# Google Play Store
## Data Pipeline Report

*Lab 2 — Data Pipeline with dbt & DuckDB*

January 2026

Imad Absri

Omar Bellmir

# Contents

# List of Figures

# List of Tables

# Introduction and Objectives

The Google Play Store hosts millions of applications and generates an enormous volume of user feedback every day. Turning this raw stream of reviews into structured, queryable, and visually explorable analytics is precisely the challenge this second data engineering lab was designed to address.

Where the first lab built the pipeline end-to-end in ad-hoc Python—mixing ingestion, transformation, and serving logic in a single script—this lab elevates the architecture by introducing industry-standard tooling. Transformation logic moves entirely into **dbt Core**, all data is stored and queried inside **DuckDB**, and the final analytics-ready tables are consumed by **Power BI**. The result is a pipeline that is modular, testable, and repeatable in a way that raw Python scripting cannot match.

The lab guidelines define six concrete objectives that drove every design decision we made:

---

**Lab Objectives**

**1.** Install and configure the development environment (DuckDB, dbt-core, dbt-duckdb adapter).

**2.** Explore dimensional data modelling for analytics (star/snowflake schema, Kimball methodology).

**3.** Structure a data pipeline using dbt model layers — staging flowing into marts.

**4.** Separate raw data, transformation logic, and serving tables with clear physical boundaries.

**5.** Apply data quality tests using dbt's native testing framework.

**6.** Prepare a stable serving layer consumable by BI dashboards (Power BI).

---

Each chapter of this report maps directly to one or more of these objectives and documents both our implementation decisions and the concrete outputs they produced. A compliance matrix in Chapter 8 cross-references every guideline requirement against its evidence in the codebase.

# Architecture and Technology Stack

## 2.1  High-Level Architecture

The pipeline separates three distinct phases, each handled by a dedicated component, with clean handoff points between them.
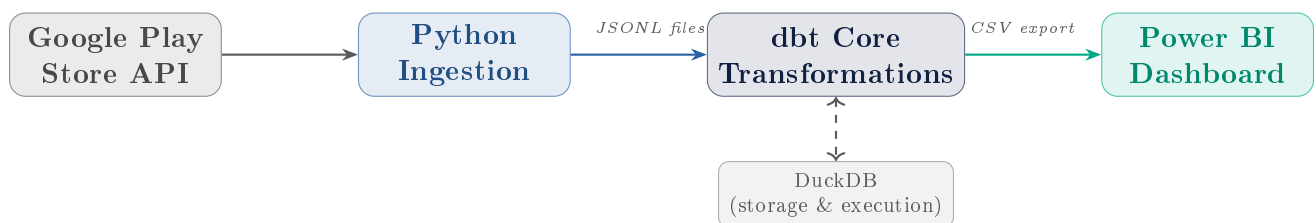
**Figure 2.1:** *End-to-end pipeline architecture showing the three phases and their data handoffs*

**Phase 1 — Ingestion.** The Python script `ingest.py` calls the `google-play-scraper` library against a curated list of 20 popular applications. For each app it collects app metadata (title, developer, category, average score, total ratings) and user reviews filtered to the last 100 days, capped at 5 000 reviews per application. Both outputs are saved as JSONL files in `data/raw/` and are never modified again—serving as an immutable raw landing zone, exactly as the guidelines specify.

**Phase 2 — Transformation.** dbt orchestrates all transformation logic, and DuckDB executes it. DuckDB reads the JSONL files directly without any intermediate loading step, then progressively cleans and materialises the data through staging views and mart tables forming a complete snowflake schema.

**Phase 3 — Serving.** The script `export_to_powerbi.py` connects to the DuckDB database, queries each mart table, and writes CSV files to `powerbi_export/`. These files are then imported into Power BI where relationships and visualisations are configured interactively.

## 2.2  Technology Stack

**Table 2.1:** *Technology stack and pinned versions used in this project*

| Component | Version | Role | Notes |
| --- | --- | --- | --- |
| **Python** | 3.13 | Ingestion & export scripting | Reused from Lab 1 |
| **google-play-scraper** | latest | Play Store API client | Handles pagination |
| **DuckDB** | 1.4.4 | Analytical database engine | Reads JSONL natively |
| **dbt-core** | 1.11.6 | Transformation orchestration | SQL-only models + YAML tests |
| **dbt-duckdb** | 1.10.1 | DuckDB adapter for dbt | Bridges dbt and DuckDB |
| **Power BI Desktop** | latest | BI dashboard | Imports CSV, defines relationships |

The choice of DuckDB as the execution engine is deliberate. It is an in-process OLAP database that requires no server setup, reads JSONL files natively, and achieves excellent analytical query performance on a standard laptop. Combined with dbt, it creates a local stack that behaves exactly like a cloud data warehouse—making it ideal for a reproducible lab environment.

