# Reuniting the two distant cousins: *Calling a Beam Pipeline from an Airflow Job*



Sadeeq Akintola

**X / Twitter:** @SadeeqAkintola
**Online:** SadeeqAkintola.com

BEΔM SUMMIT

**September 4-5, 2024**

**Sunnyvale, CA. USA**

# About the Speaker

## Sadeeq Akintola

- Customer Engineer, Data Analytics Specialist @ **Google** Cloud
- 13+ Years Industry Experience
- Previous job roles include: Software Engineer, Data Scientist & BI Analyst, Big Data Engineer
- Ex-Microsoft, ex-KPMG, ex-FMDQ Exchange
- Worked across multiple Geos: *Nigeria → Portugal → United Kingdom*
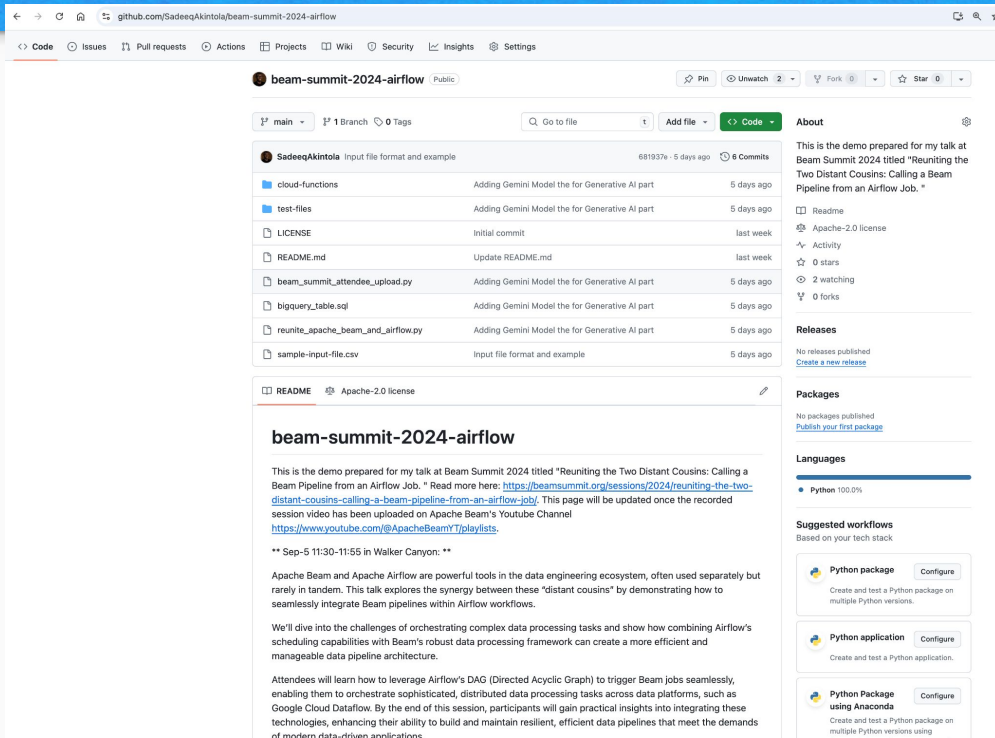- M.Sc. Data Science and Advanced Analytics @NovaIMS, Lisbon.

# Agenda

# Try the Demo, Clone the Repo.

# Motivation for this topic:



- Previous Experience as a Data Analytics Consultant
- Current Customer Interactions @ Google
- PyCon US Open Spaces Session
- People always love the idea of a Magic Wand product that can do everything, perfectly :)





5

# Understanding Important Concepts

# The Data Pipeline Platform

There are certain (other) requirements expected of a Modern Data Pipeline Platform:

| Schedule | Dependencies | Error Handling | Reporting | Connectors |

# Three Common Types of Data Pipelines

## Realtime

Realtime (or near realtime) data. An unbounded stream of data that it is desirable to process within a short period of time (mins → microseconds).

*Transport Example: Traffic sensors sending information about traffic flow. This may require quick action to be taken (e.g. open / close lanes)*

## Scheduled Batch

A bounded set of data which requires processing at regular intervals (e.g. every day, 4 hours, 5 minutes)

*Retail Example: Looking at purchasing patterns in store yesterday and comparing with online. Maybe looking for online to offline visits.*

## Triggered Batch

A bounded set of data which requires processing in a reactive manner when an event occurs (e.g. a data file is uploaded to a folder)

*Telecoms Example: Processing audio files from support call center as they are uploaded to storage. Running through voice to text conversion then processing for keywords / sentiment.*

# Data Pipeline Tools

Popular tools on Google Cloud used to Processing each Pipelines:

### Realtime

Cloud Pub/Sub

Cloud Dataflow

### Scheduled Batch

Cloud Dataflow

Apache Beam

App Engine Scheduled Tasks

### Triggered Batch

Cloud Composer

Apache Airflow

Cloud Functions

*Note: There are other tools in the open source space such as: Kafka, RabbitMQ, Luigi, Oozie, Azkaban, Hadoop etc. A combination of two or more of these might equally be suitable, depending on the use case.*

# Challenges in Orchestrating Data Processing Tasks

## Organizing Data can be such a pain…

- **Complexity of Data Workflows**
  - Managing dependencies and data flow between tasks.
  - Ensuring data consistency and integrity across distributed systems.

- **Scalability Concerns**
  - Handling large volumes of data efficiently.
  - Scaling resources dynamically based on workload.

- **Error Handling and Monitoring**
  - Detecting and recovering from failures.
  - Monitoring pipeline performance and resource utilization.

**Orchestration Need:** A pipeline orchestrator manages scheduling, monitoring, and dependencies, ensuring smooth data flow.

