# Data Pipeline
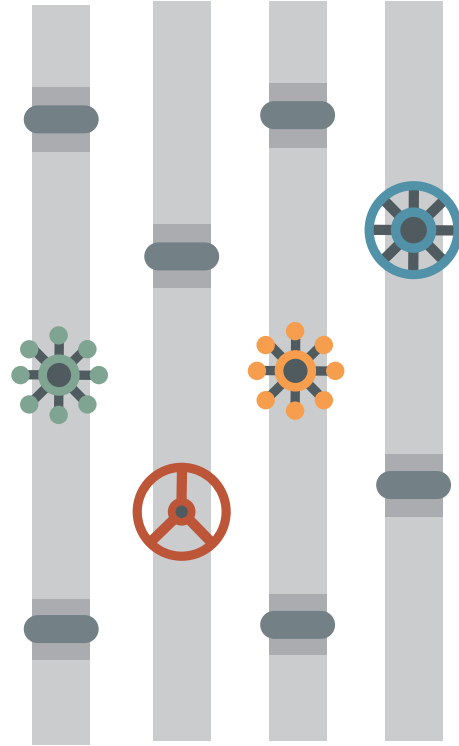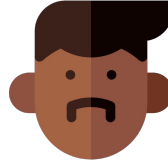
Module 2 Project - Group 4