# Dimensional Data Modelling

## 3.1 Applying the Kimball Methodology

The guidelines explicitly require applying the Kimball four-step process before writing any SQL. We worked through each step against the Google Play Store dataset and document the outcomes below.

### 3.1.1 Step 1 — Business Process

The business process under analysis is *user review submission*. Every time a user publishes a review on the Google Play Store, a discrete and measurable event occurs, carrying quantitative feedback (score, thumbs-up count) alongside rich contextual information: which application, which category, which developer, and when.

### 3.1.2 Step 2 — Grain

One row in the fact table represents **a single user review posted for one specific application at one specific point in time**. This is the finest grain available in the dataset and allows every conceivable analytical aggregation without loss of information.

### 3.1.3 Step 3 — Dimensions

Table 3.1: *Dimension identification: business questions answered and key attributes*

| Dimension | Business question | Key attributes |
|---|---|---|
| dim_apps | Which app was reviewed? | app_id, app_name, price, avg_score, total_ratings |
| dim_categories | Which category does the app belong to? | category_id, category_name |
| dim_developers | Who built the app? | developer_id, developer_name |
| dim_date | When was the review posted? | date_day, year, quarter, month, week_of_year, day_name, is_weekend |

### 3.1.4   Step 4 — Facts (Measures)

Two quantitative measures are stored in the fact table. `review_score` is an integer from 1 to 5, well-suited for averaging and distribution analysis. `thumbs_up_count` is a non-negative integer, meaningful for summing to measure community engagement with specific reviews. Both are captured directly from the raw review payload with no derivation required.

### 3.1.5   Bus Matrix

**Table 3.2:** *Kimball Bus Matrix: dimension availability per business process*

| Business Process | dim_apps | dim_categories | dim_developers | dim_date |
|---|:---:|:---:|:---:|:---:|
| User review submission | ✓ | ✓ | ✓ | ✓ |

## 3.2   Schema Design: Snowflake Variation

The schema we implemented is a **snowflake schema**. The central `fact_reviews` table joins to `dim_apps` via a surrogate key, and `dim_apps` in turn references `dim_categories` and `dim_developers` as separately normalised tables rather than embedding their attributes directly. This adds a controlled level of normalisation that keeps each dimension table focused on a single concept, while Power BI's relationship engine traverses the chain of foreign keys automatically.



**Figure 3.1:** *Power BI relationship diagram showing the complete snowflake schema with all six tables*

Figure 3.1 shows the six tables in Power BI's model view. The central `fact_reviews` table holds foreign keys to `dim_apps` (via `app_sk`) and `dim_date` (via `date_key`). The `dim_apps` dimension connects outward to `dim_categories` and `dim_developers`. The additional `dim_apps_scd` table captures SCD Type 2 history and sits alongside `dim_apps` in the model.

# Pipeline Implementation

## 4.1   Data Ingestion

The ingestion script targets 20 major Android applications spanning diverse categories: social media (WhatsApp, Instagram, Facebook, Snapchat, TikTok), entertainment (Netflix, YouTube, Spotify), productivity (Gmail, Microsoft Teams, Duolingo), commerce (Amazon Shopping, PayPal, Airbnb, Uber), and community platforms (Twitter/X, LinkedIn, Reddit, Discord).

For each application the script fetches the metadata record with a single API call, then paginates through reviews in batches of 200 sorted by most recent, stopping as soon as a review older than 100 days is encountered or the 5 000-review cap is reached. The two resulting files are saved into `data/raw/` and treated as immutable from that point forward.

## 4.2   dbt Project Structure

The dbt project follows the folder convention recommended in the guidelines, with staging views and mart tables kept both physically and logically separate:

```
playstore_pipeline/
  dbt_project.yml                  # dbt configuration and materialisation conventions
  profiles.yml                               # DuckDB connection profile
  data/
    raw/                               # Immutable JSONL files (ingest.py output)
    db/                                  # playstore.db (DuckDB database file)
  models/
    staging/                   # Views: stg_playstore_apps, stg_playstore_reviews
    marts/
      dimensions/           # dim_apps, dim_categories, dim_developers, dim_date
      facts/                                         # fact_reviews
    snapshots/                       # SCD Type 2 snapshot for dim_apps
```

The `dbt_project.yml` enforces materialisation conventions at the folder level: staging models build as lightweight `view`s (no data duplication, always reflecting the raw files),

while everything in the marts folder materialises as a physical `table`, pre-computed for fast BI queries.

