# OLIST E-COMMERCE DATA PIPELINE

Module 2 – End-to-End ELT & Analytics | DuckDB · Parquet · GCS · BigQuery · dbt · Looker Studio / Plotly Dash

# AGENDA

- Executive Summary:
    - Purpose, focus and assumption
    - Problem statement & goals
- Data Architecture
- Implementation of ELT pipeline
- Data quality & testing
- Data Insights: Data Analysis with Python and BI Interface: Looker Studio / Plotly Dash
    - Key metrics dashboard
    - E-commerce sales insights
    - Analysis and findings
- Business recommendations
- Phase 2: get ready for operation with Pipeline Orchestration

# BUSINESS SETTING

Olist is a Brazilian **e-commerce marketplace**

The challenge:

With thousands of sellers and millions of orders, Olist needs a robust **data platform** to understand customer behavior, seller performance, and logistics bottlenecks etc

# PROBLEM & GOALS

Data given: 9 Olist CSVs raw data available

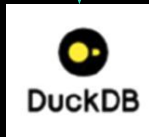Outcome required: Transformed into analytics-ready tables & insights in a dashboard

- Reproducible pipeline required: Raw → Marts
- Clean Staging: Renames, Cast Types, Dedupe, Normalization
- EDA and Quality Gates
- Star schema for BI
- BI Dashboards for easy KPI tracking & Business Decisions
- Orchestration to maintain the data in real time  (Phase 2)
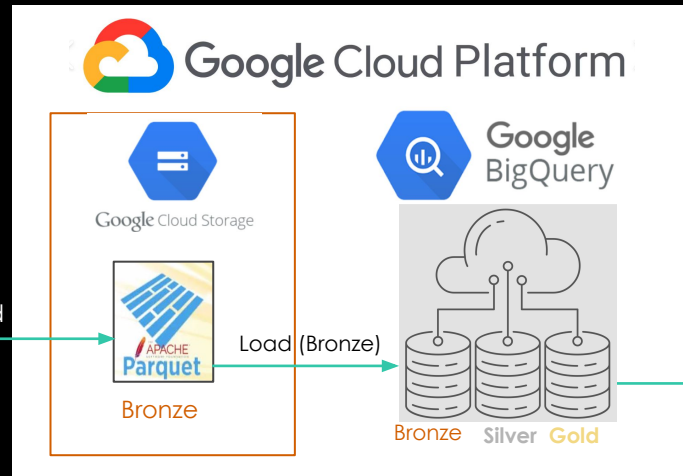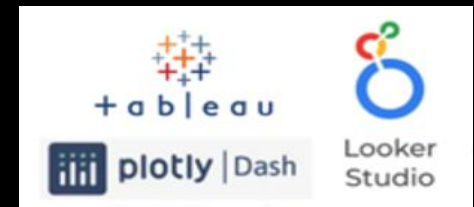
# ARCHITECTURE

**Local data source**

**Data Warehouse**

**Data Analysis**



Extract

Convert & Load

Load (Bronze)

Gold Layer (Business-ready tables)

Bronze

Bronze  Silver  Gold

Transform & Test
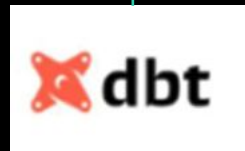
- Start with Kaggle **CSVs** locally. Using **DuckDB**, convert CSVs to **Parquet** (columnar, compressed).

- Those Parquet files are then **uploaded to GCS** as object storage for raw Bronze tables.

- GCS is the middle layer stores Parquet files and act as the source from which BigQuery ingest data.
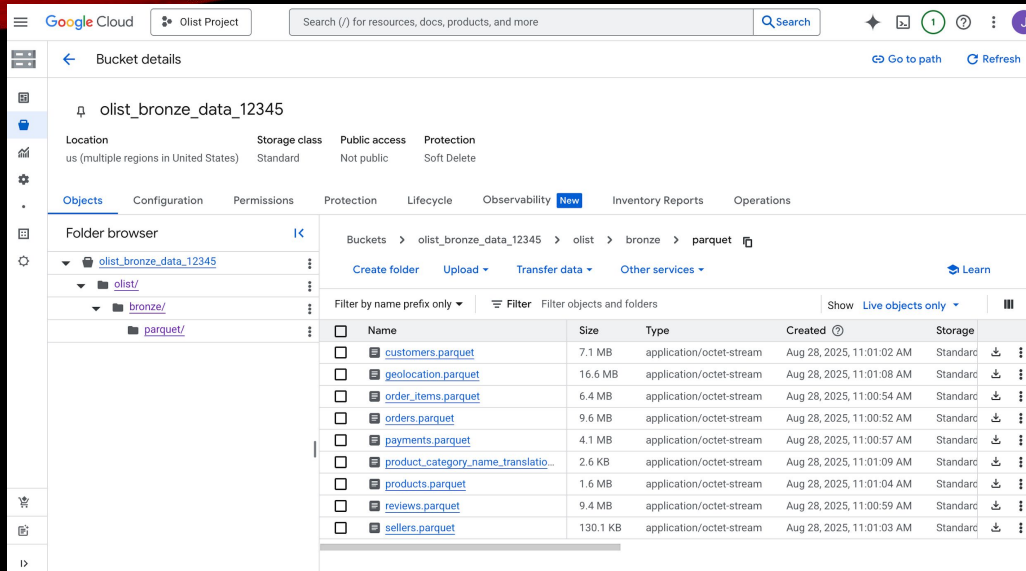
- dbt takes those raw Bronze tables and **transforms** them into clean, conformed Silver models.