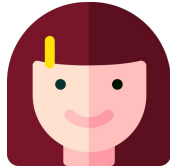# OUR TEAM

BAI HUIJING

TANG LAI LIN

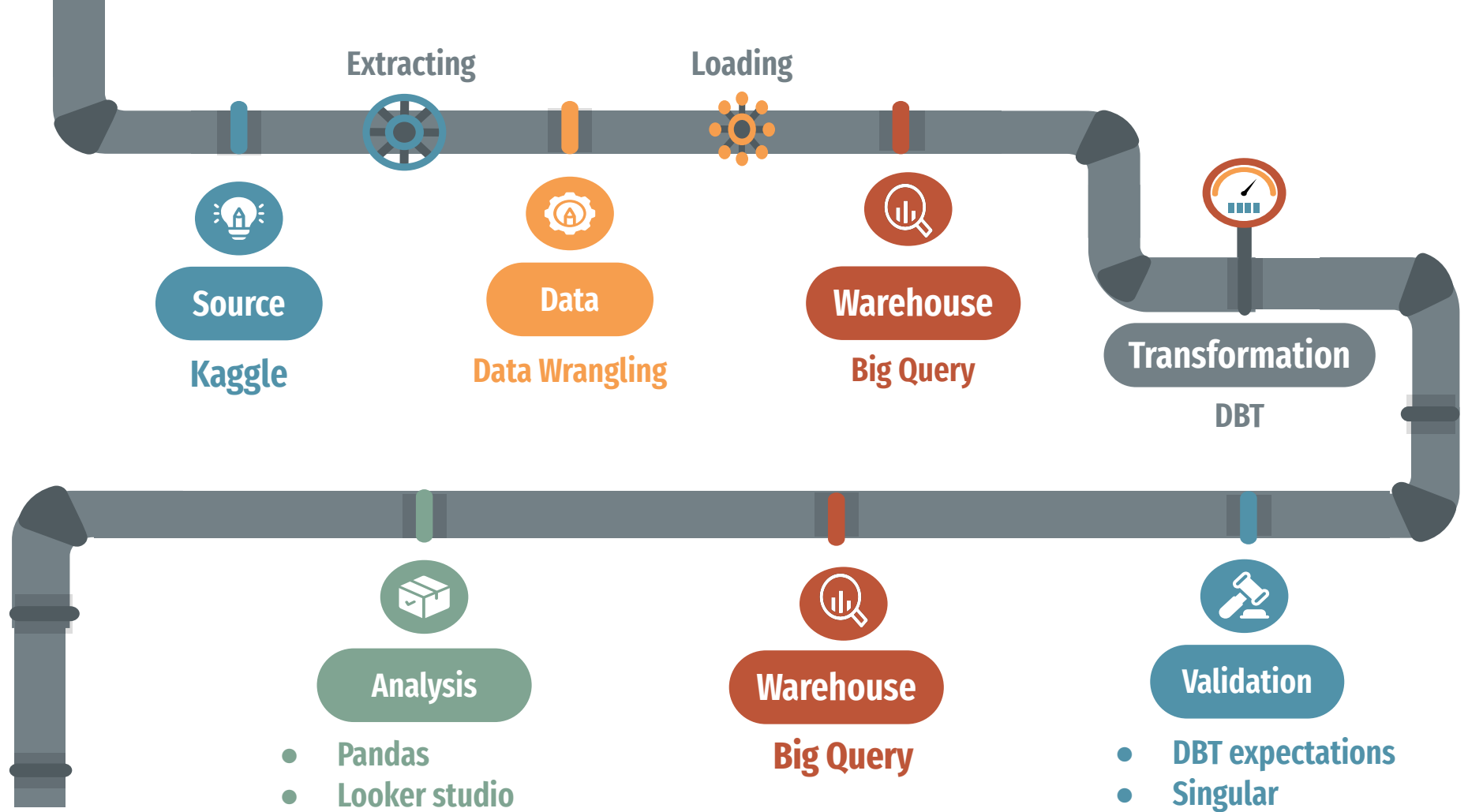MATHAVAN ARUGALAIMUTHU

LIM KAH HUI ESTHER

THOMAS TAY

**Extracting**

**Loading**

**Source**
Kaggle

**Data**
Data Wrangling

**Warehouse**
Big Query

**Transformation**
DBT

**Analysis**
- Pandas
- Looker studio

**Warehouse**
Big Query

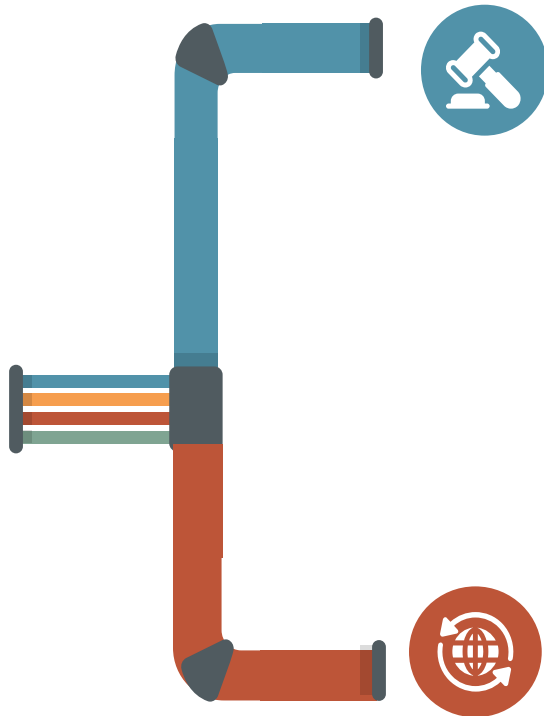**Validation**
- DBT expectations
- Singular

## Data wrangling

### olist_customers

### olist_order_items

### #1: Dropping of duplicates

- Duplicates are dropped (keep =first) from *'customer_id'* in olist_cutomers

- Duplicates are dropped (keep = last) from based on *'order_id'* as *order_item_id* is the list of quantity purchased

