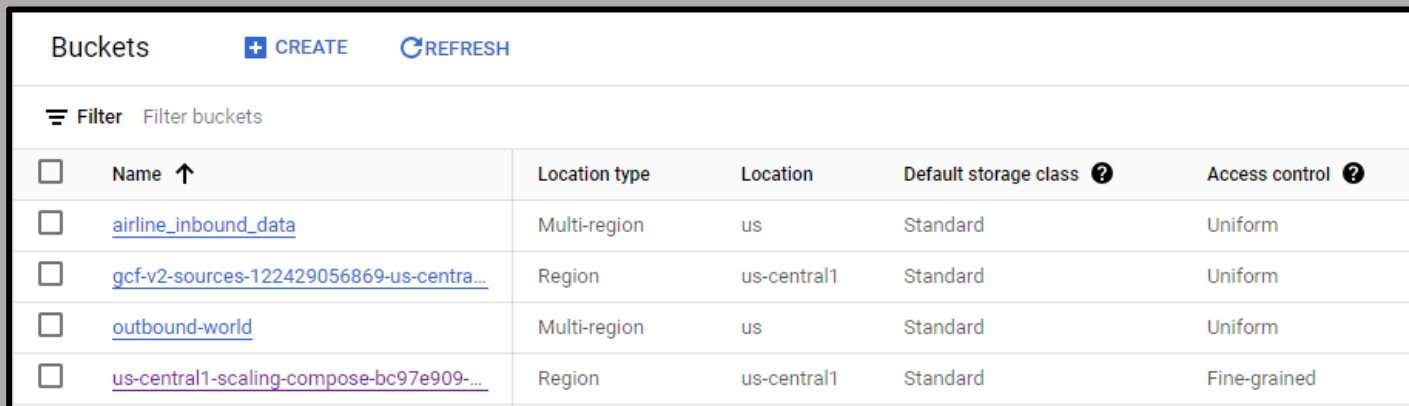
- Every day, raw data files are uploaded to our Google Cloud Storage (GCS) buckets for sales, pax, and loading data.
- The nightly scheduled DAGs orchestrate the entire ETL process, extracting data from GCS, transforming it using SQL queries in BigQuery, and loading it into BigQuery tables.
- The transformation steps include cleaning the data and performing necessary calculations to prepare it for analysis.
- Once the data is processed, it's exported back to GCS if needed, and separate Snowflake pipeline loads this curated data into warehouse.
- Throughout this process, Cloud Logging keeps track of all activities, making it easy to monitor and debug.



High Level Architecture Diagram

Components and Technologies Used:

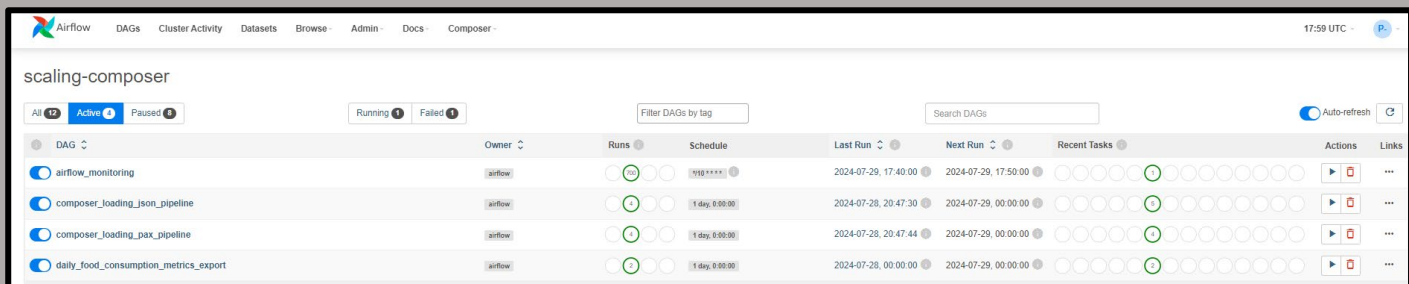
- Google Cloud Storage (GCS):
Stores raw data files (CSV, Excel, JSON) for sales, pax, and loading data.
- Google Cloud Composer (Airflow):
Orchestrates ETL processes, managing data workflows and scheduling DAGs.
- BigQuery:
Processes and transforms data using SQL, enabling efficient querying and analysis.
- Cloud Functions:
Handles ad hoc micro batching by triggering DAGs upon file arrival in GCS.
- Cloud Logging:
Monitors and logs all activities, aiding in debugging and performance tracking.
- Python:
Used in Airflow DAGs for scripting and data manipulation tasks.
- SQL:
Employed in BigQuery for data transformation, cleaning, and querying.



The screenshot shows the Google Cloud Storage Buckets interface. At the top, there are buttons for '+ CREATE' and 'REFRESH'. Below is a 'Filter' section with a 'Filter buckets' input. The main table lists several buckets with columns for Name, Location type, Location, Default storage class, and Access control. The buckets listed are 'airline_inbound_data', 'gcf-v2-sources-122429056869-us-centra...', 'outbound-world', and 'us-central1-scaling-compose-bc97e909-...'.

<input type="checkbox"/>	Name ↑	Location type	Location	Default storage class ?	Access control ?
<input type="checkbox"/>	airline_inbound_data	Multi-region	us	Standard	Uniform
<input type="checkbox"/>	gcf-v2-sources-122429056869-us-centra...	Region	us-central1	Standard	Uniform
<input type="checkbox"/>	outbound-world	Multi-region	us	Standard	Uniform
<input type="checkbox"/>	us-central1-scaling-compose-bc97e909-...	Region	us-central1	Standard	Fine-grained

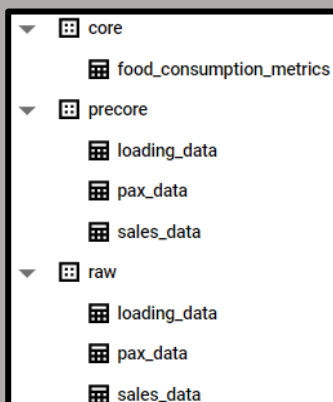
Cloud Storage buckets for : Daily source feeds | Artifact Registry | Data exports to external systems | composer bucket for dags, plugins, and logs.