# Apache BEAM (Batch + StrEAM)

# Apache Beam: Core Concepts

- Open source, unified model for batch and streaming data pipelines

- Using one of the open source SDKs, you can build a program that defines the pipeline

- The pipeline is then executed by one of Beam's supported runners - Apache Apex, Flink, Spark or Google Cloud Dataflow

**One Model**   **Multiple Modes**   **Multiple SDKs**   **Multiple Runners**

Batch

Streaming

Java

Python

**Direct**: local for testing

**Cloud Dataflow**: fully managed service on Google Cloud

**Apache Flink**: local, on-premise, cloud

**Apache Spark**: local, on-premise, cloud

**Apache Apex**: local, on-premise, cloud

Google Cloud

BEAM SUMMIT

# Apache Beam: Core Concepts

- **Pipeline:** encapsulates the entire series of computations involved in reading input data, transforming that data, and writing output data.

- **PCollection:** represents a potentially distributed, multi-element dataset that acts as the pipeline's data. Beam transforms use PCollection objects as inputs and outputs for each step in your pipeline.

- **PTransforms:** A transform represents a processing operation that transforms data. A transform takes one or more PCollections as input, performs an operation that you specify on each element in that collection, and produces one or more PCollections as output.

- **I/O Sinks and Sources** – The Source and Sink APIs provide functions to read data into and out of collections. The sources act as the roots of the pipeline and the sinks are the endpoints of the pipeline.

Input → Transform → PCollection → Transform → PCollection → Transform → Output

*Example pipeline   Apache Beam Programming Guide

# Pipelines

- In Beam, you structure your computation as a graph of transformations, which we call a **Pipeline**.

- Each box here is a transform performing massively parallel computation, which we call a **PTransform**.

- Each Transform of the Pipeline is *applied* on a **PCollection**; the result of apply() is another **PCollection**.

- Each arrow represents the data itself, being transmitted from one PTransform to the next, which we call a PCollection.



This is a **Source**

PTransform

Pipeline

PCollection
(bounded or unbounded)

This is a **Sink**

Google Cloud

# PCollection

A PCollection represents a distributed data set that can be **bounded** or **unbounded**:

- Bounded means we know the PCollection is finite, while unbounded means that it *might be infinite*, *it might be finite*, but we just don't know.

- A Directed Acyclic Graph of data transformations.

- Possibly unbounded collections of data flow on the edges.

- May include multiple sources and multiple sinks.

- Optimized and executed as a unit.

- PCollections are immutable.

# PTransforms

- **Element-wise Transformation:** let say you process element individually and do some transformation on it. For example, you have a record with the user id and transform the user id to an email address.
  - Most frequently used *Map Transform* is *ParDo* = "Parallel Do".

- **Aggregating Transform**, also called *reduce*. Where different elements are processed together.
  - The key primitive is the *GroupByKey* - which groups key-value pairs by key.
  - It takes multiple elements and combines them.

- **Composite transformation:** they're just compound operations of more primitive things. For example, you can have a *combine fn* that counts words and then extract the top-K elements.



**Element-Wise (map)**       **Aggregating (reduce)**       **Composite (reusable combinations)**

# Why run Beam on Google Cloud Dataflow?

**Cloud native, serverless, extensible solution for mission critical ingestion, ETL, and streaming analytics:**

- Fully-managed and auto-configured
  - Resource management: Spinning up and down the machines that process data.
  - Dynamic work rebalancing: Partition and spread the data so that all machines have work to do, all the time.
- Auto graph-optimized for best execution path

- Autoscaling mid-job: if the load goes up or down, adjust the infrastructure accordingly.

- Dynamic Work Rebalancing mid-job



**Dataflow: Platform Powered by Google + Rich Open Source Apache Beam SDK**

# Apache Airflow

# Apache Airflow: Core Concepts

- Apache Airflow is an open-source workflow management platform for data engineering pipelines.

- It started at Airbnb in October 2014 as a solution to manage the company's increasingly complex workflows.

- Apache Airflow is used for the scheduling and orchestration of data pipelines or workflows.

- Orchestration of data pipelines refers to the sequencing, coordination, scheduling, and managing of complex data pipelines from diverse sources.

# Apache Airflow: Workflow Principles

These are the key principles you need to know about when building an Airflow workflow

| |
|---|
| .py File |
| DAG |
| Arguments |
| Schedule |
| Tasks / Operators |
| Macros |
| Dependencies |

Each Airflow workflow is a python file that is placed in the dags folder where Airflow runs.

The python syntax used to build the workflow is very simple and makes use of the operators with simple arguments for each task.

Configuration as code instead of drag and drop UI.

# Apache Airflow: Core Concepts

- **Directed Acyclic Graphs (DAGs):**
  - Represents workflows as a collection of tasks with defined dependencies.
  - Enables workflows to be defined as code for better maintainability.
- **Tasks**:
  - Basic units of execution within a DAG.
  - Can perform various operations like data fetching and analysis.
- **Operators:**
  - Templates that define what a task does (e.g., BashOperator, PythonOperator).
  - Variety available for different use cases.
- **Hooks**:
  - Interfaces to external platforms and services.
  - Used by operators for tasks like database queries or API calls.
- **XComs:**
  - Mechanism for tasks to exchange small amounts of data.
  - Facilitates communication between tasks in a DAG.
- **Architecture:**
  - Scheduler: Manages task execution and DAG scheduling.
  - Web Server: Provides a UI for monitoring and managing workflows.
  - Metadata Database: Stores the state of tasks and workflows.
  - Executor: Executes tasks locally or on distributed systems.
- **Extensibility and Community:**
  - Highly extensible with custom operators and hooks.
  - Supported by a large, active open-source community.