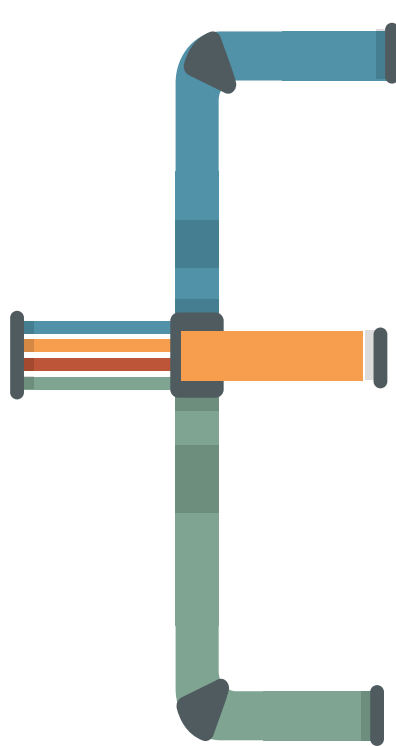### #2: Exporting dataframe to csv

- Clean data frame is exported to csv

**Data wrangling**

olist_payments

olist_orders

#1: Dropping of duplicates
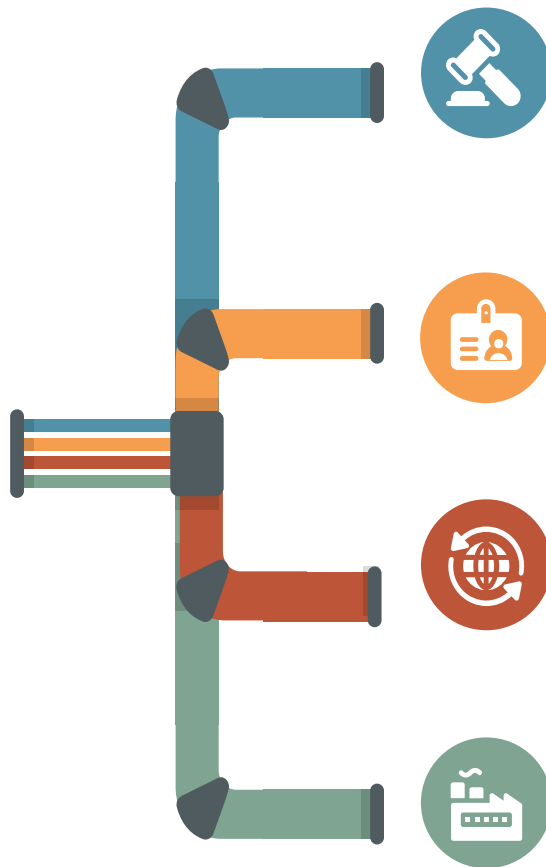
#2: Dropping rows with null values

#3: Exporting dataframe to csv

**Data wrangling**

olist_products

#1: Dropping of Rows/Columns
- Duplicated rows & rows with null values dropped
- Columns dropped:
  - *'product_name_lenght'*,
  - *'product_description_lenght'*
  - *'product_photos_qty'*

#2: Merging of dataframes
- Dataframe is merged with *product_category_name_translation.csv* using *'product_category_name'*.

#3: Replace null values in English column
- Some Portuguese words might go to English column, but can be fixed later
- Replace *'product_category_name'* with English column

#4: Exporting dataframe to csv
- Duplicates *'product_id'* are removed
- Set *'product_id'* as index
- Clean data frame is exported to csv

