

## CS 327E Final Project: Milestone 1, due Thursday 12/02

1. Choose one of these two database systems for your final project and follow the setup guide to bring up your environment on GCP:

[MySQL Server 8](#)

[Postgres Server 13](#)

2. Download the load testing dataset from Google Cloud Storage:

```
gsutil cp gs://cs327e-open-access/load_testing/us-500.zip .  
gsutil cp gs://cs327e-open-access/load_testing/us-1000000.zip .
```

3. Unzip both files. The file `us-500.csv` contains 500 rows and is used for sampling the data and testing your code. The `us-1000000.csv` contains 1,000,000 rows of data and is what we want to load into the database. The first line in both files contains the column names. The column names are:

- First Name
- Last Name
- Company
- Address
- City
- County (where applicable)
- State/Province (where applicable)
- ZIP/Postal Code
- Phone 1
- Phone 2
- Email
- Web

4. In a new Jupyter notebook named `milestone1.ipynb`, write some python code to sample the character lengths of the fields in the `us-500.zip` file. Your code should scan the contents of the file and keep track of the max length of each column.

5. Based on the results from the previous step, write the DDL to create the table.

If you are using MySQL, include logic to create the table in a database and drop the database if it already exists before creating it. Hint: review Project 1 if you don't remember how to do this step.

If you are using Postgres, include the logic to create the table in a schema and drop the schema and its constraints if it already exists before creating it. Hint: review Project 2 if you don't remember how to do this step.

The table should use VARCHAR types for variable-length fields and CHAR types for regular length fields (e.g. zip code, phone number, etc.). For each VARCHAR field in the table, add 5 characters to the length you sample in case there are longer values in the `us-1000000.csv` file.

Define a Primary Key for the table using an auto-incremented field. Refer to the documentation for instructions on how to create the auto-incremented field ([MySQL doc](#)) ([Postgres doc](#)).

Follow these naming conventions when writing the DDL script:

Database (mysql)	load_testing
Schema (postgres)	load_testing
Table	Person
Columns	first_name, last_name, etc.
Primary Key	id
DDL File	create_table.sql

6. In your notebook, install the appropriate Python connector for your chosen system.

```
pip install mysql-connector-python (documentation)  
pip install psycopg[binary] (documentation)
```

7. In your notebook, write some Python code to load the contents of `us-1000000.csv` into the Person table.

Skip the first line of the file as it contains only the column headings.

If you are using MySQL, use the connector's [execute\(\)](#) method to insert each row. See [code sample](#) for more details. Do not use the connector's [executemany\(\)](#) method as this will be done as part of Milestone 2.

If you are using Postgres, use the connector's [execute\(\)](#) method to insert each row. See

[code sample](#) for more details. Do not use the connector's `execute_batch()` method as this will be done as part of Milestone 2.

Notice the `try`, `except`, `finally` blocks, make sure to include those blocks in your solution.

Issue a commit every 5000 rows and print the number of records successfully inserted into your table.

Use the `%%timeit` magic function to measure the total execution time of your code. This function executes the code block a total of 8 times and reports the average time across 7 runs.

Your code should produce the following output:

```
5000 records inserted successfully into Person table
10000 records inserted successfully into Person table
15000 records inserted successfully into Person table
20000 records inserted successfully into Person table
25000 records inserted successfully into Person table
30000 records inserted successfully into Person table
35000 records inserted successfully into Person table
...
1000000 records inserted successfully into Person table
MySQL/Postgres connection is closed
4min 45s ± 19.6 s per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

Note that your `%%timeit` results may differ from those reported here.

8. Query the Person table to get the record count. You should end up with exactly 8,000,000 records in the table after loading the `us-1000000.csv` file under the `%%timeit` function. If you have only 7,999,992 rows, you may be missing the last line of the file.

CS 327E Milestone 1 Rubric

**Due Date: 12/02/21**

Found functional database environment on Compute Engine instance, as per Final Project setup guide (MySQL 8 or Postgres 13)	10
Installed python database connector for MySQL 8 ( <code>mysql-connector-python</code> ) or Postgres 13 ( <code>psycopg2</code> ) on Jupyter instance <ul style="list-style-type: none"> <li>-5 no connector found on Jupyter instance</li> <li>-5 incorrect connector found on Jupyter instance</li> </ul>	4
Code block that scans sample file ( <code>us-500.csv</code> ) and outputs max length of each field <ul style="list-style-type: none"> <li>-5 used wrong file for sampling</li> <li>-5 incorrect calculation of max field length</li> <li>-incorrect output produced by code block: <ul style="list-style-type: none"> <li>-3 one or more fields missing from output</li> <li>-3 incorrect length for one or more fields</li> </ul> </li> </ul>	16
Write a DDL script that creates the database/schema and table with the Primary Key. <ul style="list-style-type: none"> <li>-2 for each missing column or incorrect column name</li> <li>-2 for each incorrect data type used (e.g. <code>varchar</code> in place of <code>char</code>, etc.)</li> <li>-2 for each incorrect field length used (e.g. <code>varchar(12)</code> in place of <code>varchar(20)</code>)</li> <li>-2 incorrect table name</li> <li>-3 incorrect primary key specification (sequence generator missing or incorrect)</li> <li>-2 missing or incorrect database/schema creation and drop</li> </ul>	20
Code block that loads full file ( <code>us-1000000.csv</code> ) into table and measures average load time <ul style="list-style-type: none"> <li>-10 code has syntax errors or doesn't run</li> <li>-10 incorrect CSV file parsing logic used</li> <li>-10 incorrect insert logic (e.g. not using <code>execute()</code> method)</li> <li>-10 incorrect or incomplete output produced (e.g. missing <code>%%timeit</code> function)</li> <li>-10 incorrect commit interval</li> <li>-5 missing or incorrect <code>try/except/finally</code> blocks</li> </ul>	50
<code>milestone1.ipynb</code> and <code>create_table.sql</code> pushed to your group's private repo on GitHub. Your project <b>will not</b> be graded without this submission.	<b>Required</b>
<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:  <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",</pre>	<b>Required</b>

<pre>"project-id": "some-project-id" }</pre>	
<b>Total Credit:</b>	<b>100</b>