HOW APACHE AIRFLOW WORKS

# Apache Airflow: Directed Acyclic Graphs (DAGs)

A DAG in Apache Airflow is a central concept that represents a workflow of tasks organized in a way that defines their execution order. Here's a detailed explanation:

● **Directed**: The connections between tasks have a direction, meaning Task A must complete before Task B starts, establishing a clear sequence of execution.

● **Acyclic**: means there are no loops. Once a task is executed, the workflow doesn't return to that task; it progresses forward.

● **Graph**:The DAG is essentially a collection of nodes (tasks) and edges (dependencies) that represent the workflow. It visualizes the entire pipeline of tasks from start to finish.

● **Components of a DAG**
　○　Dependencies
　○　Schedule
　○　Operators
　○　Execution

```
from airflow import DAG
dag = DAG('my-new-dag', schedule_interval='0 0 * * *')
```

22

# Why run Airflow on Google Cloud Composer?

## Apache Airflow

+ Python/open source
+ Well-established interfaces
  (*CLI / Webserver / API server*)
+ Flexible scheduling and dependency management with retries
+ Rich library of connectors
+ Active community support
- Non-trivial setup & management efforts
- Logging/debugging
+ Apache Airflow supports many executors:
local, Celery, Kubernetes, Mesos, Dask

benefits from

contributes to

## Google Cloud Platform

+ Ease of use
  + Integration with other GCP services
  + One-click deployment
  + Create/Update/Delete
  + UI/gcloud/API/Terraform
+ Infrastructure scalability (GKE, GAE, CloudSQL)
+ Logging and Monitoring
+ Security
  + Webserver: behind IAP, Network ACLs
  + Shared VPC, Private IP
  + VPC SC (Beta)

+ Operability & Maintainability
  + Python dependency management
  + Airflow config update propagation
  + In-place version upgrades (Beta)

- Some Airflow capabilities are not available

# Code Review (Beam Pipeline)

# Importing Libraries and Reading Files

```
1   import csv
2   import apache_beam as beam
3   from apache_beam.options.pipeline_options import PipelineOptions, GoogleCloudOptions, StandardOptions, SetupOptions
4   from apache_beam.io.fileio import MatchFiles, ReadMatches
5   from apache_beam.io.gcp.bigquery import WriteToBigQuery, BigQueryDisposition
6   import datetime
7   import logging
8
9   class ReadAndValidateCSV(beam.DoFn):
10      def process(self, file):
11          file_path = file.metadata.path
12          logging.info(f"Processing file: {file_path}")
13          try:
14              with beam.io.filesystems.FileSystems.open(file_path) as f:
15                  text_io = f.read().decode('utf-8').splitlines()
16                  reader = csv.reader(text_io)
17                  header = next(reader, None)
18                  if header and len(header) == 3:
19                      for row in reader:
20                          if len(row) == 3:
21                              # Prepare data for BigQuery
22                              yield {
23                                  'name': row[0],
24                                  'email': row[1],
25                                  'location': row[2],
26                                  'file_location': file_path   # Add file location here
27                              }
28                          else:
29                              logging.error(f"Invalid row in file {file_path}: {row}")
30                  else:
31                      logging.error(f"Skipping invalid file (wrong number of columns): {file_path}")
32          except Exception as e:
33              logging.error(f"Error processing file {file_path}: {e}")
```

beam-summit-2024-airflow / **sample-input-file.csv**

SadeeqAkintola  Input file format and example

Preview | Code | Blame | 3 lines (3 loc) · 108 Bytes

🔍 Search this file

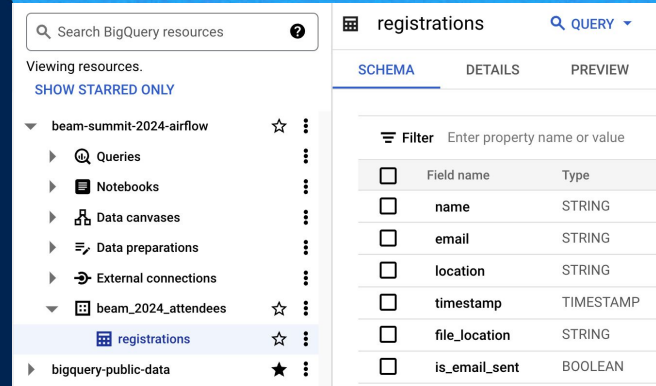| | name | email | location |
|---|---|---|---|
| 1 | | | |
| 2 | Your-name | your-email@example.com | your-city |
| 3 | Sadeeq | datatalkswithsadeeq@gmail.com | Lisbon |

25

# Building the Functions for the Pipeline

```python
class AddTimestamp(beam.DoFn):
    def process(self, element):
        element['timestamp'] = datetime.datetime.utcnow().strftime('%Y-%m-%d %H:%M:%S.%f UTC')
        return [element]

def run(argv=None):
    pipeline_options = PipelineOptions(argv)
    google_cloud_options = pipeline_options.view_as(GoogleCloudOptions)
    google_cloud_options.project = 'beam-summit-2024-airflow'
    #google_cloud_options.job_name = f"beam-summit-2024-run-{datetime.datetime.now().strftime('%Y%m%d-%H%M%S')}"
    google_cloud_options.region = 'us-central1'
    pipeline_options.view_as(StandardOptions).runner = 'DataflowRunner'
    pipeline_options.view_as(SetupOptions).save_main_session = True

    p = beam.Pipeline(options=pipeline_options)

    input_pattern = 'gs://beam-summit-2024/initiated-runs/*.csv'
    table_schema = {
        'fields': [
            {'name': 'name', 'type': 'STRING', 'mode': 'NULLABLE'},
            {'name': 'email', 'type': 'STRING', 'mode': 'NULLABLE'},
            {'name': 'location', 'type': 'STRING', 'mode': 'NULLABLE'},
            {'name': 'timestamp', 'type': 'TIMESTAMP', 'mode': 'NULLABLE'},
            {'name': 'file_location', 'type': 'STRING', 'mode': 'NULLABLE'}  # Add schema for file_location
        ]
    }
```

## Worth Noting

- **beam.DoFn :** The DoFn object that you pass to ParDo contains the processing logic that gets applied to the elements in the input collection.