**ELT infrastructure**

# Data Warehouse : BigQuery vs DuckDb

## #1: Concurrency
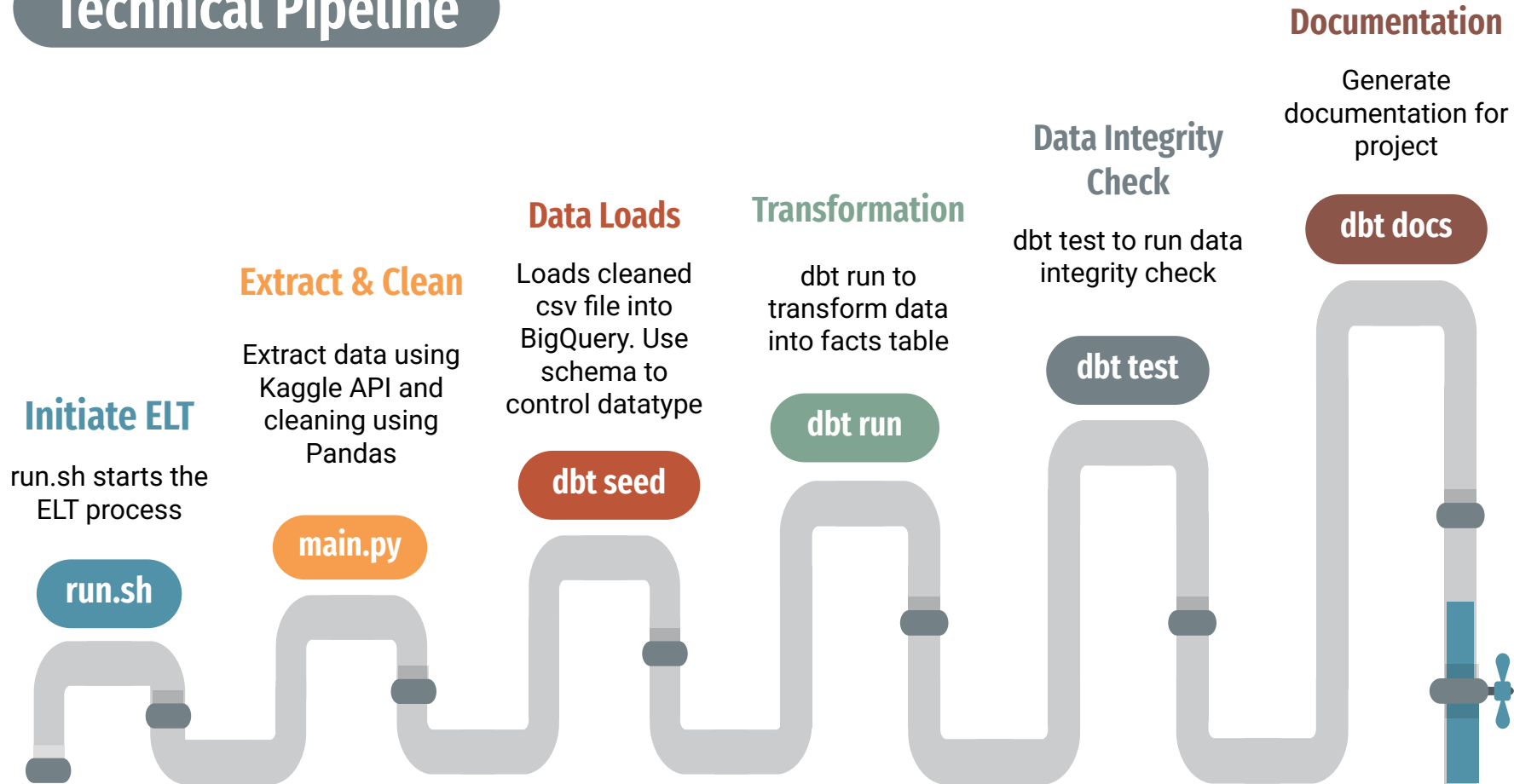- Challenge in managing DuckDb file lock

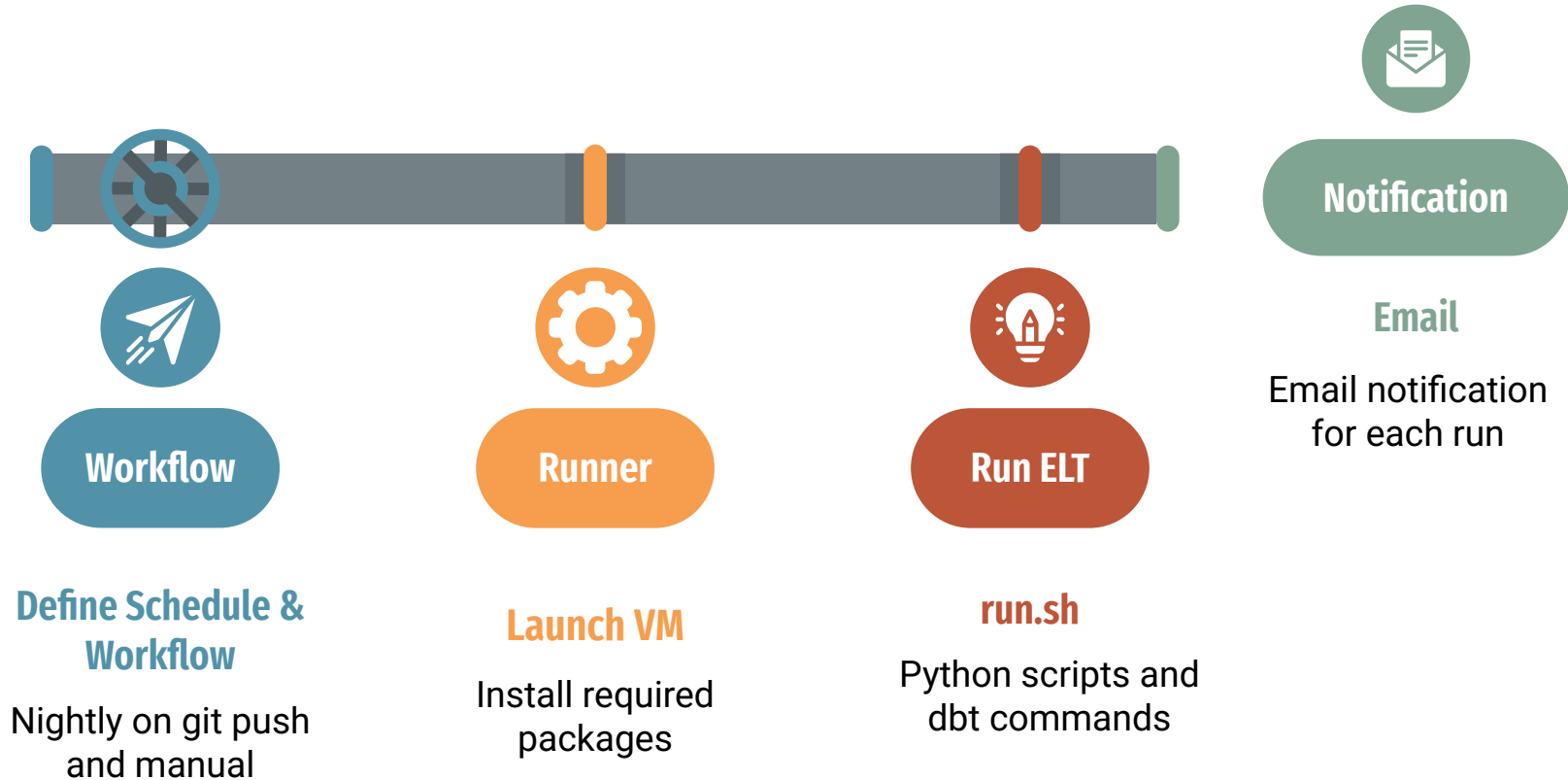## #2:  Scalability
- No scalability concern for BigQuery