The screenshot shows the Airflow DAGs interface. At the top, there are tabs for 'DAGs', 'Cluster Activity', 'Datasets', 'Browse', 'Admin', 'Docs', and 'Composer'. Below is a 'scaling-composer' section with a 'Filter DAGs by tag' input and a 'Search DAGs' input. The main table lists several DAGs with columns for DAG, Owner, Runs, Schedule, Last Run, Next Run, Recent Tasks, and Actions. The DAGs listed are 'airflow_monitoring', 'composer_loading_json_pipeline', 'composer_loading_pax_pipeline', and 'daily_food_consumption_metrics_export'.

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions
airflow_monitoring	airflow	1	*/15 * * * *	2024-07-29, 17:40:00	2024-07-29, 17:50:00	1	▶ ⏏
composer_loading_json_pipeline	airflow	1	1 day, 0:00:00	2024-07-28, 20:47:30	2024-07-29, 00:00:00	1	▶ ⏏
composer_loading_pax_pipeline	airflow	1	1 day, 0:00:00	2024-07-28, 20:47:44	2024-07-29, 00:00:00	1	▶ ⏏
daily_food_consumption_metrics_export	airflow	1	1 day, 0:00:00	2024-07-28, 00:00:00	2024-07-29, 00:00:00	1	▶ ⏏

Daily batch pipelines on **Airflow** for pax | sales | loading ETL workloads.



The screenshot shows the BigQuery datasets interface. It displays a hierarchical view of datasets under the 'core' schema. The datasets are organized into three categories: 'core', 'precore', and 'raw'. Each category contains three datasets: 'loading_data', 'pax_data', and 'sales_data'.

Dataset
core
food_consumption_metrics
precore
loading_data
pax_data
sales_data
raw
loading_data
pax_data
sales_data

BigQuery datasets

raw: truncate and load tables for pax | sales | loading

precore:
preprocessed and cleaned data

core:
curated layer to calculate metrics

Eventarc trigger ⓘ

Name

[trigger-sales-pipeline-064267](#)

Event provider

Cloud Storage

Event type

google.cloud.storage.object.v1.finalized

Receive events from

[airline_inbound_data](#) (us)

Service account

[122429056869-compute@developer.gserviceaccount.com](#)

Retry on failure

Disabled

The **EventArc** trigger named 'trigger-sales-pipeline-064267' is set up to activate the 'trigger-sales-pipeline' **Cloud Function** whenever a new file is finalized/created in the airline_inbound_data bucket. It uses the default compute service account to handle these events from Google Cloud Storage.

Eventarc lets you asynchronously deliver events between decoupled microservices while managing delivery, security, authorization, observability, and error-handling for you.

Cloud Logging plays a crucial role in monitoring and troubleshooting by capturing detailed logs of Cloud Function executions. For example, logs show requests to the 'triggersalespipeline' function, including HTTP request details, response sizes, and execution latencies, which helps ensure the function is operating correctly and identifies potential issues.

- Logs provide information on HTTP request methods, status codes, and response sizes for monitoring function performance.
- Resource labels identify the specific Cloud Run service and revision associated with the function, aiding in pinpointing issues.

The screenshot displays the Google Cloud Logs Explorer interface. At the top, there's a 'Logs Explorer' header with a 'Refine scope' button and a 'Project' dropdown. Below this, a 'Query' section shows a filter: `(resource.type = "cloud_function" AND resource.labels.function_name = "trigger_sales_pipeline" AND resource.labels.region = "us-central1")`. The interface includes a search bar, filters for 'All resources', 'All log names', and 'Default+', along with a '+1 filter' button and a 'Show query' toggle. The main area shows a timeline view with a search bar and a 'Log fields' sidebar. The timeline shows a series of log entries, with the selected entry expanded to show details. The log entry is a POST request to the trigger-sales-pipeline function, with a status of 200, a response size of 130 B, and a latency of 606 ms. The log entry is expanded to show the full HTTP request details, including the request URL, request method, request size, response size, server IP, status, and user agent.

Log fields

Search fields and values

RESOURCE TYPE

- Cloud Run Revision 14
- Cloud Function 2

SEVERITY

- Info 12
- Notice 4

Timeline

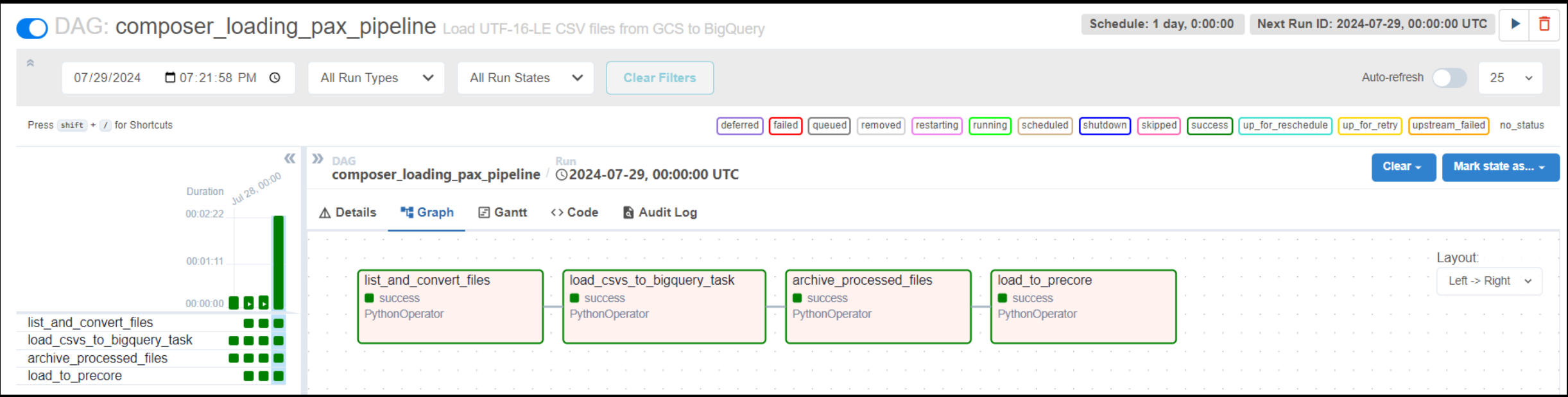
16 results

SEVERITY	TIME	IST	SUMMARY
Info	2024-07-28 20:43:35.061		run.googleapis.com v1 _30409-b3/services/trigger-sales-pipeline Ready condition status changed to True for Service trigger-sales-
Info	2024-07-28 20:43:51.822		cloudfunctions.googleapis.com _s.v2.FunctionService.CreateFunction _entral1/functions/trigger_sales_pipeline prashant.dixit2307@ma
Info	2024-07-28 20:47:54.896		POST 200 130 B 606 ms APIs-Google; https://trigger-sales-pipeline-31mc6ruhsq-uc.a.run.app/?__GCP_CloudEventsMode=GCS_NOTIFICATION