- **PipelineOptions :** are used to configure Pipelines. You can extend PipelineOptions to create custom configuration options.

- **beam.Pipeline :** is the entry point for constructing and running a data processing pipeline, defining the series of transformations and operations that will be executed on the input data.

# Bringing it all together, and Run!

```python
 1    (p
 2        | 'MatchFiles' >> MatchFiles(input_pattern)
 3        | 'ReadMatches' >> ReadMatches()
 4        | 'ReadAndValidateCSV' >> beam.ParDo(ReadAndValidateCSV())
 5        | 'AddTimestamp' >> beam.ParDo(AddTimestamp())
 6        | 'WriteToBigQuery' >> WriteToBigQuery(
 7                'beam-summit-2024-airflow:beam_2024_attendees.registrations',
 8            schema=table_schema,
 9            create_disposition=BigQueryDisposition.CREATE_IF_NEEDED,
10            write_disposition=BigQueryDisposition.WRITE_APPEND
11        )
12    )
13
14    result = p.run()
15    result.wait_until_finish()
16
17 if __name__ == '__main__':
18     run()
```

Search BigQuery resources

Viewing resources.
SHOW STARRED ONLY

- beam-summit-2024-airflow
  - Queries
  - Notebooks
  - Data canvases
  - Data preparations
  - External connections
  - beam_2024_attendees
    - registrations
- bigquery-public-data

registrations          QUERY

SCHEMA      DETAILS      PREVIEW

Filter   Enter property name or value

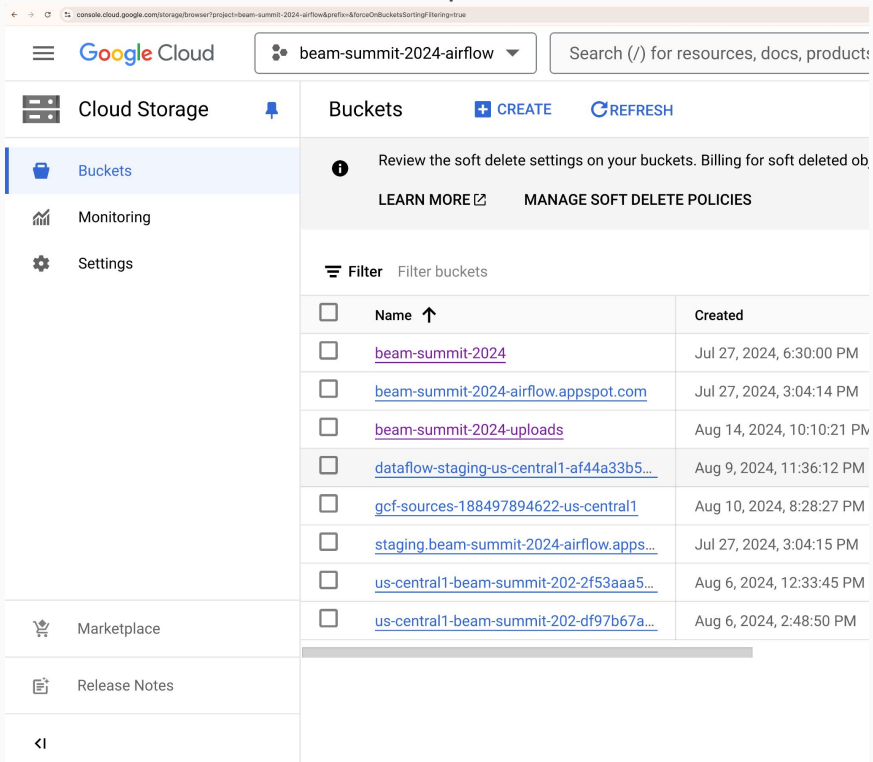| Field name | Type |
| --- | --- |
| name | STRING |
| email | STRING |
| location | STRING |
| timestamp | TIMESTAMP |
| file_location | STRING |
| is_email_sent | BOOLEAN |

## Worth Noting

- **The '|' symbol** is used as an operator to apply transformations to a PCollection

- **WriteToBigQuery :** is a transform used to write data from a PCollection to a BigQuery table.
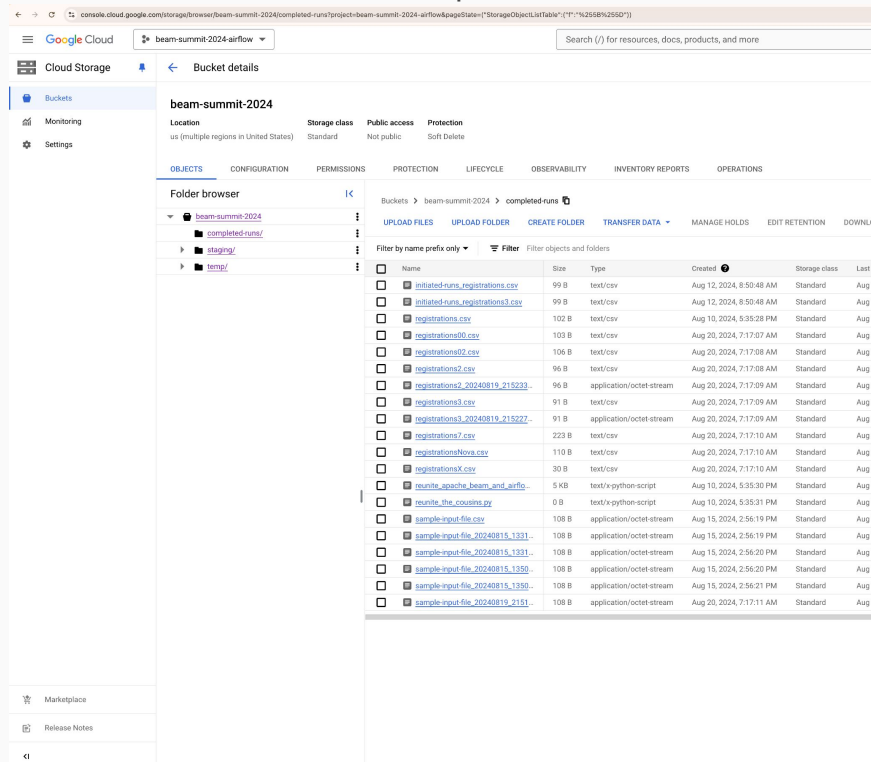
BEAM SUMMIT