## #3: Availability
- Available 24/7, include friendly tools: Looker Studio

# Technical Pipeline

## Initiate ELT

run.sh starts the ELT process

**run.sh**

## Extract & Clean

Extract data using Kaggle API and cleaning using Pandas

**main.py**

## Data Loads

Loads cleaned csv file into BigQuery. Use schema to control datatype

**dbt seed**

## Transformation

dbt run to transform data into facts table

**dbt run**

## Data Integrity Check

dbt test to run data integrity check

**dbt test**

## Documentation

Generate documentation for project

**dbt docs**

# Github Actions - CI/CD Platform



**Notification**

**Email**

Email notification for each run

**Workflow**

**Define Schedule & Workflow**

Nightly on git push and manual

**Runner**

**Launch VM**

Install required packages

**Run ELT**

**run.sh**

Python scripts and dbt commands

# GitHub Actions Implementations and Challenges

## #1: Ease of Implementation

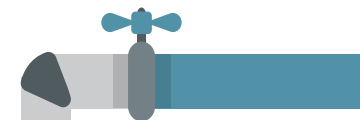- Only one yaml file required

## #2: Schedule Not Guaranteed

- Not suitable for mission critical job, need external scheduler

## #3: Challenge to Maintain Github Secrets

- Cannot edit secrets, app password can only be copied directly.

