# CSCI 5408: Assignment 2

## Problem 2:

## Task 1:

1. **olist_customers_dataset.csv**
   a. There are no Null, Blank or NaN values in this file.
   b. So, considering the top 500 rows for further operations.

2. **olist_geolocation_dataset.csv**
   a. There are no Null, Blank or NaN values in this file.
   b. So, considering the top 500 rows for further operations.

3. **olist_order_items_dataset.csv**
   a. There are no Null, Blank or NaN values in this file.
   b. So, considering the top 500 rows for further operations.

4. **olist_order_payments_dataset.csv**
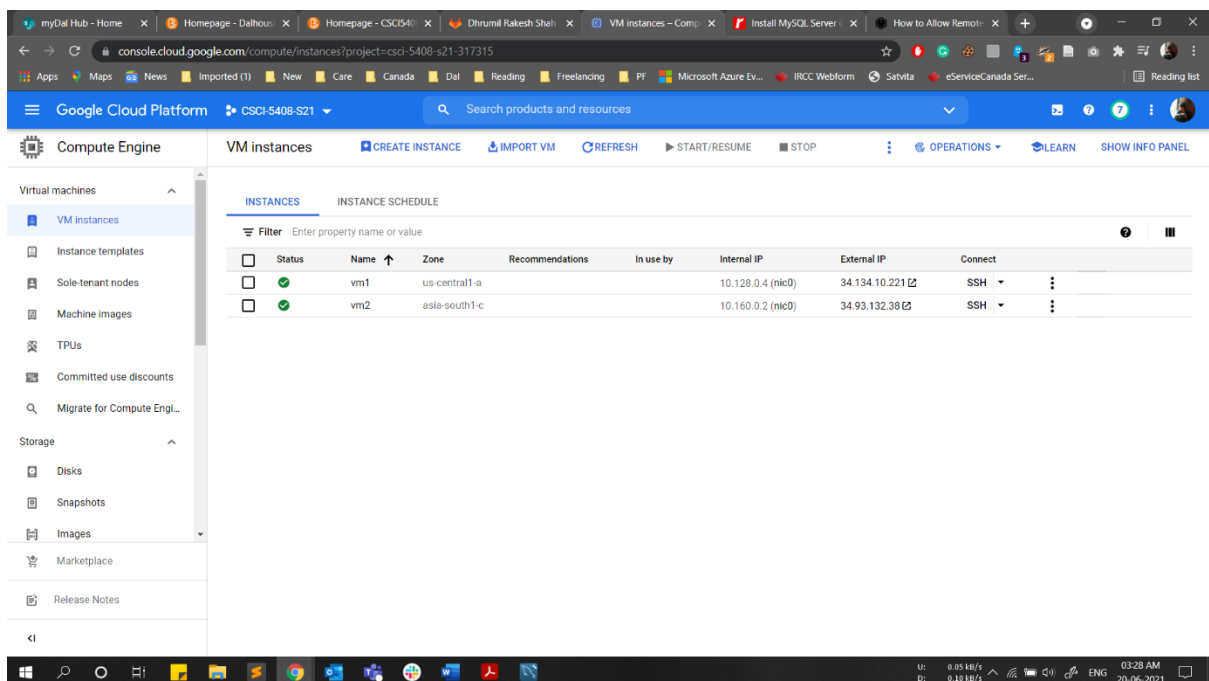   a. There are no Null, Blank or NaN values in this file.
   b. So, considering the top 500 rows for further operations.

5. **olist_order_reviews_dataset.csv**
   a. Considering the top 500 rows for further operations.
   b. Replacing all the blank values with NA under column **review_comment_title** and **review_comment_message**
   c.

6. **olist_orders_dataset.csv**
   a. Replacing all the blank values with NaN under column **order_delivered_carrier_date** and **order_delivered_customer_date**
   b. Considering the top 500 rows for further operations.

7. **olist_products_dataset.csv**
   a. Considering the top 500 rows for further operations.

     b. Rectified the name of columns **product_name_lenght** and **product_description_lenght** to **product_description_length** and **product_description_length**

     c. Replacing blank values in **product_description_length**, **product_name_length** and **product_photos_qty** with 0.

     d. Replacing blank values in **product_category_name** with NA.