# Results (Beam) - GCS Folders

Before the Pipeline Run



After the Pipeline Run

# Results (Beam) - Dataflow and BigQuery

# Code Review (Airflow DAG)

# Import Libraries, Create the DAG, and Input the Environment Variables in Airflow UI

```python
1   from airflow import DAG
2   from airflow.operators.python_operator import PythonOperator, BranchPythonOperator
3   from airflow.operators.dummy_operator import DummyOperator
4   from airflow.utils.dates import days_ago
5   from google.cloud import storage
6   from airflow.providers.google.cloud.operators.dataflow import DataflowCreatePythonJobOperator
7   from airflow.providers.google.cloud.hooks.bigquery import BigQueryHook
8   import sendgrid
9   from sendgrid.helpers.mail import Mail
10  import os
11  from datetime import datetime
12  import vertexai
13  from vertexai.generative_models import GenerativeModel, SafetySetting
14
15  # Define the DAG
16  dag = DAG(
17      'reunite_apache_beam_and_airflow',
18      default_args={
19          'owner': 'airflow',
20          'start_date': days_ago(1),
21      },
22      schedule_interval=None,
23  )
```

Airflow
DAGs    Cluster Activity    Datasets    Browse    Admin    Docs    Composer

## Edit Variable

**Key \***                                          SENDGRID_API_KEY

**Val**                                             ********

**Description**                                     Description

Save  ←

Google Cloud    beam-summit-2024-airflow ▾

Composer    ←    Environment details    ⬈ OPEN AIRFLOW UI

✅ beam-summit-2024-composer-instance    This environment is running

MONITORING    LOGS    DAGS    ENVIRONMENT CONFIGURATION    AIRFLO

✏ EDIT

Required libraries from the Python Package Index (PyPI)

| Name | Version |
| --- | --- |
| sendgrid | - |
| google-cloud-aiplatform | - |

# Operators! Operators!! Operators!!!

- **The PythonOperator**: executes Python functions as tasks in a DAG, allowing for flexible workflow management and integration with other tasks. It supports arguments, retries, and logging.

- **The BranchPythonOperator:** allows you to conditionally direct the execution flow of a DAG. It runs a Python function that returns the ***taskID of the next task to execute***, effectively branching the workflow. Only the branch selected runs, and downstream tasks are determined by this choice, allowing for ***dynamic workflows*** based on runtime conditions.

- **The DummyOperator:** is a no-op operator used primarily as a placeholder in DAGs. ***It doesn't perform any action*** but can be useful for organizing complex workflows, acting as a boundary, or grouping tasks without executing any tasks itself. It's often used for joining or splitting task flows or as a marker in a DAG's structure.

```python
1   # Function to check for new files in the uploads bucket
2   def check_for_new_files(**kwargs):
3       client = storage.Client()
4       bucket_name = 'beam-summit-2024-uploads'
5       bucket = client.get_bucket(bucket_name)
6       blobs = list(bucket.list_blobs())  # List all blobs in the bucket
7
8       if len(blobs) ≥ 5:
9           # If there are 5 or more files, proceed with the DAG
10          file_names = [blob.name for blob in blobs]
11          kwargs['ti'].xcom_push(key='file_names', value=file_names)
12          return 'move_files_to_initiated_runs'
13      else:
14          # If there are fewer than 5 files, branch to print message and stop
15          return 'insufficient_files_task'
16
17  check_files_task = BranchPythonOperator(
18      task_id='check_files_task',
19      python_callable=check_for_new_files,
20      provide_context=True,
21      dag=dag,
22  )
23
24  # Task to print insufficient files message
25  insufficient_files_task = PythonOperator(
26      task_id='insufficient_files_task',
27      python_callable=lambda: print("Insufficient number of files. DAG execution stopped."),
28      dag=dag,
29  )
30
31  # DummyOperator to end the branch when files are insufficient
32  end_task = DummyOperator(
33      task_id='end_task',
34      dag=dag,
35  )
```

# Explaining the *xcom* Libraries

The *xcom_pull* is a method in Apache Airflow that allows a task to retrieve data (XCom) pushed by a previous task. It can pull specific data based on task ID, key, and execution date. It's typically used for inter-task communication within a DAG to share information across tasks.

