

CS 327E Milestone 4 due Thursday, 12/10 at **12pm**.

## Part 1:

Convert your previously written Beam pipeline(s) to Dataflow. Run them on Dataflow over the entire input data; debug and fix as necessary.

### General Coding Conventions:

- Create a new notebook `milestone4.ipynb` and call the pipelines from the notebook.
- The code should be commented sufficiently to follow the main logic of the transforms.

### Dataflow Coding Conventions:

- A Beam pipeline should transform a single source table.
- All transforms applied to a source table should be placed in the same Beam pipeline.
- A pipeline script should be named `<table>_dataflow.py`.
- A table should be named `<table>_Dataflow` when produced by the Dataflow Runner.

## Part 2:

Verify that each BigQuery output table (e.g. `<table>_Dataflow`) contains a valid primary key. Child tables must also have a valid foreign key. Run the appropriate SQL statements within your `milestone4` notebook to verify these constraints.

Update your ERD to reflect the schema of your transformed tables:

- Diagram should capture only the latest version of each table (e.g. `<table>_Dataflow`).
- Entity types should specify field names, data types, and keys for each table.
- Diagram should visually indicate the source of each entity (e.g. entities from `dataset1` can use one background color while entities from `dataset2` can use a different background color).
- Draw the relationships between the entities within `dataset2` as well as **across** the two datasets.
- Name your ERD file `final_unified_model.pdf`.

## Part 3:

1. Implement your cross-dataset queries:

- Develop and run three cross-dataset queries from your `milestone4` notebook
- Queries should use the modeled tables from both refined datasets

- Each query should be wrapped into a view, created in your `insights` dataset
- A short comment should appear above each SQL statement to describe its function

2. Create visualizations in Data Studio:

- Create a data visualization from each cross-dataset query
- Data Sources query the SQL views from the previous section.
- Charts should visualize the data in a compelling way.
- Add the 3 charts to your existing Data Studio report (aka dashboard).
- Download the report and save it as `final_dashboard.pdf`.

<p><b>Part 1</b> - Convert your Beam pipelines to Dataflow. Each Beam pipeline should have two Python scripts, <code>&lt;table&gt;_beam.py</code> and <code>&lt;table&gt;_dataflow.py</code> per source table.</p> <ul style="list-style-type: none"> <li>-<b>X</b> for each missing <code>&lt;table&gt;_dataflow.py</code> where X is dependent on the number of Beam pipelines. If you have 2, <b>-20</b> each. 3, <b>-13.3</b> each, and so on.             <ul style="list-style-type: none"> <li>-10 Beam pipelines not using DataflowRunner</li> <li>-10 Beam pipelines do not execute properly</li> <li>-10 Beam pipelines not writing to output table</li> </ul> </li> <li><code>&lt;table&gt;_Dataflow</code></li> <li>-10 Beam pipeline run calls missing from <code>milestone4.ipynb</code></li> </ul> <p><i>(points will be broken based on number of pipelines)</i></p>	40
<p><b>Part 2</b> - Verify primary key constraints on tables transformed by Beam. Verify foreign key constraints if those tables are also child tables. Add this logic to your notebook.</p> <ul style="list-style-type: none"> <li>-10 missing or incorrect primary key verification on final output tables</li> <li>-10 missing or incorrect foreign key verification on final child output tables</li> </ul> <p>Create an updated ERD that finalizes your table schema after Beam transforms have been applied.</p> <ul style="list-style-type: none"> <li>-10 <code>./final_unified_model.pdf</code> not found in repository</li> <li>-5 ERD is missing one or more entities</li> <li>-5 ERD is missing one or more primary keys</li> <li>-5 ERD is missing one or more foreign keys</li> <li>-5 ERD is missing or incorrect relationship between entities</li> </ul>	20
<p><b>Part 3</b> - Implement and run your three cross-dataset queries. Comment each query with the function it performs.</p> <ul style="list-style-type: none"> <li>-5 each missing or erroneous query, up to <b>-15</b></li> <li>-5 each missing or incorrect comment, up to <b>-15</b></li> <li>-5 each query not on a transformed table, up to <b>-15</b></li> </ul> <p>Create 3 data visualizations and add them to your existing Data Studio report. The visualizations should represent the results from the three BQ views.</p> <p>The Data Studio report should contain a total of <b>5 charts</b>, 2 from Milestone 2 and 3 from the current milestone. Each chart should have a relevant title describing the dataset.</p> <ul style="list-style-type: none"> <li>-20 <code>./final_dashboard.pdf</code> not found in repository</li> <li>-10 each missing chart, up to <b>-20</b></li> <li>-10 each chart created from a BQ table instead of a BQ view, up to <b>-20</b></li> <li>-5 each missing title, up to <b>-15</b></li> </ul>	40
<p><code>submission.json</code> submitted into Canvas. Your project <b>will not</b> be graded without this</p>	<b>Required</b>

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