## 4.3  Staging Layer

The staging layer performs type casting, column renaming, surrogate key generation using `md5()`, and null filtering. No business logic or aggregation is applied here—the clean staging layer principle from the guidelines is strictly respected.

Listing 4.1: *stg_playstore_reviews.sql — staging model for user reviews*

```
1  WITH source AS (
2      SELECT * FROM {{ ref('src_playstore_reviews') }}
3  )
4  SELECT
5      md5(CAST(reviewId AS VARCHAR))  AS review_sk,
6      reviewId::VARCHAR               AS review_id,
7      appId                           AS app_id,
8      userName                        AS user_name,
9      CAST(score AS INT)              AS review_score,
10     content                         AS review_text,
11     CAST("at" AS TIMESTAMP)         AS review_at,
12     CAST(thumbsUpCount AS INT)      AS thumbs_up_count
13 FROM source
14 WHERE reviewId IS NOT NULL
```

Surrogate keys are generated with `md5()`, producing a deterministic 32-character hash from the natural key. This is consistent with Kimball's recommendation: surrogate keys must be stable, unique, and independent of the source system. On each full refresh, the same input produces the same surrogate key, ensuring idempotent joins downstream.

## 4.4  Dimension Tables

### 4.4.1  dim_apps

`dim_apps` is the central dimension. It joins the apps staging model with the category and developer surrogate keys through `LEFT JOIN`s on their natural keys.

Listing 4.2: *dim_apps.sql — central dimension model joining apps, categories, and developers*

```
1  WITH apps AS (SELECT * FROM {{ ref('stg_playstore_apps') }}),
2      categories AS (SELECT * FROM {{ ref('dim_categories') }}),
3      developers AS (SELECT * FROM {{ ref('dim_developers') }})
4  SELECT
5      a.app_sk, a.app_id, a.app_name,
```

```
6        a.price, a.avg_score, a.total_ratings,
7        c.category_sk,
8        d.developer_sk
9   FROM apps a
10  LEFT JOIN categories c ON a.category_id = c.category_id
11  LEFT JOIN developers d ON a.developer_id = d.developer_id
```



**Figure 4.1:** `dim_apps` *in Power BI: 20 rows with MD5 surrogate keys, average scores, and total rating counts*

Figure 4.1 shows the 20 application records loaded in Power BI. Average scores range from approximately 3.65 (Twitter/X) to 4.65 (Duolingo), and total rating counts span from 3.4 M (PayPal) to 225 M (WhatsApp Messenger)—reflecting the enormous variation in user base sizes across the dataset.

## 4.4.2   dim_categories and dim_developers

These two dimensions are simple deduplication tables derived from the staging apps model. `dim_categories` resolves 12 unique Play Store category identifiers into human-readable names, while `dim_developers` holds 17 distinct developer entities.



**(a)** `dim_categories`: *12 unique categories*



**(b)** `dim_developers`: *17 unique developer entities*

**Figure 4.2:** *Category and developer dimension tables loaded in Power BI*

### 4.4.3 dim_date

Per the Kimball best practice noted in the guidelines, the date dimension uses an integer key in `YYYYMMDD` format (e.g., `20251113`). The table is generated programmatically using a SQL range spanning the minimum and maximum review dates found in the staging model, producing one row per calendar day with attributes for year, quarter, month, week of year, day name, and a weekend flag.



**Figure 4.3:** *`dim_date`: continuous calendar from November 2025 through February 2026 with full time attributes*

Figure 4.3 shows the date dimension spanning the review period. The integer key format makes joins efficient and values human-readable directly in the fact table without any additional date formatting in the BI tool.

## 4.5   Fact Table

`fact_reviews` is materialised at the declared grain—one row per review event—and stores foreign keys to `dim_apps` and `dim_date` alongside the two measures. Joining to `dim_date` requires converting the review timestamp to the integer `YYYYMMDD` key format:

**Listing 4.3:** *Excerpt from `fact_reviews.sql` showing key extraction and grain enforcement*

```sql
SELECT
    r.review_sk,
    r.review_id,
    a.app_sk,
    CAST(STRFTIME(r.review_at, '%Y%m%d') AS INTEGER)  AS date_key,
    r.review_at,
    r.review_score,
    r.thumbs_up_count,
    r.user_name,
```

```
10        r.review_text
11  FROM stg_reviews r
12  INNER JOIN dim_apps a ON r.app_id = a.app_id
13  WHERE a.app_sk IS NOT NULL
```



**Figure 4.4:** `fact_reviews` *in Power BI: timestamps from 18 February 2026 confirm live data capture from the Play Store*

The fact table (Figure 4.4) contains timestamps from 18 February 2026, confirming that the pipeline successfully captured very recent data from the live Play Store API. The rows visible correspond to WhatsApp Messenger reviews with `review_score = 5`, representing the freshest submissions that had not yet accumulated community votes.