# DBT transformation | validation | documentation

**Projects**
dbt_ecomm | dbt_expectations

**Models (8)**
staging (5) | dimensions (2) | facts (1)

**Tests / data tests**
Generic tests (48) | singular tests (3)

**DBT Docs**
dbt docs generate | dbt docs serve

## dbt

Overview

Project    Database

Group

**Sources**
ecomm_raw

**Projects**
dbt_ecomm
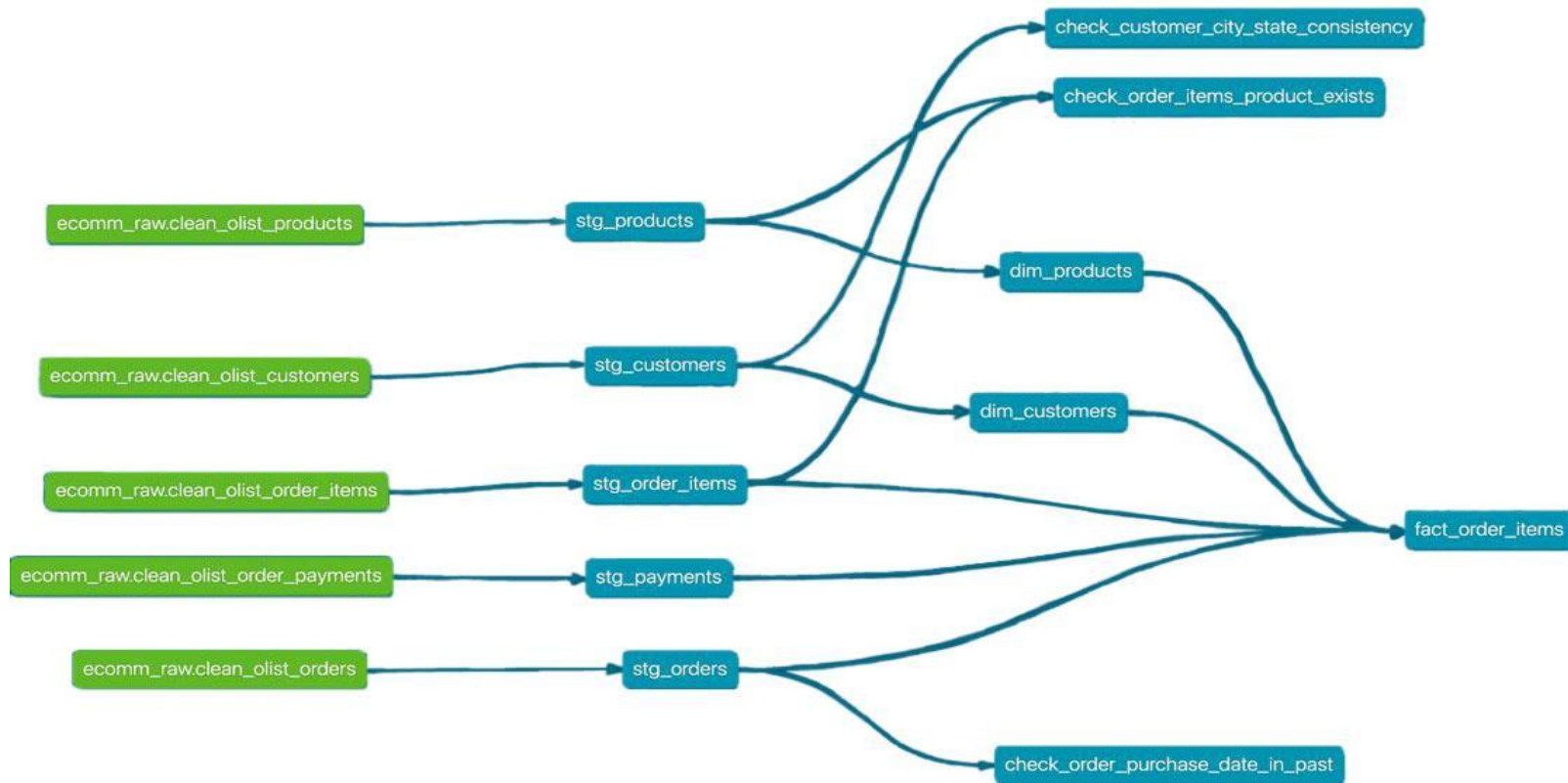models
seeds
tests
dbt_utils
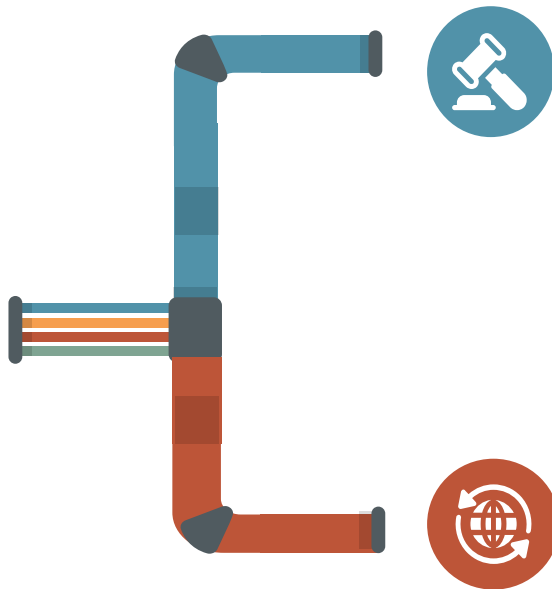dbt_date
dbt_expectations

# DBT lineage graph

**Source** · **Models** · **Tests**

- ecomm_raw.clean_olist_products → stg_products
- ecomm_raw.clean_olist_customers → stg_customers
- ecomm_raw.clean_olist_order_items → stg_order_items
- ecomm_raw.clean_olist_order_payments → stg_payments
- ecomm_raw.clean_olist_orders → stg_orders

Models: stg_products, stg_customers, stg_order_items, stg_payments, stg_orders, dim_products, dim_customers, fact_order_items

Tests: check_customer_city_state_consistency, check_order_items_product_exists, check_order_purchase_date_in_past

# DBT validation

## Test cases (51)

- Generic test cases: 48 (Unique, not_null, accepted_values, dbt expectation) focus on validation
- Singular test case: 3 focus on integrity

## Data quality

## Data integrity

## Failed test case (4)*

- Accepted_values mismatch
- Order lists product_id mismatch with product's product_id

```
Finished running 48 data tests in 0 hours 0 minutes and 34.36 seconds (34.36s).

Completed with 2 errors, 0 partial successes, and 0 warnings:

Failure in test check_customer_city_state_consistency (tests/singular/check_customer_city_state_consistency.sql)
  Got 3015 results, configured to fail if != 0

  compiled code at target/compiled/dbt_ecomm/tests/singular/check_customer_city_state_consistency.sql

Failure in test check_order_items_product_exists (tests/singular/check_order_items_product_exists.sql)
  Got 19730 results, configured to fail if != 0

  compiled code at target/compiled/dbt_ecomm/tests/singular/check_order_items_product_exists.sql
```