- **Purpose**: xcom_pull retrieves shared data (XComs) between tasks in a DAG.
- **Communication**: Allows tasks to share data using xcom_push and xcom_pull.
- **Task ID**: Pull data from a specific task by specifying its task_id.
- **Key-Value:** Retrieve specific data by providing a key; defaults to all if not specified.
- **Execution Date:** Access data from a specific execution date if needed.
- **Return Value:** Returns the relevant data or a list if multiple records match.
- **Usage:** Commonly used in task functions/operators for dynamic workflows.
- **Limitations:** Best for small data; use external storage for large data transfers.

```python
# Function to move files from the uploads bucket to initiated-runs
def move_files_to_initiated_runs(**kwargs):
    client = storage.Client()
    source_bucket_name = 'beam-summit-2024-uploads'
    destination_bucket_name = 'beam-summit-2024'
    file_names = kwargs['ti'].xcom_pull(task_ids='check_files_task', key='file_names')

    if file_names:
        source_bucket = client.get_bucket(source_bucket_name)
        destination_bucket = client.get_bucket(destination_bucket_name)
        for file_name in file_names:
            blob = source_bucket.blob(file_name)
            # Ensure the file is placed under initiated-runs/
            new_name = f'initiated-runs/{file_name.split("/")[-1]}'
            blob_copy = source_bucket.copy_blob(blob, destination_bucket, new_name)
            print(f'Copied {file_name} to {new_name}')
            # Optionally delete the file from the uploads bucket after copying
            blob.delete()
            print(f'Deleted {file_name} from source bucket')
    else:
        print('No files found to move.')

move_files_task = PythonOperator(
    task_id='move_files_to_initiated_runs',
    python_callable=move_files_to_initiated_runs,
    provide_context=True,
    dag=dag,
)
```

# This code block is the reason for this talk!

```python
1   # Task 2: Run the Beam pipeline on Dataflow
2   beam_task = DataflowCreatePythonJobOperator(
3       task_id='run_beam_pipeline',
4       py_file='gs://beam-summit-2024/beam_summit_attendee_upload.py',
5       py_options=[],
6       job_name=f'beam_pipeline_run_by_airflow_at_{datetime.now().strftime("%Y%m%d%H%M%S")}',
7       dataflow_default_options={
8           'project': 'your-gcp-project',
9           'region': 'us-central1',
10          'stagingLocation': 'gs://beam-summit-2024/staging',
11          'tempLocation': 'gs://beam-summit-2024/temp'
12      },
13      location='us-central1',
14      dag=dag,
15  )
```

# Fetching data from BigQuery in Airflow

```python
1  # Task 3: Fetch records to send emails to, including location
2  def query_bigquery():
3      query = """
4      SELECT name, email, location
5      FROM `beam-summit-2024-airflow.beam_2024_attendees.registrations`
6      WHERE (is_email_sent IS NULL OR is_email_sent = FALSE)
7      AND (email <> "")
8      """
9      hook = BigQueryHook(gcp_conn_id='google_cloud_default', use_legacy_sql=False)
10     return hook.get_pandas_df(query)
11
12 query_bigquery_task = PythonOperator(
13     task_id='query_bigquery',
14     python_callable=query_bigquery,
15     dag=dag,
16 )
```

# Creating a Gemini Flash GenAI Model  in Airflow

```python
1  # Function to describe the location using Vertex AI
2  def describe_this_location(location):
3      vertexai.init(project="beam-summit-2024-airflow", location="us-central1")
4      model = GenerativeModel("gemini-1.5-flash-001")
5      response = model.generate_content(
6          f"Tell me some fun facts about {location} in 150 words",
7          generation_config={
8              "max_output_tokens": 8192,
9              "temperature": 1,
10             "top_p": 0.95,
11         },
12         safety_settings=[
13             SafetySetting(
14                 category=SafetySetting.HarmCategory.HARM_CATEGORY_HATE_SPEECH, threshold=SafetySetting.HarmBlockThreshold.BLOCK_MEDIUM_AND_ABOVE
15             ),
16             SafetySetting(
17                 category=SafetySetting.HarmCategory.HARM_CATEGORY_DANGEROUS_CONTENT, threshold=SafetySetting.HarmBlockThreshold.BLOCK_MEDIUM_AND_ABOVE
18             ),
19             SafetySetting(
20                 category=SafetySetting.HarmCategory.HARM_CATEGORY_SEXUALLY_EXPLICIT, threshold=SafetySetting.HarmBlockThreshold.BLOCK_MEDIUM_AND_ABOVE
21             ),
22             SafetySetting(
23                 category=SafetySetting.HarmCategory.HARM_CATEGORY_HARASSMENT, threshold=SafetySetting.HarmBlockThreshold.BLOCK_MEDIUM_AND_ABOVE
24             ),
25         ],
26         stream=False
27     )
28     return response.text.strip()
```

```python
1   # Task 4: Send email to the identified users
2   def send_email(**context):
3       sg = sendgrid.SendGridAPIClient(api_key=os.environ['SENDGRID_API_KEY'])
4       results = context['task_instance'].xcom_pull(task_ids='query_bigquery')
5
6       for _, row in results.iterrows():
7           name = row['name']
8           email = row['email']
9           location = row['location']
10
11          # Generate fun facts about the location
12          location_funfact = describe_this_location(location)
13
14          # Convert the fun facts into bullet points
15          bullet_points = ''.join([f'<li>{sentence.strip()}</li>' for sentence in location_funfact.split('.') if sentence.strip()])
16
17          # Create the email content in HTML format
18          subject = 'Welcome to BEAM Summit 2024 Demo by Sadeeq Akintola'
19          content = f"""
20          <html>
21          <body>
22              <p>Dear {name},</p>
23
24              <p>Thank you for attending my BEAM Summit 2024 session
25              <a href="https://beamsummit.org/sessions/2024/reuniting-the-two-distant-cousins-calling-a-beam-pipeline-from-an-airflow-job/">
26              https://beamsummit.org/sessions/2024/reuniting-the-two-distant-cousins-calling-a-beam-pipeline-from-an-airflow-job/</a>,
27              and testing out the demo! It means a lot to me!! Here is the link to the source code:
28              <a href="https://github.com/SadeeqAkintola/beam-summit-2024-airflow">https://github.com/SadeeqAkintola/beam-summit-2024-airflow</a>.
29              Fork it, Star it, and Share it, please.</p>
30
31              <p>Please be informed that your CSV file containing <strong>{email}</strong> has been successfully uploaded.</p>
32
33              <p><em>By the way, here are some fun facts about your {location} (generated, with love, by Google's powerful gemini-1.5-flash model):</em></p>
34              <ul>{bullet_points}</ul>
35
36              <p>Sincerely yours in Data Engineering,</p>
37              <p><strong>Sadeeq</strong></p>
38              <p>Follow on X for more: <a href="https://x.com/SadeeqAkintola">https://x.com/SadeeqAkintola</a></p>
39          </body>
40          </html>
41          """
42
43          # Send the email
44          mail = Mail(from_email='datatalkswithsadeeq@gmail.com', to_emails=email, subject=subject, html_content=content)
45          response = sg.send(mail)
46
47          print(f'Sent email to: {email} | Status Code: {response.status_code}')
48
49  send_email_task = PythonOperator(
50      task_id='send_email',
51      provide_context=True,
52      python_callable=send_email,
53      dag=dag,
54  )
```