Copy Similar entries Collapse nested fields Hide log summary

```
{
  httpRequest: {
    latency: "0.606683273s"
    protocol: "HTTP/1.1"
    remoteIp: "64.233.172.161"
    requestMethod: "POST"
    requestSize: "3198"
    requestUrl: "https://trigger-sales-pipeline-31mc6ruhsq-uc.a.run.app/?__GCP_CloudEventsMode=GCS_NOTIFICATION"
    responseSize: "130"
    serverIp: "216.239.34.53"
    status: 200
    userAgent: "APIs-Google; (+https://developers.google.com/webmasters/APIs-Google.html)"
  }
}
```

3. Pipeline Details [pax, sales, loading]



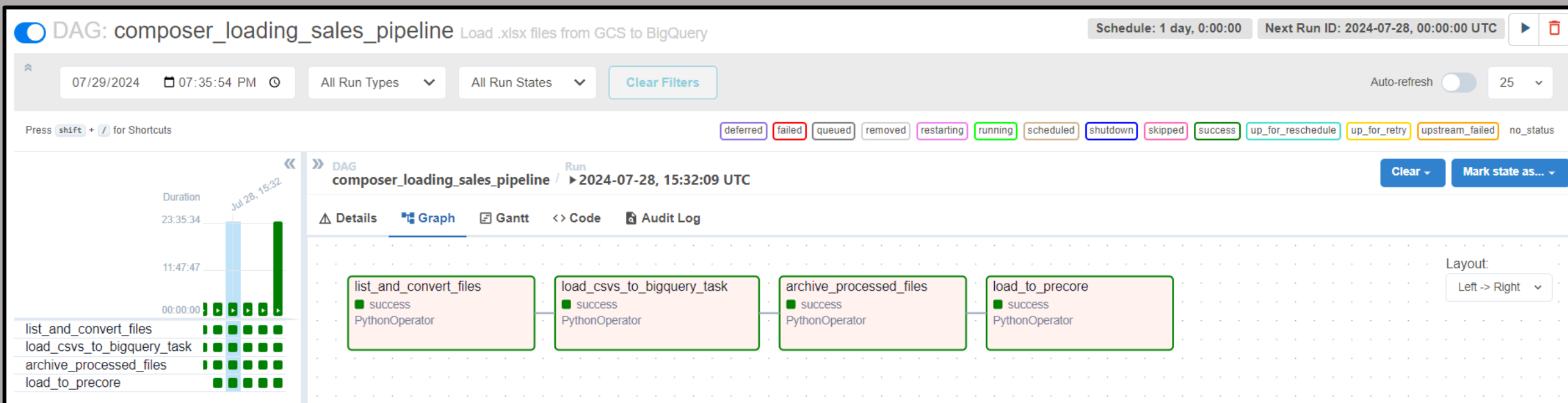
Source Data:
CSV files stored in GCS bucket ``gs://airline_inbound_data/pax/`` with utf-16-le encoding.
File pattern: `pax_%m_%Y.csv` [for ex: `pax_03_2019.csv` and `pax_04_2019.csv`]

Destination:
BigQuery table: ``sunlit-analyst-430409-b3.raw.pax_data``.

Steps and Transformations:
Loading: CSV files are read directly from GCS using Google Cloud Storage Client in Composer DAGs.

Transformations: Data is cleaned and transformed using SQL queries in BigQuery. The transformation process includes converting the UTF16LE encoded text to a readable format and handling data type conversions.

Load Mode: Data is loaded into raw BigQuery tables in `write_truncate` mode, then processed to move transformed data into precore tables, including adding a timestamp field.



Source Data:

Excel files stored in GCS bucket `gs://airline_inbound_data/sales/`,
File pattern: Sales_Report_%B_%Y.xlsx [Sales_Report_April_2019.xlsx, Sales_Report_March_2019.xlsx]

Destination:

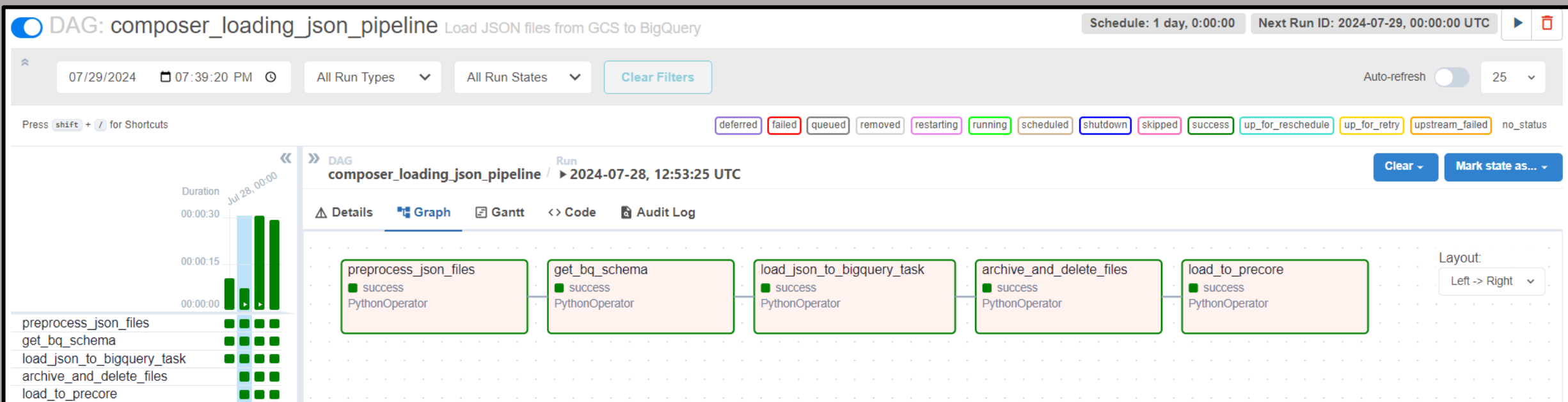
BigQuery table: `sunlit-analyst-430409-b3.raw.sales_data`.

Steps and Transformations:

Loading: Excel files are loaded from GCS, with processing to handle Excel specific challenges, such as formatting issues and data type recognition.