- dbt runs on local machine, Transformations inside BigQuery, which is data warehouse.

- Keep all **transformations in BigQuery** for performance, governance, and BI connectivity.

At this stage, we are working with static Kaggle CSVs

- Load CSV into **DuckDB** for type-checking and quick validation

- Convert to **Parquet** for efficient columnar storage and store in **GCS**

- Upload Parquet files to **BigQuery** (`olist_bronze` dataset)

Result: *olist_bronze* dataset in BigQuery with 9 raw files

Why Parquet?

- Compressed, efficient for analytics

- Schema preserved, faster to load into BigQuery

- Scales better than raw CSV

### Role of GCS in ingestion
**1.Landing zone for Bronze**
- Start with Kaggle CSVs locally.
- Using DuckDB/Python, we convert them to **Parquet** (columnar, compressed).
- Those Parquet files are then **uploaded to GCS**.
- At this stage, GCS is acting as your data lake for raw Bronze data.

**2.Staging point for BigQuery**
- BigQuery can't read local **files** directly — it integrates naturally with GCS.
- Our notebook uses the **BigQuery Python client** (or CLI) to LOAD or create Native Tables in BigQuery pointing at those Parquet files in GCS.

**3.Benefit: Decoupling compute & storage: GCS stores the data files**, **BigQuery ingests them**
- By keeping Bronze in GCS, we don't overload BigQuery with raw data. If something goes wrong, we can always re-load from GCS.
- We can reprocess/reload from Parquet in GCS at any time without re-downloading the Kaggle dataset.
- It also sets us up for **incremental ingestion** later (partition by date in GCS, load only new files).

9 dbt **staging views** in `olist_silver` and 1 intermediate layer: `int_zip_prefixes` (aggregates geolocation ZIP codes to lat/lng)

Cleaning steps:
- Normalized column names & types
- Deduplication (`stg_reviews`)
- Payment normalization (`stg_payments` → bucket rare values to `other`/`unknown`)
  - Provides a clean, consistent interface for downstream models.

**Why needed intermediate layer ?**

•**Raw geolocation (Bronze)**:
  •every customer/seller record comes with a **zip prefix** (lots of repeats, sometimes dirty, inconsistent).
•**Intermediate layer (int_zip_prefixes)**:
  •Deduplicates all unique zip_prefix values
  •Joins them to their **city, state, lat/lng** from the geolocation dataset
  •Produces a **clean reference table** with one row per zip_prefix

- Customers and sellers both point to zip prefixes → instead of repeating city/state in every record, we normalize it.
- This makes **joins simpler and consistent** in Gold (facts/dims).
- Keeps **dim_customer** and **dim_seller** slim and tidy.

👉 So the intermediate layer is just a **clean bridge** between raw zip prefixes and Gold dimensions.

