**Executive Summary**

The whole objective of this assignment was to implement a data-ware house using oracle cloud and Apachehop to connect my local database to Oracle Cloud

This proceeded in the following steps:
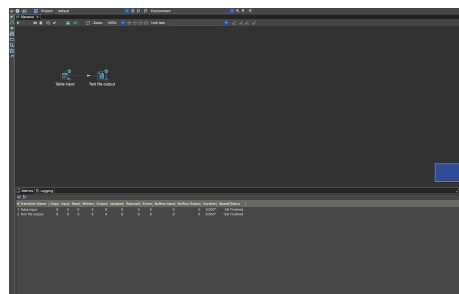
**Database Creation on Oracle and Apache Hop Setup:**
First I created a database in Oracle Cloud with Product_Dimension, Customer_Dimension, Date_Dimension, and FactSales tables. These tables were created using predefined scripts. Next, we populated the tables with data.

I setup Apache hop to automate data loading from local machine to the oracle data warehouse. This process involved installing javascript on the local machine, and then installing the oracle database's wallet into the java script files, furthermore, i  integrated Apachehop  to the oracle cloud using the wallet.
Once the connection was setup, Apachehop was launched and checked for successful data connection.

**Building my First Dimensional Loader on Apache Hop:**

I tested this setup by using a node on Apache hop to view data in our data warehouse. Using a table input field, i set up a connection to our data warehouse. From this connection, i was able to use SQL statements to view rows we previously added to the data warehouse. Next, I added a text file output which allows to create a text output to save the table on my local machine. I connected the text file output to the table input field using a hop, then ran the entire pipeline to call DIM_Customer to view and save on my local machine.





**Loading Slowly Changing Dimensions**

The next task in the assignment involved loading slowly changing dimensions (SCDs) into the data warehouse. I updated the data warehouse with the `Product_Dimension` and `Customer_Dimension` tables, which included new fields.

To begin, I created an input node to read the source update file and transfer the data into the node. After verifying that all columns were correct and had the appropriate data types, the input node was connected to a Dimension Lookup node.

The Dimension Lookup node was configured to match the business key between the update data and the data warehouse. Additionally, we defined the type of slowly changing dimension operation to be performed for each row in the dimension table. Apache Hop provides various SCD options, such as:

- SCD Type 1 (Punch Through): Overwrites all records for the business key.
- SCD Type 2 (Insert): Adds a new row to represent the updated data.
- Update (Similar to SCD Type 1): Updates only the most recent record with the new data while leaving prior records unchanged.

A new field called `ISCURRENT` was also added to indicate whether a record is active or not.
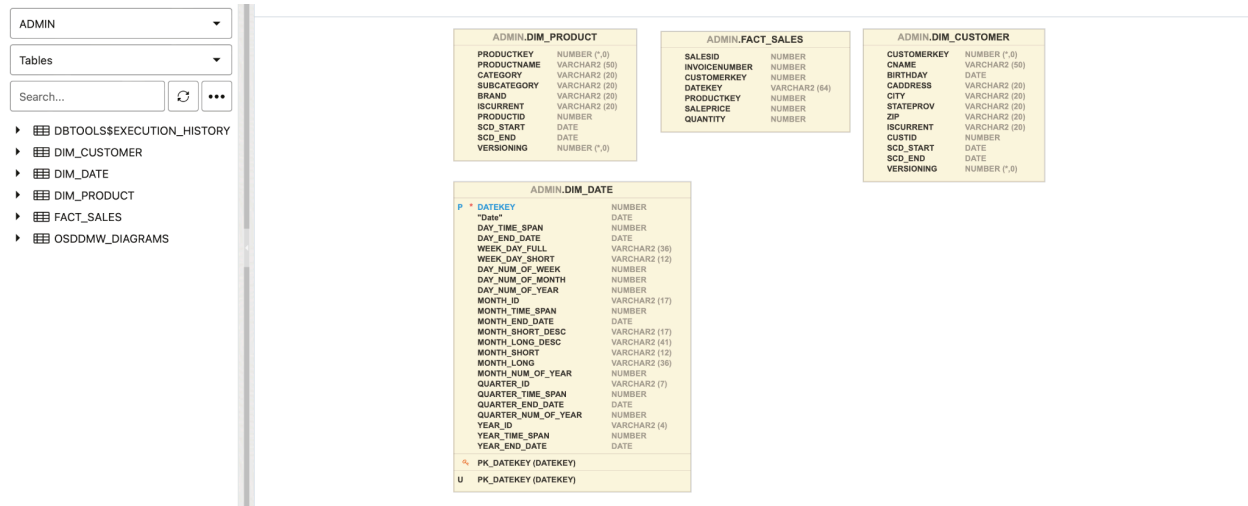
### Loading fact tables

After updating the dimension tables, I proceeded to load the Sales fact table. To accomplish this, I set up a pipeline that first verified a business key in the source sales update file. It then checked whether this record was current in the dimension tables. If the record was current, the pipeline retrieved the corresponding primary key and inserted it as a foreign key in the fact table.

**Introduction:**

**Generate a diagram of the database you created and attach a screenshot to your report.**



**Based on the diagram generated, what is this database missing that you'd expect to see? Why might it be missing this component? Use materials we've discussed in class and research (citing sources) to write no more than 2-4 sentences in response.**

There are a few more tables that would make the database complete such as the Store_Dimension table. This will be useful in determining which stores are performing interms of sales and profitability. We could also have tables on inventory which can help categorise items sales by categories. Structure-wise wise there are missing relationships or connections between primary and secondary keys between the tables and there are no pre-assigned Product keys.