Transformations: Data is transformed using SQL queries in BigQuery (data cleaning, formatting adjustments, and conversion into a consistent structure)

Load Mode: Data is initially loaded into raw BigQuery tables in write_truncate mode and then moved to precore tables with additional timestamp fields.



Source Data:

JSON files stored in GCS bucket `gs://airline_inbound_data/loading/``.
File pattern: `loading_%m_%Y.json` [`loading_03_2019.json`, `loading_04_2019.json`].

Destination:

BigQuery table: ``sunlit-analyst-430409-b3.raw.loading_data``

Steps and Transformations:

Loading: JSON files are read from GCS with special handling for problematic keys and escape characters, ensuring proper parsing.

Transformations: Data is cleaned and structured using SQL queries in BigQuery. Transformation involves resolving issues with JSON formatting, such as handling spaces in keys.

Load Mode: Data is first loaded into raw BigQuery tables in write truncate mode, then processed to move transformed data into ``precore`` tables, including adding a timestamp field.

4. Common Pipeline Components

Data Extraction using storage client

In the pipelines, I've used the Google Cloud Storage Client Library to efficiently interact with GCS buckets. By importing the library and instantiating a storage client, the process of accessing and managing files is streamlined. This allows to easily read data from various sources, such as CSVs, Excel files, and JSON files, directly from GCS, facilitating smooth data extraction and integration into BigQuery. The client library simplifies authentication and reduces the amount of code needed, making it easier to handle different file formats and manage data across the pipelines.

```
def preprocess_json_files(bucket_name, input_folder, output_folder):
    client = storage.Client()
    bucket = client.get_bucket(bucket_name)
    blobs = bucket.list_blobs(prefix=input_folder)

    for blob in blobs:
        if blob.name.endswith('.json'):
            print(f"Processing file: {blob.name}")
            json_file_path = f'/tmp/{os.path.basename(blob.name)}'
            blob.download_to_filename(json_file_path)

            with open(json_file_path, 'r', encoding='utf-8') as f:
                records = json.load(f)

            processed_records = [map_fields(record) for record in records]
            processed_json_file_path = f'/tmp/processed_{os.path.basename(blob.name)}'

            with open(processed_json_file_path, 'w', encoding='utf-8') as f:
                for record in processed_records:
                    f.write(json.dumps(record) + '\n')

            processed_blob = bucket.blob(output_folder + os.path.basename(processed_json_file_path))
            processed_blob.upload_from_filename(processed_json_file_path)
            print(f"Uploaded processed file to: {processed_blob.name}")
```

Snippet from loading pipeline - mapping json records as per table schema

```
def list_and_convert_csv_files(bucket_name, input_folder, converted_folder):
    print("Starting to list and convert CSV files from the GCS bucket...")
    client = storage.Client()
    bucket = client.get_bucket(bucket_name)
    blobs = client.list_blobs(bucket_name, prefix=input_folder)
    csv_files = [blob.name for blob in blobs if blob.name.endswith('.csv')]

    for csv_file in csv_files:
        print(f"Converting {csv_file} to UTF-8 encoding...")
        blob = bucket.blob(csv_file)
        csv_file_path = f'/tmp/{os.path.basename(csv_file)}'
        blob.download_to_filename(csv_file_path)
        df = pd.read_csv(csv_file_path, encoding='utf-16-le')

        converted_file_path = f'/tmp/{os.path.basename(csv_file)}'
        df.to_csv(converted_file_path, index=False, encoding='utf-8')
        print(f"Converted {csv_file} to UTF-8 CSV at {converted_file_path}")

        converted_file = csv_file.replace(input_folder, converted_folder)
        converted_blob = bucket.blob(converted_file)
        converted_blob.upload_from_filename(converted_file_path)
        print(f"Uploaded converted CSV file to GCS at {converted_file}")

    return csv_files
```

Snippet from pax pipeline – utf-16-le encoded csv processing

Loading Data into BigQuery

Here I've used Airflow's 'GCSToBigQueryOperator' to load the data into BigQuery raw tables.

The schema to be used for the BigQuery table may be specified in one of two ways. We may either directly pass the schema fields in or point the operator to a Google Cloud Storage object name. The object in Google Cloud Storage must be a JSON file with the schema fields in it.

We are loading multiple objects from a single bucket using the 'source_objects' parameter.

JSON Load Method:

- Uses NEWLINE_DELIMITED_JSON format for JSON files.
- Schema is retrieved from previous task via XCom and applied during the load.
- WRITE_TRUNCATE option overwrites existing data in the BigQuery raw table.

```
def load_to_bq_callable(**kwargs):
    schema = kwargs['ti'].xcom_pull(key='bq_schema', task_ids='get_bq_schema')
    load_to_bq_task = GCSToBigQueryOperator(
        task_id='load_json_to_bigquery',
        bucket=BUCKET_NAME,
        source_objects=[f'{PROCESSED_JSON_FOLDER}*.json'],
        destination_project_dataset_table=f'{BQ_PROJECT_NAME}.{BQ_DATASET_NAME}.{BQ_TABLE_NAME}',
        schema_fields=schema,
        source_format='NEWLINE_DELIMITED_JSON',
        write_disposition='WRITE_TRUNCATE',
        autodetect=False,
        dag=dag,
    )
    load_to_bq_task.execute(context=kwargs)
    print("Loading JSON files to BigQuery...")
```

```
def load_to_bq_callable(**kwargs):
    schema = fetch_schema_from_bigquery(BQ_PROJECT_NAME, BQ_DATASET_NAME, BQ_TABLE_NAME)
    load_to_bq_task = GCSToBigQueryOperator(
        task_id='load_csvs_to_bigquery',
        bucket=BUCKET_NAME,
        source_objects=[f'{CONVERTED_FOLDER}*.csv'],
        destination_project_dataset_table=f'{BQ_PROJECT_NAME}.{BQ_DATASET_NAME}.{BQ_TABLE_NAME}',
        schema_fields=schema,
        source_format='CSV',
        skip_leading_rows=1,
        write_disposition='WRITE_TRUNCATE',
        autodetect=False,
        dag=dag)
    load_to_bq_task.execute(context=kwargs)
    print("Loading CSV files to BigQuery...")
```

CSV Load Method:

- Uses CSV format for CSV files.
- Skips the header row with skip_leading_rows=1.
- Schema is fetched directly from BigQuery, and WRITE_TRUNCATE option is used for raw loads.