*bugs will be raised for the 4 failed test cases

# DBT documentations

**Docs generation**

**Docs serve**

**#1: Docs generation (catalog.json)**

## ecomm_raw.clean_olist_products source table

Details    Description    Columns    Referenced By    SQL

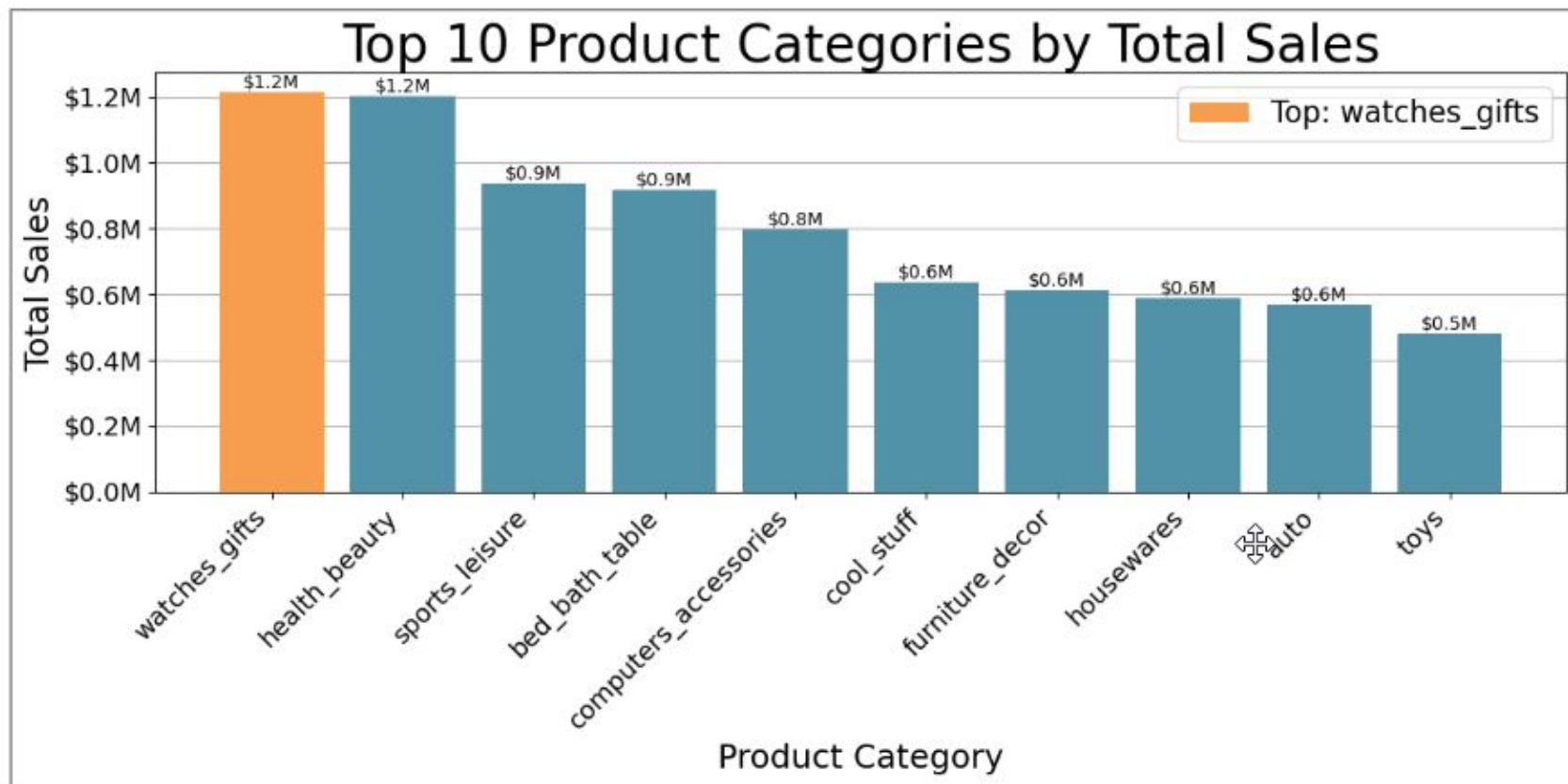### Details

| TAGS | OWNER | TYPE | PACKAGE | RELATION | CONTRACT | LOADE |
|------|-------|------|---------|----------|----------|-------|
| untagged | | table | dbt_ecomm | sctp-data-eng-ecomm.ecomm_raw.clean_olist_products | Not Enforced | |

