

CS 329E Project 7, due Thursday, 03/21.

In this project, we complete the automation sprint by implementing referential integrity checks, creating the final target tables, and orchestrating the end-to-end pipeline.

### Objectives

- Create Cloud Composer environment
- Develop an Airflow pipeline that creates the primary keys and foreign keys on the staging tables
- Develop an Airflow pipeline that validates the primary key and foreign key constraints on the staging tables
- Develop an Airflow pipeline that creates the final target tables in the consumption layer
- Develop an Airflow controller that calls each one of the sub controllers to automate the entire pipeline from ingest to consumption.
- Successfully execute the end-to-end pipeline. Reduce the level of task parallelism if necessary.
- Delete and re-create your Cloud Composer environment to reduce billing charges

### Implementation Guidelines

Please follow these guidelines when implementing your solution:

- Create a custom Composer cluster with **3.25G** of RAM for the worker nodes (instead of the default 2G).
- Store the target tables in a new dataset in BigQuery. The name of the consumption dataset should follow our convention of **[domain]\_csp\_af**.
- Ensure that the target tables generated through Airflow match the ones created from Colab in the consumption layer. Both the schema and contents should match.
- Use the 5 provided code samples as a starting point for your own DAGs: **p7-key-controller.py**, **p7-create-pk.py**, **p7-create-fk.py**, **p7-target-controller.py**, and **p7-master-controller.py**.
- If you encounter non-determinist failures when running your DAG, it is likely a resource issue. You can reduce the number of concurrent tasks using the `DAG.currency` parameter as shown [here](#) (line 29).
- Take a screenshot of your master controller execution run showing that all tasks completed successfully. Name the file **master-controller-run.png**.
- When not actively developing, delete your Composer instance to avoid burning through your GCP credits. Note: there is no way to stop and restart a Composer instance.
- Publish to your repo: **key-controller.py**, **create-pk.py**, **create-fk.py**, **target-controller.py**, **master-controller.py**, and **master-controller-run.png**.

CS 329E Project 7 Rubric

**Due Date: 03/21/24**

<p>key-controller.py has all required info and correctly populates keys</p> <ul style="list-style-type: none"> <li>-10 did not update global variables</li> <li>-10 for each TriggerDagRunOperator object missing</li> <li>-15 if upload .ipynb instead of .py</li> <li>-40 missing file</li> </ul>	40
<p>target-controller.py creates dataset and populates all tables correctly</p> <ul style="list-style-type: none"> <li>-5 did not update global variables</li> <li>-5 for each table missing</li> <li>-5 for each BigQueryInsertJobOperator object missing</li> <li>-10 if upload .ipynb instead of .py</li> <li>-25 missing file</li> </ul>	25
<p>create-pk.py, create-fk.py, master-controller.py have all required info</p> <ul style="list-style-type: none"> <li>-5 missing _create_fk and _check_ref_integrity methods in create-fk.py</li> <li>-5 missing _create_pk and _check_pk methods in ,create-pk.py</li> <li>-5 missing TriggerDagRunOperator objects in master-controller.py for all controllers</li> <li>-5 for each missing file</li> </ul>	10
<p>Google Cloud BigQuery bucket has properly loaded all keys</p> <ul style="list-style-type: none"> <li>-5 for each missing key or table</li> <li>-10 if tables not under "csp_af" dataset</li> <li>-20 missing file</li> </ul>	15
<p>master-controller-run.png shows proper proof of Airflow controller execution</p> <ul style="list-style-type: none"> <li>-10 missing file</li> </ul>	10
<p>submission.json 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

<b>Total Credit:</b>	<b>100</b>
----------------------	------------