Data cleaning and transformations

In our data pipelines for `pax`, `sales`, and `loading`, data cleaning and transformation processes were crucial to ensure data quality and integrity. During initial data profiling, we identified several instances of junk data in the raw tables, such as invalid dates, unexpected text entries, and inconsistencies in data formats. To address these issues and prepare the data for analytical use, we implemented specific cleaning and transformation steps for each pipeline.

For the pax pipeline, we found that the `Date` column, expected to contain dates in the format 'DD.MM.YYYY', occasionally had invalid entries such as ('Gesamtsumme'). We filtered these out using a regular expression to match the correct date format, ensuring only valid date entries were included in the `precore` table.

In the sales pipeline, we encountered issues with various data formats and inconsistencies within the JSON files. We employed a series of transformations to parse these JSON files correctly, handling escape characters and spaces in keys.

The loading pipeline had similar challenges, particularly with the `Flight_Month_Year` column, which sometimes contained irrelevant text entries like 'ZG Orders and Invoices' instead of the expected month-year values (e.g., '42019', '82024').

```
def load_to_precore_callable(**kwargs):
    client = bigquery.Client()
    query = f"""
    INSERT INTO `{PRECORE_PROJECT_NAME}`.{PRECORE_DATASET_NAME}.{PRECORE_TABLE_NAME}
    SELECT *, CURRENT_TIMESTAMP() as ingestion_time
    FROM `{BQ_PROJECT_NAME}`.{BQ_DATASET_NAME}.{BQ_TABLE_NAME}`
    WHERE REGEXP_CONTAINS(Flight_Month_Year, r'^\d{1,2}\d{4}$')
    """
    client.query(query).result()
    print("Loaded data into precore table with an additional timestamp.")
```

Loading pipeline - load cleaned data to precore layer

```
def load_to_precore_callable(**kwargs):
    client = bigquery.Client()
    query = f"""
    INSERT INTO `{PRECORE_PROJECT_NAME}`.{PRECORE_DATASET_NAME}.{PRECORE_TABLE_NAME}`
    SELECT *, CURRENT_TIMESTAMP() as ingestion_time
    FROM `{BQ_PROJECT_NAME}`.{BQ_DATASET_NAME}.{BQ_TABLE_NAME}`
    WHERE REGEXP_CONTAINS(Date, r'^\d{2}\.\d{2}\.\d{4}$')
    """
    client.query(query).result()
    print("Loaded data into precore table with an additional timestamp.")
```

pax pipeline - load cleaned data to precore layer

5. Micro Batching using Cloud functions

Deploying a cloud function

In our pipelines, we used Cloud Functions to trigger Cloud Composer DAGs whenever storage created/finalized events occurred. By following GCP documentation, we set up Cloud Functions to monitor changes in a Cloud Storage bucket. When a new file was added or an existing one was updated, the Cloud Function triggered the corresponding DAG through the Airflow REST API. This method allowed us to process data promptly and keep our pipelines responsive.

- Enabled necessary APIs, including the Airflow REST API.
- Configured Cloud Functions to trigger on changes in the Cloud Storage bucket.
- Uploaded and set up DAGs to be triggered by these functions.
- Ensured API access by configuring Webserver Access Control.
- Tested by uploading files and checking the Airflow interface for results.

requirements.txt

```
functions-framework==3.*  
google-auth==2.19.1  
requests==2.32.2
```

composer2_airflow_rest_api.py

Pass below url for composer web server and dag_id we want to trigger from main.py

```
web_server_url =  
("https://c67920cb09d149e5885eb  
bd0cc6c6b07-dot-us-  
central1.composer.googleusercontent.com")
```

```
dag_id =  
'composer_loading_sales_pipeline'
```

The screenshot displays the Google Cloud Functions interface for a function named 'trigger_sales_pipeline' (2nd generation). The function is deployed at 8:43:51 PM on July 28, 2024, with a URL pointing to the Cloud Functions endpoint. The 'SOURCE' tab is active, showing the Python code for the function. The code imports 'google.auth', 'google.auth.transport.requests', and 'requests'. It defines 'AUTH_SCOPE' and 'CREDENTIALS'. The 'make_composer2_web_server_request' function is defined to make a GET request to the Airflow REST API. The 'trigger_dag' function is defined to trigger a DAG by making a POST request to the Airflow REST API. The main logic is in the 'trigger_sales_pipeline' function, which calls 'make_composer2_web_server_request' and 'trigger_dag'.

```
3 import google.auth  
4 from google.auth.transport.requests import AuthorizedSession  
5 import requests  
6  
7 AUTH_SCOPE = "https://www.googleapis.com/auth/cloud-platform"  
8 CREDENTIALS, _ = google.auth.default(scopes=[AUTH_SCOPE])  
9  
10 def make_composer2_web_server_request(url: str, method: str = "GET", **kwargs: Any) -> google.auth.transport.Response:  
11     auth_session = AuthorizedSession(CREDENTIALS)  
12     if "timeout" not in kwargs:  
13         kwargs["timeout"] = 90  
14     return auth_session.request(method, url, **kwargs)  
15  
16 def trigger_dag(web_server_url: str, dag_id: str, data: dict) -> str:  
17     """Make a request to trigger a dag using the stable Airflow 2 REST API. https://airflow.apache.org/docs/apache-airflow/stable/stable-rest-api-ref.html  
18     Args: web_server_url: The URL of the Airflow 2 web server | dag_id: The DAG ID. | data: Additional configuration parameters for the DAG run (json)."""  
19  
20     endpoint = f"api/v1/dags/{dag_id}/dagRuns"  
21     request_url = f"{web_server_url}/{endpoint}"  
22     json_data = {"conf": data}  
23  
24     response = make_composer2_web_server_request(request_url, method="POST", json=json_data)  
25  
26     if response.status_code == 403:  
27         raise requests.HTTPError(  
28             "You do not have a permission to perform this operation. Check Airflow RBAC roles for your account."  
29             f"{response.headers} / {response.text}"  
30         )  
31     elif response.status_code != 200:  
32         response.raise_for_status()  
33     else:  
34         return response.text
```

Trigger based on Storage Event

This JSON message is a log entry generated when a Cloud Storage event triggers a Cloud Run service, specifically for the "trigger-sales-pipeline."

Key Points:

1. Event Source: A POST request from the Google Cloud Storage event notification system.
2. Target: The Cloud Run service "trigger-sales-pipeline" in the "us-central1" region.
3. Event Details: The event was logged at "2024-07-28T15:17:54.896939Z" with a status code of 200 (success).
4. Context: The request triggered the Cloud Function configured to handle **sales** pipeline events.