# Sending Emails in Airflow using SendGrid

- **Create an account at** https://sendgrid.com/ and register your API key in the Airflow Environment Variable.

- **Call the Generative AI Model** `describe_this_location(location)` function to generate fun facts about the location entered in the csv file earlier.

- Infuse the results returned with a preconfigured text to for the email body. The use the Mail function to send the email.

# Update the Email Flag column in BigQuery from Airflow

```python
1   # Task 5: Update the is_email_sent flag in BigQuery
2   def update_bigquery():
3       update_query = """
4       UPDATE `beam-summit-2024-airflow.beam_2024_attendees.registrations`
5       SET is_email_sent = TRUE
6       WHERE (is_email_sent IS NULL OR is_email_sent = FALSE)
7       AND (email <> "")
8       """
9       hook = BigQueryHook(gcp_conn_id='google_cloud_default', use_legacy_sql=False)
10      hook.run_query(update_query)
11
12  update_bigquery_task = PythonOperator(
13      task_id='update_bigquery',
14      python_callable=update_bigquery,
15      dag=dag,
16  )
```

# Final Steps: Setup the Task Dependencies

- **DAG Structure:** Tasks are organized in a DAG, with nodes representing tasks and edges representing dependencies.

- **Upstream/Downstream:** Tasks must complete upstream tasks before downstream tasks can start.

- **Setting Dependencies:** Use **>>, <<,** *set_upstream()*, and *set_downstream()* to define task order.

- **Trigger Rules:** Control task execution based on upstream task outcomes (e.g., all_success, one_success).

- **Cross-DAG Dependencies**: Enable tasks in one DAG to trigger tasks in another using sensors.

- **Task Groups:** Group tasks for better management and visualization of dependencies.

```python
 1  # Task 6: Move processed files to completed-runs
 2  def move_files_to_completed(**kwargs):
 3      client = storage.Client()
 4      bucket_name = 'beam-summit-2024'
 5      blobs = client.list_blobs(bucket_name, prefix='initiated-runs/')
 6      for blob in blobs:
 7          new_name = blob.name.replace('initiated-runs/', 'completed-runs/')
 8          bucket = client.get_bucket(bucket_name)
 9          new_blob = bucket.rename_blob(blob, new_name)
10          print(f'Moved {blob.name} to {new_name}')
11
12  move_files_to_completed_task = PythonOperator(
13      task_id='move_files_to_completed',
14      python_callable=move_files_to_completed,
15      dag=dag,
16  )
17
18  # Setting up the task dependencies for the entire DAG
19  check_files_task >> insufficient_files_task >> end_task
20  check_files_task >> move_files_task >> beam_task >> query_bigquery_task >>
    send_email_task >> update_bigquery_task >> move_files_to_completed_task
```

