

CS 327E Final Project: Milestone 2, due Wednesday 12/06. **Due date is not flexible.**

This milestone has three sections. The first asks you to set up your environment for Trino. The second asks you to implement three federated queries that span multiple databases. The third asks you to evaluate Trino and summarize its trade-offs.

## Part 1: Deployment

- Follow our [setup guide](#) to deploy Trino and MongoDB. Note: We are deploying MongoDB instead of using Atlas because Trino is currently incompatible with Atlas.
- Go back to `final-project-mongodb.ipynb` from Milestone 1. Change the `HOST` and `url` to point to your new MongoDB server. The `HOST` should be set to the Private IP of your Trino VM. The `url` should look like this:  

```
url = "mongodb://{host}:{port}".format(HOST, PORT)
```

Do **not** change the port number. It needs to be 27017. Refer to [this notebook](#) for more details.
- Re-run the contents of your notebook to create the `final_project` collection and populate it with the ticketing data from Milestone 1.
- If you still have your college datasets in MySQL, Postgres, and BigQuery, go through [this notebook](#) to verify that you can federate to each of your databases. This is an optional step.

## Part 2: Federated Queries

Create a Jupyter notebook and name it `final-project-trino.ipynb`. Using the Trino cli, implement the following three federated queries in your notebook. Refer to [this notebook](#) for code samples.

*Q1: Who are the shoppers who have reservations in a different city from their own city and who paid for their reservations in US dollars? For those shoppers, return their `cust_id`, email, along with their reservation count. Order the results by `cust_id`. Limit the results to 5 records.*

*Q2: Who are the shoppers who have bought non-stop tickets on American Airlines with a departure date between 01/01/2020 and 12/31/2024? Return the shopper's `cust_id`, email, `dep_airport` and `dep_date`. Order the results by `dep_date`. Limit the results to 5 records.*

*Q3: Which shoppers have reservations and tickets with matching dates such that the `arr_date` and `dep_date` are equal between reservations and ticketing? For such shoppers, sum up their `pmt_amt` and `tik_amt` and return this sum as `trip_amount`. Also, return the shoppers' `cust_id`, `email`, `pmt_amt`, `tik_amt`. Order the results by `trip_amount` in descending order. Limit the results to 5 records.*

### **Part 3: Trino Evaluation**

Now that you have seen how to implement cross-database joins, it's time to conduct a more thorough evaluation of Trino. Think of other tests you can run to gain a more complete picture of Trino's functionality. Come up with some tests and run them from your notebook. You want to test a variety of scenarios that include DDL, DML, and SELECT statements. Then, write a short paragraph to summarize your findings.

Your summary needs to be specific and supported by your test results. For example: "CREATE TABLE statement X against MySQL resulted in Y" or "BEGIN TRANSACTION and COMMIT" against Postgres resulted in Z". Write your summary as a Markdown comment on the last cell of the notebook.

CS 327E Final Project Milestone 2 Rubric

**Due Date: 12/06/23**

<p>Part 1: Deploy Trino and MongoDB on a Compute Engine VM as per <a href="#">our guide</a></p> <ul style="list-style-type: none"> <li>-10 Trino service does not start on VM</li> <li>-10 mongod is not active on VM</li> </ul>	20
<p>Part 1: Update the HOST and url variables in <code>final-project-mongodb.ipynb</code> and re-run the cells to populate the ticketing collection with 100 documents.</p> <ul style="list-style-type: none"> <li>-10 MongoDB server on Trino VM does not contain a ticketing collection</li> <li>-8 Ticketing collection is empty or is missing documents</li> </ul>	10
<p>Part 2: Implement federated queries Q1, Q2, and Q3.</p> <ul style="list-style-type: none"> <li>-3 for each missing or incorrect join clause</li> <li>-2 for each missing or incorrect where clause</li> <li>-2 for each missing or incorrect select clause</li> <li>-2 for each missing or incorrect order by clause</li> <li>-1 for each missing or incorrect limit clause</li> </ul>	30
<p>Part 3: Test a wide variety of scenarios with query federation to understand trade-offs.</p> <ul style="list-style-type: none"> <li>-10 for ignoring or missing at least one major scenario</li> <li>-10 for ignoring or missing at least one database system</li> <li>-10 for lack of logical progression or organization</li> <li>-10 for summary that is vague or unsupported</li> </ul>	40
<p><code>final-project-trino.ipynb</code> and <code>final-project-mongodb.ipynb</code> pushed to your group's private repo on GitHub. Your milestone <b>will not</b> be graded without this submission.</p>	<b>Required</b>
<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>	<b>Required</b>
<p><b>Total Credit:</b></p>	<b>100</b>