Summary: This log entry records a successful trigger of the sales pipeline via a Cloud Storage event notification, processed by a Cloud Run service in response to a storage event.

```
{
  "insertId": "66a661230007f69a4cd11adb",
  "httpRequest": {
    "requestMethod": "POST",
    "requestUrl": "https://trigger-sales-pipeline-3imc6ruhsq-uc.a.run.app/?_GCP_CloudEventsMode=GCS_NOTIFICATION",
    "requestSize": "3198",
    "status": 200,
    "responseSize": "130",
    "userAgent": "APIs-Google; (+https://developers.google.com/webmasters/APIs-Google.html)",
    "remoteIp": "64.233.172.161",
    "serverIp": "216.239.34.53",
    "latency": "0.606683273s",
    "protocol": "HTTP/1.1"
  },
  "resource": {
    "type": "cloud_run_revision",
    "labels": {
      "location": "us-central1",
      "project_id": "sunlit-analyst-430409-b3",
      "service_name": "trigger-sales-pipeline",
      "configuration_name": "trigger-sales-pipeline",
      "revision_name": "trigger-sales-pipeline-00001-tux"
    }
  },
  "timestamp": "2024-07-28T15:17:54.896939Z",
  "severity": "INFO",
  "labels": {
    "instanceId": "0087244a8001d26599356b69828a21ec1062340ff618c8ae2a51edb3671c3dc6950579eaa862604f6614c07d2c92a84dde86c76b07d2b063bf8025e8c944aa7afc87",
    "goog-managed-by": "cloudfunctions"
  },
  "logName": "projects/sunlit-analyst-430409-b3/logs/run.googleapis.com%2Frequests",
  "trace": "projects/sunlit-analyst-430409-b3/traces/a89a3fcc5d39a9083737b61f6eec36a0",
  "receiveTimestamp": "2024-07-28T15:17:55.524780864Z",
  "spanId": "9770553742873758851",
  "traceSampled": true
}
```

6. Infrastructure setup and configuration

External package dependencies on Composer

- Cloud Composer images contains both preinstalled and custom PyPI packages. Preinstalled PyPI packages are packages that are included in the Cloud Composer image of our environment. Each Cloud Composer image contains PyPI packages that are specific for our version of Cloud Composer and Airflow.
- Custom PyPI packages are packages that we can install in our environment in addition to preinstalled packages.

IAM considerations:

Cloud Composer uses Identity and Access Management (IAM) for access control. We control access to different Cloud Composer features by granting roles and permissions both for IAM service accounts and for user accounts in your Google Cloud project. Cloud Composer uses two types of IAM service accounts:

- Cloud Composer Service Agent account
- Environment's service account

In addition to these two types of service account, Google APIs Service Agent runs internal Google processes on your behalf.

Composer

← Environment details

✓ scaling-composer

This environment is running

MONITORING

LOGS

DAGS

ENVIRONMENT

EDIT

Required libraries from the Python Package Index (PyPI)

Name	Version
pandas	-
fsspec	-
gcsfs	-
openpyxl	-
snowflake-connector-python	-
cryptography	-

Used default compute service account to manage permissions for Google Cloud Composer. This account has been granted permissions to access cloud storage, BigQuery and other GCP services.

For production environments, it's recommended that we set up a user-managed service account for Cloud Composer environments. Grant a role that is specific for Cloud Composer to this account. Afterwards, specify this service account when creating new environments.

VIEW BY PRINCIPALS

VIEW BY ROLES

GRANT ACCESS

REMOVE ACCESS

Filter

Enter property name or value

Type	Principal ↑	Name	Role	Security insights
<input type="checkbox"/>	<div>122429056869-compute@developer.gserviceaccount.com</div>	Compute Engine default service account	<div>Cloud Composer v2 API Service Agent Extension</div> <div>Composer Worker</div> <div>Eventarc Event Receiver</div>	

7. Leveraging Snowflake for Machine Learning Workloads (optional)

In a modern enterprise, data is generated and stored across multiple cloud platforms, and this fragmented data storage poses a challenge for data scientists who need to access, aggregate, and analyze data from various sources to develop and deploy machine learning (ML) models effectively.

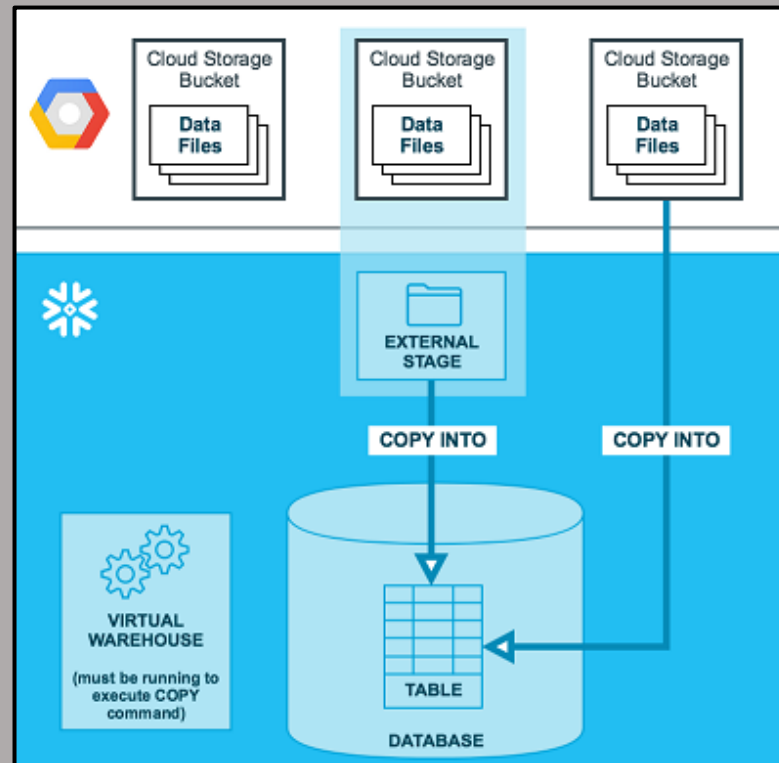
Why Snowflake as the Enterprise Cloud Data Warehouse ?