**Loading your dimension tables:**

 **Open a new query and type "SELECT * FROM Dim_Date" and click "Execute". Paste a screenshot of the result in your report (with a sample of rows).**

| | DATEKEY | DATE | DAY_TIME_SPAN | DAY_END_DATE | WEEK_DAY_FULL | WEEK_DAY_SHOR | DAY_NUM_OF_WE | DAY_NUM_OF_MC | DAY_NUM_OF_YE/ | MONTH_ID | MONTH_TIME_SP, | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 20180101 | 1/1/2018, 12:00:00 A | 1 | 1/1/2018, 12:00:00 A | Monday | MON | 2 | 1 | 1 | JAN-2018 | 31 | 1, |
| 2 | 20180102 | 1/2/2018, 12:00:00 A | 1 | 1/2/2018, 12:00:00 A | Tuesday | TUE | 3 | 2 | 2 | JAN-2018 | 31 | 1, |
| 3 | 20180103 | 1/3/2018, 12:00:00 A | 1 | 1/3/2018, 12:00:00 A | Wednesday | WED | 4 | 3 | 3 | JAN-2018 | 31 | 1, |
| 4 | 20180104 | 1/4/2018, 12:00:00 A | 1 | 1/4/2018, 12:00:00 A | Thursday | THU | 5 | 4 | 4 | JAN-2018 | 31 | 1, |
| 5 | 20180105 | 1/5/2018, 12:00:00 A | 1 | 1/5/2018, 12:00:00 A | Friday | FRI | 6 | 5 | 5 | JAN-2018 | 31 | 1, |

**Run the "SELECT * FROM Dim_Date" query again and paste a screenshot in your report, along with a copy of the full updated script you ran.**

| | DATEKEY | DATE | DAY_TIME_SPAN | DAY_END_DATE | WEEK_DAY_FULL | WEEK_DAY_SHOR | DAY_NUM_OF_WE | DAY_NUM_OF_MC | DAY_NUM_OF_YE/ | MONTH_ID | MONTH_TIME_SP, | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 20160101 | 1/1/2016, 12:00:00 A | 1 | 1/1/2016, 12:00:00 A | Friday | FRI | 6 | 1 | 1 | JAN-2016 | 31 | 1, |
| 2 | 20160102 | 1/2/2016, 12:00:00 A | 1 | 1/2/2016, 12:00:00 A | Saturday | SAT | 7 | 2 | 2 | JAN-2016 | 31 | 1, |
| 3 | 20160103 | 1/3/2016, 12:00:00 A | 1 | 1/3/2016, 12:00:00 A | Sunday | SUN | 1 | 3 | 3 | JAN-2016 | 31 | 1, |
| 4 | 20160104 | 1/4/2016, 12:00:00 A | 1 | 1/4/2016, 12:00:00 A | Monday | MON | 2 | 4 | 4 | JAN-2016 | 31 | 1, |
| 5 | 20160105 | 1/5/2016, 12:00:00 A | 1 | 1/5/2016, 12:00:00 A | Tuesday | TUE | 3 | 5 | 5 | JAN-2016 | 31 | 1, |

**Once complete, run a "Select * from Dim_Product" query, take a screenshot, and paste the result in your report.**

🗑 ⓘ ↗ Download ▾ Execution time: 0.002 seconds

| | PRODUCTKEY | PRODUCTNAME | CATEGORY | SUBCATEGORY | BRAND | ISCURRENT | PRODUCTID | SCD_START | SCD_END | VERSIONING |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | *(null)* | Milk | Diary | Liquid | Buffalo Farms | Y | 2 | 2/1/2024, 12:00:00 A | 12/31/2099, 12:00:0( | 1 |
| 2 | *(null)* | Chocolate Chip Cool | Candy | Cookies | Nothing Breader | Y | 3 | 3/1/2024, 12:00:00 A | 12/31/2099, 12:00:0( | 1 |
| 3 | *(null)* | Eggs | Diary | Solid | Rochester Farms | Y | 4 | 4/1/2024, 12:00:00 A | 12/31/2099, 12:00:0( | 1 |
| 4 | *(null)* | Rotini | Wheat | Pasta | Buffalo Farms | Y | 5 | 5/1/2024, 12:00:00 A | 12/31/2099, 12:00:0( | 1 |
| 5 | *(null)* | Cinnamon Bread | Wheat | Bread | Nothing Breader | Y | 1 | 1/1/2024, 12:00:00 A | 12/31/2099, 12:00:0( | 1 |

**Back in Oracle Cloud, execute the following query: "SELECT * FROM Dim_Product". Take a screenshot of the results and place it into your report.**

**Loading your fact table:**

Once you think you have the lookups done correctly, you'll need to take a screenshot of your entire flow, along with screenshots of your second table input and second stream input nodes (from edit mode in each) and append them to your report.

**You might have noticed we're not doing a lookup for the date dimension. Write 1-2 sentences detailing why we don't need to. You'll be able to figure this out likely by looking at the data in the fact CSV and the date dimension.**

This is due to the format the data for the date_dimension table which will make it challenging to do lookup.

**Finally, take a screenshot of your output from the select query you ran, a screenshot of your completed pipeline, and attach it to your submission.**