# Use Cloud Functions to Trigger the Airflow DAG once there's a new *.csv file in GCS

```python
1  from __future__ import annotations
2
3  from typing import Any
4
5  import google.auth
6  from google.auth.transport.requests import AuthorizedSession
7  import requests
8
9
10 # Following GCP best practices, these credentials should be
11 # constructed at start-up time and used throughout
12 # https://cloud.google.com/apis/docs/client-libraries-best-practices
13 AUTH_SCOPE = "https://www.googleapis.com/auth/cloud-platform"
14 CREDENTIALS, _ = google.auth.default(scopes=[AUTH_SCOPE])
15
16
17 def make_composer2_web_server_request(
18     url: str, method: str = "GET", **kwargs: Any
19 ) -> google.auth.transport.Response:
20     """
21     Make a request to Cloud Composer 2 environment's web server.
22     Args:
23       url: The URL to fetch.
24       method: The request method to use ('GET', 'OPTIONS', 'HEAD', 'POST', 'PUT',
25         'PATCH', 'DELETE')
26       **kwargs: Any of the parameters defined for the request function:
27                 https://github.com/requests/requests/blob/master/requests/api.py
28                 If no timeout is provided, it is set to 90 by default.
29     """
30
31     authed_session = AuthorizedSession(CREDENTIALS)
32
33     # Set the default timeout, if missing
34     if "timeout" not in kwargs:
35         kwargs["timeout"] = 90
36
37     return authed_session.request(method, url, **kwargs)
```

```python
1  def trigger_dag(web_server_url: str, dag_id: str, data: dict)
   → str:
2      """
3      Make a request to trigger a dag using the stable Airflow 2
   REST API.
4      https://airflow.apache.org/docs/apache-airflow/stable/stable-rest-api-ref.html
5
6      Args:
7        web_server_url: The URL of the Airflow 2 web server.
8        dag_id: The DAG ID.
9        data: Additional configuration parameters for the DAG run (json).
10     """
11
12     endpoint = f"api/v1/dags/{dag_id}/dagRuns"
13     request_url = f"{web_server_url}/{endpoint}"
14     json_data = {"conf": data}
15
16     response = make_composer2_web_server_request(
17         request_url, method="POST", json=json_data
18     )
19
20     if response.status_code == 403:
21         raise requests.HTTPError(
22             "You do not have a permission to perform this operation. "
23             "Check Airflow RBAC roles for your account."
24             f"{response.headers} / {response.text}"
25         )
26     elif response.status_code ≠ 200:
27         response.raise_for_status()
28     else:
29         return response.tex
```

```python
1  """
2  Trigger a DAG in a Cloud Composer 2 environment in response to an event,
3  using Cloud Functions.
4  """
5
6  from typing import Any
7
8  import composer2_airflow_rest_api
9
10 def trigger_dag_gcf(data, context=None):
11     """
12     Trigger a DAG and pass event data.
13
14     Args:
15       data: A dictionary containing the data for the event. Its format depends
16         on the event.
17       context: The context object for the event.
18
19     For more information about the arguments, see:
20     https://cloud.google.com/functions/docs/writing/background#function_parameters
21     """
22
23     # TODO(developer): replace with your values
24     # Replace web_server_url with the Airflow web server address. To obtain this
25     # URL, run the following command for your environment:
26     # gcloud composer environments describe example-environment \
27     #   --location=your-composer-region \
28     #   --format="value(config.airflowUri)"
29     web_server_url = (
30         "https://26d527f1cb8b4cbdb8c3bd839a44e35a-dot-us-central1.composer.googleusercontent.com"
31     )
32     # Replace with the ID of the DAG that you want to run.
33     dag_id = 'reunite_apache_beam_and_airflow'
34
35     composer2_airflow_rest_api.trigger_dag(web_server_url, dag_id, data)
```

40

# Results (Airflow and Email)

# Remember: "Airflow is (just) an Orchestrator"



PREMIUM
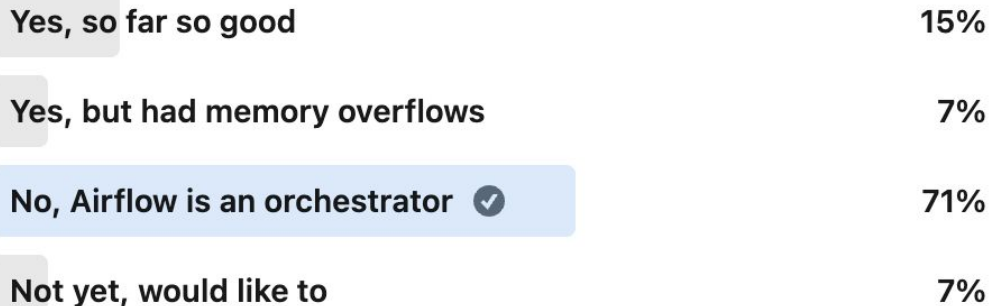AIRFL... ...E MOON 🚀

**Marc Lamberti** 🔗
Head of Customer
Education @Astronomer |
Best Selling Instructor
@Udemy

Followers          63,251

**Message**

## Do you process data in Airflow?
The author can see how you vote. **Learn more**

| | |
|---|---|
| Yes, so far so good | **15%** |
| Yes, but had memory overflows | **7%** |
| No, Airflow is an orchestrator ✓ | **71%** |
| Not yet, would like to | **7%** |

772 votes • Poll closed • **Remove vote**

👍😄 15                                                          11 comments

# Demo!

# Have you tried the Demo? Starred/Cloned the Repo? Oya, Do it now!!

# Useful Resources

**Official Documentation:**

1. **Apache Beam Documentation:** https://beam.apache.org/documentation/
2. **Apache Airflow Documentation:** https://airflow.apache.org/docs/
3. **Google Dataflow Documentation:** https://cloud.google.com/dataflow/docs/
4. **Triggering Beam Pipelines with Cloud Composer (Google Documentation):** https://cloud.google.com/composer/docs/how-to/using/triggering-with-gcf

**Popular Medium Posts:**

5. **Event-Based Dataflow Job Orchestration with Cloud Composer, Airflow, and Cloud Functions:**
   https://gulia.medium.com/event-based-dataflow-job-orchestration-with-cloud-composer-airflow-and-cloud-functions-b61219f9aeaf
6. **Launching Dataflow Pipelines via Cloud Composer (Airflow):**
   https://medium.com/@kolban1/cloud-composer-launching-dataflow-pipelines-38cd29e970d4
7. **Launch an Apache Beam Pipeline with Apache Airflow — Part 1/2: Setting up the Airflow Environment with Docker-Compose:**
   https://medium.com/@carmelwenga/launch-an-apache-beam-pipeline-with-apache-airflow-part-1-setting-up-the-airflow-environment-d97dd64ded18

**YouTube Videos:**

8. **Apache Beam: A Unified Model for Batch and Streaming Data Processing:** https://www.youtube.com/watch?v=7DZ8ONmeP5A
9. **Flexible, Easy Data Pipelines on Google Cloud with Cloud Composer (Cloud Next '18):** https://www.youtube.com/watch?v=GeNFEtt-D4k
10. **Cloud Composer - Orchestrating an ETL Pipeline Using Cloud Dataflow:** https://www.youtube.com/watch?v=PCg9AQNGX3E

**Also, join us for Airflow Summit Next Week:** https://airflowsummit.org