- Unified Data Platform (SaaS for all public clouds)
- Scalability and Performance Data (autoscaling, performance comes In T-shirt Sizes Like S, M, L, XL...)
- Sharing and Collaboration (modern features like Cloning, Data Shares, AI & ML Products, extensive range of connectors)
- Advanced Analytics Capabilities (in-built change data capture (CDC), serverless automations through Tasks & Streams, Partner Connect, BI, CI/CD)

To address the data integration challenge and enable effective ML workloads, we set up a pipeline that exports data from Google Cloud Storage (GCS) to Snowflake. This allows data scientists to access and analyze data stored in GCP seamlessly

To address the data integration challenge and enable effective ML workloads, we set up a pipeline that exports data from Google Cloud Storage (GCS) to Snowflake. This allows data scientists to access and analyze data stored in GCP seamlessly.

1. Data Aggregation in BigQuery
2. Export Data to GCS
3. Setup External Stage in Snowflake
4. Load Data into Snowflake
5. Data Transformation and Analysis in Snowflake
6. Machine Learning Workloads



As illustrated in the diagram below, loading data from a Cloud Storage bucket is performed in two steps:

Step 1

Snowflake assumes the data files have already been staged in a Cloud Storage bucket. If they haven't been staged yet, use the upload interfaces/utilities provided by Google to stage the files.

Step 2

Use the `COPY INTO <table>` command to load the contents of the staged file(s) into a Snowflake database table. You can load directly from the bucket, but Snowflake recommends creating an external stage that references the bucket and using the external stage instead.

←	snow_role	✎ EDIT ROLE	📄 CREATE FROM ROLE
ID	projects/sunlit-analyst-430409-b3/roles/CustomRole		
Role launch stage	Alpha		

Loading Data into a Snowflake from an external GCP bucket involves following

```
CREATE STORAGE INTEGRATION GCP_INTEGRATION
TYPE = EXTERNAL_STAGE
STORAGE_PROVIDER = 'GCS'
ENABLED = TRUE
STORAGE_ALLOWED_LOCATIONS = ('gcs://outbound-world/metrics');
```

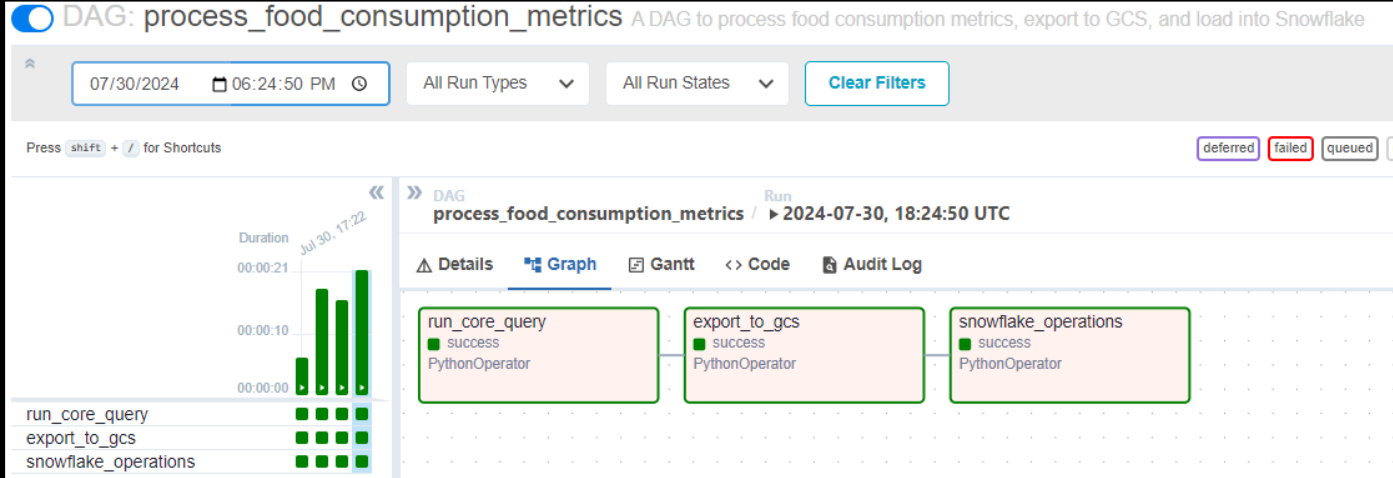
```
DESC STORAGE INTEGRATION GCP_INTEGRATION;  
--STORAGE_GCP_SERVICE_ACCOUNT=kooaoyasmn@azcentralindia-1-5514.iam.gserviceaccount.com
```

Authorize this service account on gcp IAM & admin portal with required roles - **snow_role** with specific permissions for storage listed in the image)

```
CREATE OR REPLACE STAGE APP_DB.APP_SCHEMA.GCS_STAGE
URL = 'gcs://outbound-world/metrics'
STORAGE_INTEGRATION = GCP_INTEGRATION;
```

```
COPY INTO AIRLINE.INSIGHTS.FOOD_CONSUMPTION_METRICS
FROM @GCS_STAGE
FILE_FORMAT =
  (TYPE = 'CSV'
   FIELD_OPTIONALLY_ENCLOSED_BY = '"'
   FIELD_DELIMITER = ';'
   SKIP_HEADER = 1)
PATTERN = '.*snow-exports_.*\\.csv';
```

	file	status	rows_parsed	rows_loaded	error_limit	errors_seen
1	gcs://outbound-world/metrics/snow-exports_07302024.csv	LOADED	711	711	1	0



Daily pipeline to load data to Snowflake from cloud storage

- # 1. Exports core table data to GCS staging location in outbound bucket
- # 2. Creates files with prefix_{date_part}.csv pattern

```
def export_to_gcs():  
    from google.cloud import storage  
    import datetime  
  
    client = bigquery.Client()  
    bucket_name = 'outbound-world'  
    file_name = f'snow-exports_{datetime.datetime.now().strftime("%m%d%Y")}.csv'  
  
    query = "SELECT * FROM `sunlit-analyst-430409-b3.core.food_consumption_metrics`"  
    df = client.query(query).to_dataframe()  
  
    local_file_path = '/tmp/metrics.csv'  
    df.to_csv(local_file_path, index=False)  
  
    storage_client = storage.Client()  
    bucket = storage_client.bucket(bucket_name)  
    blob = bucket.blob(f'metrics/{file_name}')  
    blob.upload_from_filename(local_file_path)
```