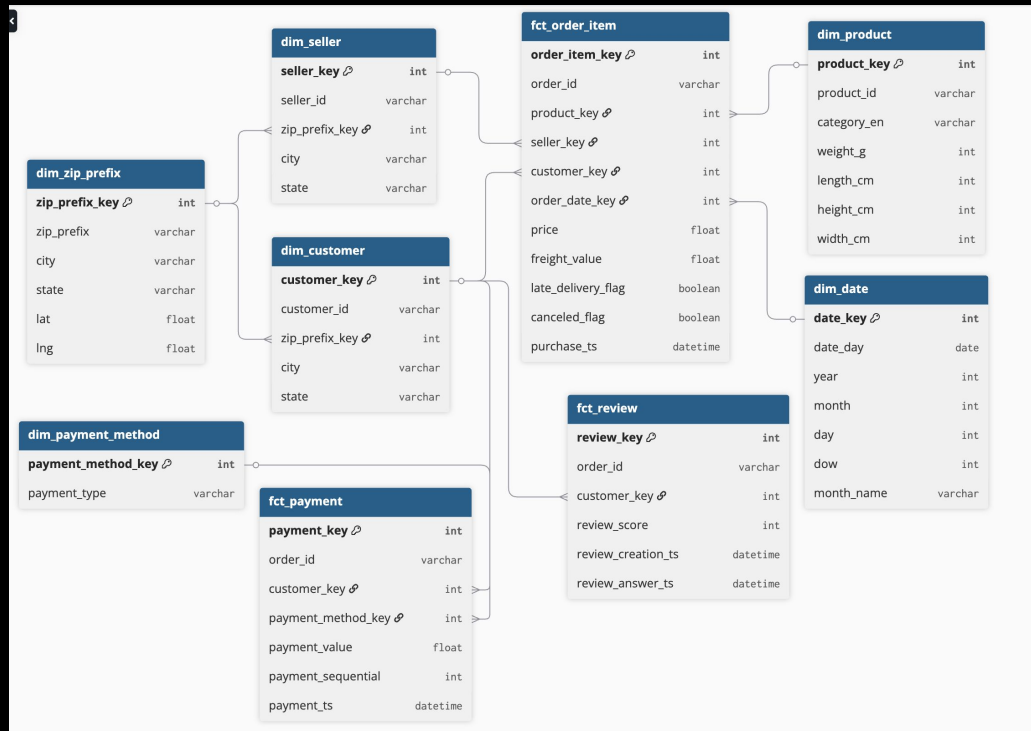**Gold (Business Layer)**
**6 Dimensions**
`dim_customer`, `dim_product`, `dim_seller`,
`dim_zip_prefix`, `dim_date`, `dim_payment_method`
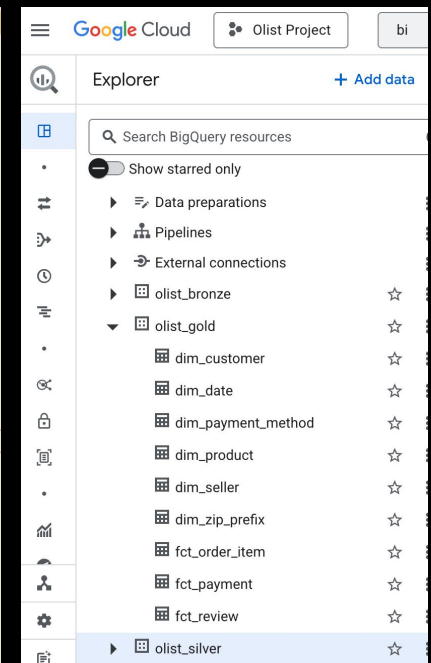
**3 Facts**:
`fct_order_item`, `fct_payment`, `fct_review`

Materialized as **tables** in `olist_gold`
Star schema design supports BI dashboards &
analytics (refer to star schema in next slide)

| Finding (Initial EDA) | Data Cleaning Action |
|---|---|
| Repeated / inconsistent ZIP prefixes | Built int_zip_prefixes → dim_zip_prefix for normalized reference |
| Null / invalid payment types | Created dim_payment_method, flagged invalid/null values |
| Duplicate reviews, skewed ratings | Deduplicated review_id; flagged outliers in review scores |
| Canceled orders & invalid date sequences | Added canceled_flag, late_delivery_flag; validated date fields |
| Missing customer/seller location values | Standardized ZIP joins; filled/handled missing city/state |
| Outliers in freight values | Applied dbt-expectations range tests; flagged anomalies |

## Ensuring Data Quality & Testing

**Data Testing with dbt:** **tests** are in YAML files (dims.yml, facts.yml, staging.yml)

**Not Null & Unique**
    Primary keys in dim_customer, dim_product, dim_seller, fct_order_item, etc.
**Relationships**
    Foreign keys like fct_order_item.customer_key → dim_customer.customer_key
    fct_order_item.product_key → dim_product.product_key
**Accepted Values**
    payment_type in stg_payments (credit_card, boleto, voucher, debit_card)
**dbt-expectations**
    run expect_column_values_to_be_in_set on payment_type
    had tests for late_delivery_flag values (0/1)

✅ Result: A **trusted Silver/Gold layer** with enforced integrity and test coverage for downstream analysis & BI.

# BI DASHBOARD 2 (PLOTLY DASH)

# TABLEAU

# INSIGHTS & RECOMMENDATIONS

1) Limited repeat customers

   - Loyalty program to manage customer ecosystem

 2) Logistics bottlenecks

   - Partnerships with delivery partners

   - Tracking of service levels

 3) 18 Sellers generate 80% of GMV – concentration risk

   - Develop seller ecosystem and incentives to diversify

 4) Regional concentration of GMV in 3 states

   - Market study to develop market in rest of states

# BUSINESS VALUE

Data-driven approach

- Ability to track KPIs in real time
- Receive feedback from the market in real time
- Provide benchmark for future strategy testing

When we move to live or updated sources, we will :

- Add Snapshot, allow us to capture changes in slowly changing dimensions (SCD Type 2).
- Use Meltano to manage extract & load into BigQuery.
- Orchestrate Pipeline with Dagster to ensure:
  - Bronze ingestion → Silver cleaning → Gold marts run in sequence
  - Snapshots run before transformations
  - Automated tests + alerts on failures
- Continue connecting gold layer to BI Tools for updated sources

**Meltano** to manage extract & load into BigQuery

**Local data source**

**Data Warehouse**

**Data Analysis**

Extract

Convert & Load

Load (Bronze)

Bronze

Bronze Silver Gold

Business table (Gold)

Transform & Test

**Dagster** handles: **when** to run, **in what order**, retries, sensors, backfills.

**dbt** handles: **what** SQL to run, model lineage, tests, snapshots, docs