| # ROWS | APPROXIMATE SIZE |
|--------|------------------|
| 25,981 | 2 MB |

**#2: Docs serve (localhost)**

Top 10 Product Categories by Total Sales

Average Installments per Payment Type

# Total Sales by Day of the Week



Total Sales by Day of the Week. Bar chart showing Total Sales (y-axis) by Day of the Week (x-axis). Sunday $1.5M, Monday $2.1M, Tuesday $2.1M, Wednesday $2.0M, Thursday $1.9M, Friday $1.8M, Saturday $1.5M. Highest Sales: Monday.

# Looker Studio

# Looker Studio



- Web-based data visualization tool
- create interactive dashboards and reports from various data sources
- To explore, visualize, and share insights easily.

# Improvements

**1** **Github Secrets**
- For better security convert service key to Github Secrets
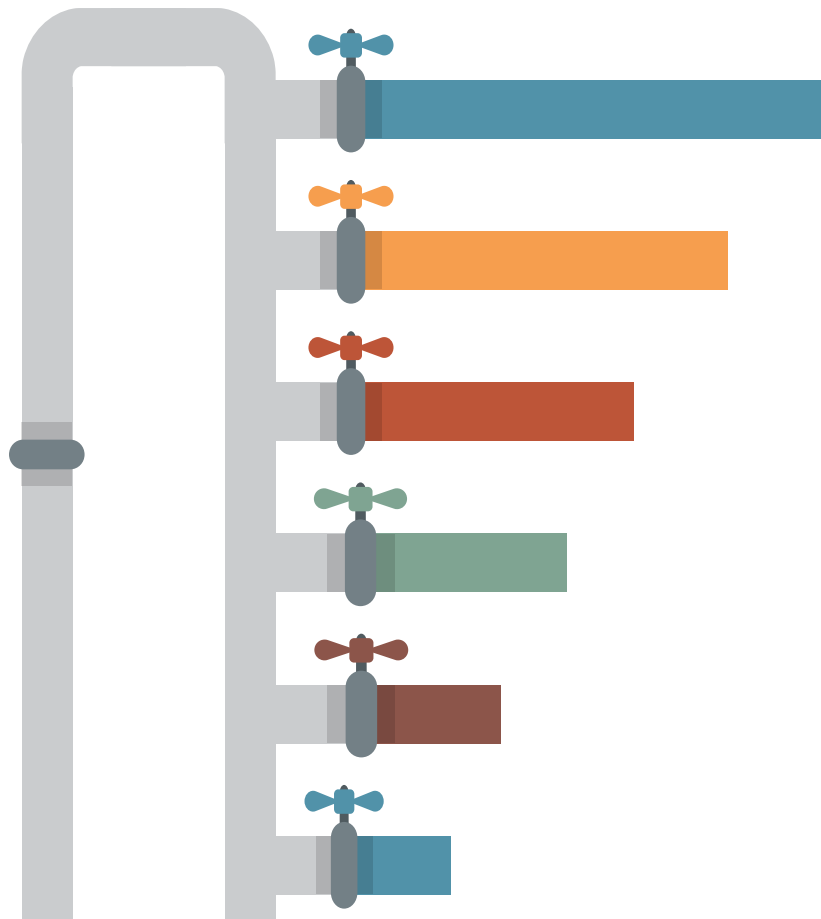
**2** **Google Job Scheduler**
- Use Google job scheduler to guarantee workflow runs on time

**3** **dbt Macro and dbt Unit Test**
- Improve data quality and integrity

# References

**rudderstack**
https://www.rudderstack.com/guides/how-to-access-and-query-your-bigquery-data-using-python-and-r/

**kaggle**
https://www.kaggle.com/datasets/olistbr/brazilian-ecommerce?select=olist_customers_dataset.csv

**mbvyn.medium**
https://mbvyn.medium.com/understanding-dbt-seeds-and-sources-c5611be17d32

**drlee.io**
https://drlee.io/unlock-the-power-of-kaggle-how-to-call-datasets-directly-with-python-and-rule-the-data-science-2303cbae4e8a

**Chatgpt**
https://chatgpt.com/share/67e4e54f-b340-8000-8d01-fce47c0f4870

**Slides template**
Business Pipeline Infographics by Slidesgo

# THANKS!