This daily pipeline calculates the metric in core layer, exports the data to cloud storage bucket, and then finally the data is written to snowflake db table `AIRLINE.INSIGHTS.FOOD_CONSUMPTION_METRICS` using an external stage.

```
def get_snowflk_ctx():  
    local_rsa_key_path = '/tmp/rsa_key.p8'  
    rsa_key_path = 'gs://airline_inbound_data/requirements/rsa_key.p8'  
    storage_client = storage.Client()  
    bucket = storage_client.bucket('airline_inbound_data')  
    blob = bucket.blob('requirements/rsa_key.p8')  
    blob.download_to_filename(local_rsa_key_path)  
  
    with open(local_rsa_key_path, "rb") as key:  
        p_key = serialization.load_pem_private_key(key.read(), password=None, backend=default_backend())  
  
    pkb = p_key.private_bytes(encoding=serialization.Encoding.DER,  
                             format=serialization.PrivateFormat.PKCS8,  
                             encryption_algorithm=serialization.NoEncryption())  
  
    ctx = snowflake.connector.connect(  
        user='VISUALINSIGHTS',  
        password='DataMatters123',  
        account='gtor200-yw14931',  
        private_key=pkb,  
        warehouse='COMPUTE_WH',  
        database='APP_DB',  
        schema='APP_SCHEMA')  
    return ctx  
  
conn = get_snowflk_ctx()  
  
copy_sql = f"""  
COPY INTO AIRLINE.INSIGHTS.FOOD_CONSUMPTION_METRICS  
FROM @GCS_STAGE  
FILE_FORMAT = (TYPE = 'CSV' FIELD_OPTIONALLY_ENCLOSED_BY = '"' FIELD_DELIMITER = ',' SKIP_HEADER = 1)  
PATTERN = '.*snow-exports_.*\\.csv';  
"""  
  
df = pd.read_sql_query(copy_sql, conn)  
print(df)  
conn.close()
```

Streamlit is an open-source Python library that makes it easy to create and share custom web apps for machine learning and data science. By using Streamlit we can quickly build and deploy powerful data applications. Streamlit in Snowflake helps developers securely build, deploy, and share Streamlit apps on Snowflake’s data cloud. Using Streamlit in Snowflake, we can build applications that process and use data in Snowflake without moving data or application code to an external system.

<

Streamlit Apps

Airline_ML_Workloads ▾

Share

▶ Run

⋮

Packages ▾

```
1 # Import required packages
2 import streamlit as st
3 from snowflake.snowpark.context import get_active_session
4 from snowflake.snowpark.functions import col
5 import pandas as pd
6
7 # Set the title of the Streamlit app
8 st.title("Food Consumption Metrics Dashboard")
9
10 # Initialize the Snowflake session
11 session = get_active_session()
12
13 # Define a function to query data from Snowflake
14 def get_food_consumption_data():
15     query = session.table("AIRLINE.INSIGHTS.FOOD_CONSUMPTION_METRICS")
16     df = query.select(
17         col("FLIGHTNUMBER"),
18         col("FLIGHTDATE"),
19         col("TOTALPASSENGERS"),
20         col("TOTALSOLDITEMS"),
21         col("TOTALSALES"),
22         col("AVERAGESOLDITEMSPERPASSENGER"),
23         col("AVERAGESALESPERPASSENGER")
24     ).to_pandas()
25     return df
26
27 # Get the food consumption data from Snowflake
28 food_consumption_data = get_food_consumption_data()
29
30 # Display the data in a table
31 st.subheader("Food Consumption Metrics")
32 st.dataframe(food_consumption_data, use_container_width=True)
33
34 # Create a bar chart for TotalSoldItems
35 st.subheader("Total Sold Items per Flight")
36 st.bar_chart(data=food_consumption_data.rename(columns={"FLIGHTNUMBER": "FlightNumber", "TOTALSOLDITEMS": "TotalSoldItems"}))
37
38 # Create a bar chart for TotalSales
39 st.subheader("Total Sales per Flight")
40 st.bar_chart(data=food_consumption_data.rename(columns={"FLIGHTNUMBER": "FlightNumber", "TOTALSALES": "TotalSales"}))
41
42 # Create a bar chart for AverageSoldItemsPerPassenger
43 st.subheader("Average Sold Items per Passenger")
```

Total Sold Items per Flight

FlightNumber	TotalSoldItems
453	100
733	400
1497	100
1580	900
2235	100
2483	1100
2844	600
2964	500
3135	200
3634	200
3743	100
4413	100
5284	100
5472	200
5817	3800
5846	100
5888	100
6080	100
6330	100
6966	100
7039	1400
7327	100
7520	1100
7913	400
8335	100
8545	100
8647	900
9164	100
9331	100
9382	300
9931	200

Total Sales per Flight

FlightNumber	TotalSales
453	2000
733	12000
1497	2000
1580	30000
2235	2000
2483	38000
2844	22000
2964	18000
3135	8000
3634	5000
3743	1000
4413	1000
5284	1000
5472	5000
5817	112000
5846	1000
5888	1000
6080	2000
6330	2000
6966	5000
7039	42000
7327	5000
7520	38000
7913	12000
8335	3000
8545	1000
8647	30000
9164	5000
9331	5000
9382	8000
9931	8000

Feedback

Ⓢ

8. Summary and Appendices

Summary:

1. Data Loading: Loaded sales, passenger (pax), and loading data from Google Cloud Storage (GCS) to BigQuery tables for March and April 2019.
2. Data Transformation: Created SQL queries to transform data in BigQuery, including aggregating sales and passenger data and adding a timestamp field.
3. Pipeline Design: Developed Airflow DAGs in Google Cloud Composer for data ingestion, transformation, & loading into BigQuery.
4. Data Export: Built a pipeline to export processed metrics from BigQuery to GCS & integrated it with Snowflake for further analysis, handling RSA key connections.
5. Architecture: Infrastructure setup and configurations, discussed IAM considerations and external packaged dependencies and service accounts.
6. Future plans: Enhance the Streamlit app using Snowpark to visualize metrics data directly from Snowflake, keeping it separate from Composer workflows.

Reference Links:

Cloud Composer: <https://cloud.google.com/composer?hl=en>

Cloud Storage: <https://cloud.google.com/storage?hl=en>

Trigger DAGs with Cloud functions: <https://cloud.google.com/composer/docs/composer-3/triggering-with-gcf>

Snowflake-GCP integration: <https://docs.snowflake.com/en/user-guide/data-load-gcs-config>

Streamlit - <https://docs.snowflake.com/developer-guide/streamlit/about-streamlit>

This exercise is confidential and copyright (c) zeroG GmbH 2024. Please do not share this exercise, the solution, or any details with anyone.

Developed by Prashant Dixit for the ZeroG Data Engineering Case Study.

Thank you.