**8. olist_sellers_dataset.csv**

     a. There are no Null, Blank or NaN values in this file.

     b. So, considering the top 500 rows for further operations.

**9. product_category_name_translation.csv**

     a. There are no Null, Blank or NaN values in this file.

**If the datasets are converted to database tables, and database(s), how will it be placed, state the reasons? (E.g. why did you consider specific Fragmentation, transparency etc.)**

- Firstly, I have created two MySQL instances on vm1 and vm2 virtual machines on Google Cloud Platform.
- Vm1 is located in USA whereas Vm2 is located in India. Achieving the essence of Distributed Database System.
- The dataset is downloaded from the **Brazilian E-Commerce Public Dataset by Olist website** on Kaggle.
- There are 9 csv files containing data about **customers, geolocation, order_items, order_payments, order_reviews, orders, products, sellers, product_category_name**.
- The fragmentation is done based on multiple factors:
  - From the database tables point of view keeping the Customers data and Geolocation data together at the same location gives more meaning.
    - As there are very high chances of retrieving the customers and geolocation data simultaneously.
  - On the other hand, keeping orders, order items, order payments, order reviews, products, product categories and the respective seller's data together at the same location.
    - Keeping the order related data with the products and their respective sellers gives more insight upon being fetched.
- Transparency in distributed database is very important, as it gives a clear logical picture about the stored data.

**Global Data Dictionary (GDD):**

- GDD is an essential element for implementing a fully distributed DB environment.
- This file maps the client's logical DB names to the physical DB names on different servers being used.
- The GDD can be considered to be an address book used by the client.
- Using GDD the client will be able to reference the DB via logical names instead of the physical DB names.
- An GDD.xml file is created storing the details about the two databases and their respective tables stored at different locations.
- The GDD file is stored in both the VM instances.

**VM Instances:**



**VM1 Instance:**

**VM2 Instance:**

## Task 2:

### Before Locking:

The first part where before applying the exclusive locks on transactions T1, T2 and T3, the sequence of execution of transactions is random and we don't have any control over it. This randomness causes the table to be in an inconsistent state because of the switching between

the transactions. This causes the problems like dirty read and no locking is being done to prevent such issues. Every time the transactions are run, different order is followed.

```
T1 started
T3 started
T2 started
SELECT query executed of T1
SELECT query executed of T2
SELECT query executed of T3
Read 1 customers data successfully for T3.
Read 1 customers data successfully for T2.
Read 1 customers data successfully for T1.
Updated city of 1 customers data successfully for T2.
T2 committed
Updated city of 1 customers data successfully for T1.
T1 committed
Updated city of 1 customers data successfully for T3.
T3 committed
```

```
T2 started
T1 started
T3 started
SELECT query executed of T1
SELECT query executed of T2
SELECT query executed of T3
Read 1 customers data successfully for T2.
Read 1 customers data successfully for T3.
Read 1 customers data successfully for T1.
Updated city of 1 customers data successfully for T2.
T2 committed
Updated city of 1 customers data successfully for T3.
T3 committed
Updated city of 1 customers data successfully for T1.
T1 committed
```

The point being that, exclusive locks on transactions is very essential to order the transaction execution and avoid problems like dirty read and write causing inconsistency in the database.

**After Locking:**

After using the exclusive locks on the transactions being executed in the given sequence, we can clearly see the difference in the execution sequence, commit sequence and complete sequence. The use of exclusive lock on transactions basically grants the transactions the permission to perform their operations while preventing other transactions to gain access to the locked resource, which here in this case in the customers table.

The exclusive lock helps avoid dirty read and write problems as it locks the tables while some transaction is performing operations on that table. In the code, exclusive read and write locks have been used for the sequence of executions mentioned in the pdf. So, before using exclusive locks, any of the three transactions would commit, but after using the exclusive locks, the commits are also properly sequenced with either T3 or T1 committing before T2.

```
T2 started
T1 started
T3 started
SELECT query executed of T3
Read 1 customers data successfully for T3.
Updated city of 1 customers data successfully for T3.
T3 committed
SELECT query executed of T1
Read 1 customers data successfully for T1.
Updated city of 1 customers data successfully for T1.
T1 committed
SELECT query executed of T2
Read 1 customers data successfully for T2.
Updated city of 1 customers data successfully for T2.
T2 committed
```