

CS 329E Project 6, due Thursday, 03/07.

In this project, we continue our automation effort and focus on orchestrating the pipelines in the staging layer.

### Objectives

- Create Cloud Composer environment, a managed Airflow service on Google Cloud
- Develop an Airflow pipeline that creates and populates the tables in staging
- Execute the Airflow pipeline
- Delete and re-create your Cloud Composer environment to reduce billing charges

### Implementation Guidelines

Please follow these guidelines when developing your Airflow DAGs:

- Store the staging tables in a new dataset in BigQuery. The name of the staging dataset should follow the convention of **[domain]\_stg\_af**.
- Use the provided code samples **p6-model-controller.py**, **p6-create-bird-airports.py**, and **p6-create-meal-snack-service.py** as a starting point for your own DAGs.
- Ensure that the resulting tables generated through Airflow match the staging tables created from your notebook. Both the record count and contents should match although in this iteration, you may use a more precise schema (with not null constraints and non-string types).
- Take a [screenshot](#) of your model controller execution showing that all tasks completed successfully and name it **model-controller-run.png**.
- When not actively using Composer, delete the environment to avoid burning through all of your GCP credits. Note: there is no way to stop a running Composer instance. Follow this guide to spin up a new instance.
- Publish to your repo: **model-controller.ipynb**, **model-controller-run.png**, and any additional sub-DAGs triggered from your main DAG.

CS 329E Project 6 Rubric

**Due Date: 03/07/24**

<p><code>model-controller.py</code> has all required info and correctly populates tables</p> <ul style="list-style-type: none"> <li>-10 did not update global variables</li> <li>-10 for each table missing</li> <li>-10 for each BigQueryInsertJobOperator object missing</li> <li>-15 if upload .ipynb instead of .py</li> <li>-40 missing file</li> </ul>	40
<p>At least one file for separate sub-DAGs is present and has all required information</p> <ul style="list-style-type: none"> <li>- If multiple files, divide points by amount of DAG files</li> </ul> <ul style="list-style-type: none"> <li>-10 does not have <code>serialize_datetime</code> or <code>remove_none_values</code> methods</li> <li>-20 does not load records properly</li> <li>-20 missing record insertion method (will use <code>@task</code> decorator)</li> <li>-40 missing file</li> </ul>	35
<p><code>model-controller-run.png</code> shows proper proof of Airflow controller execution</p> <ul style="list-style-type: none"> <li>-10 missing file</li> </ul>	10
<p>Google Cloud BigQuery bucket has properly loaded all tables</p> <ul style="list-style-type: none"> <li>-5 for each missing table</li> <li>-10 if tables not under "stg_af" dataset</li> <li>-20 missing file</li> </ul>	15
<p><code>submission.json</code> submitted into Canvas. Your project <b>will not</b> be graded without this submission. The file should have the following schema:</p> <pre>{   "commit-id": "your most recent commit ID from Github",   "project-id": "your project ID from GCP" }</pre> <p>Example:</p> <pre>{   "commit-id": "dab96492ac7d906368ac9c7a17cb0dbd670923d9",   "project-id": "some-project-id" }</pre>	Required
<p><b>Total Credit:</b></p>	<b>100</b>