## 4.6   Executing the Full Pipeline

The complete pipeline is reproduced by executing five commands in sequence from the project root:

**Listing 4.4:** *Full pipeline execution sequence*

```
1  # 1. Scrape raw data from the Google Play Store
2  python ingest.py
3
4  # 2. Navigate into the dbt project directory
5  cd playstore_pipeline
6
7  # 3. Build all dbt models (full refresh rebuilds all tables)
8  dbt run --full-refresh --profiles-dir .
9
10 # 4. Run the SCD Type 2 snapshot against stg_playstore_apps
```

```
11   dbt snapshot --profiles-dir .
12
13   # 5. Build the clean SCD dimension model from the snapshot
14   dbt run --select dim_apps_scd --profiles-dir .
15
16   # 6. Export all mart tables to CSV files for Power BI
17   cd ..
18   python export_to_powerbi.py
```

# Advanced Features

## 5.1 Slowly Changing Dimensions — SCD Type 2

One of the more advanced requirements of the lab was implementing SCD Type 2 on the app dimension. The motivation is straightforward: if an app's category changes over time, any fact record linked only to the current dimension row would retroactively misattribute historical reviews to the wrong category. SCD Type 2 solves this by preserving every version of the row with validity timestamps, so a review written when WhatsApp was categorised as "Social Networking" can still be correctly linked to that historical version even after the category changes to "Communication."

We implemented this using dbt's native **snapshot** feature with the *check strategy*, where dbt compares a defined set of columns on each run and inserts a new row whenever a change is detected, automatically managing `dbt_valid_from` and `dbt_valid_to`:

**Listing 5.1:** *Snapshot definition implementing SCD Type 2 on `dim_apps`*

```
1  {% snapshot dim_apps_snapshot %}
2  {{
3    config(
4      target_schema='main',
5      unique_key='app_id',
6      strategy='check',
7      check_cols=['category_id', 'avg_score', 'total_ratings'],
8      invalidate_hard_deletes=True
9    )
10 }}
11 SELECT * FROM {{ ref('stg_playstore_apps') }}
12 {% endsnapshot %}
```

After running `dbt snapshot`, a downstream model `dim_apps_scd` reads from the snapshot table and adds an `is_current` flag (true where `dbt_valid_to` is null) for convenient current-record filtering in BI tools.

**Figure 5.1:** `dim_apps_scd`: *WhatsApp Messenger appears twice with different category values, confirming SCD Type 2 captured the change from "Social Networking" to "Communication"*

Figure 5.1 makes the SCD behaviour visible. WhatsApp Messenger appears in two rows with different `category_name` values—one historical (closed) record and one current (open) record. The `dbt_valid_from` and `dbt_valid_to` columns are present and correctly populated, as the guidelines require.

## 5.2   Incremental Loading

Beyond SCD, the guidelines required transforming the pipeline to support incremental ingestion. Rather than rebuilding the entire fact table on every run, the updated `fact_reviews` model uses dbt's `incremental` materialisation strategy, filtering new records on `review_at` and deduplicating on `review_id`. On a first run the table is built in full; on subsequent runs only reviews newer than the maximum already-stored `review_at` are processed and appended. This represents a significant efficiency gain as the dataset grows: rebuilding 40 000-plus rows from scratch becomes unnecessary once the pipeline enters production.

# Data Quality and Testing

Data quality was enforced at every layer using dbt's native schema tests. Tests were defined in `schema.yml` files placed alongside each model and executed with `dbt test` after each `dbt run`. The table below documents every test applied across the pipeline:

Table 6.1: *dbt schema tests applied across all pipeline layers*

| Layer | Model | Test applied | Column |
|---|---|---|---|
| Staging | stg_playstore_apps | not_null, unique | app_sk |
| Staging | stg_playstore_apps | not_null | app_id, app_name |
| Staging | stg_playstore_reviews | not_null, unique | review_sk |
| Staging | stg_playstore_reviews | not_null | review_id, app_id |
| Staging | stg_playstore_reviews | accepted_values (1–5) | review_score |
| Dimension | dim_apps | not_null, unique | app_sk |
| Dimension | dim_categories | not_null, unique | category_sk |
| Dimension | dim_date | not_null, unique | date_key |
| Dimension | dim_apps | relationships | category_sk, developer_sk |
| Fact | fact_reviews | not_null, unique | review_sk |
| Fact | fact_reviews | relationships | app_sk, date_key |
| Fact | fact_reviews | not_null | review_score, thumbs_up_count |

All tests pass on the final pipeline run, confirming referential integrity across the snowflake schema and the absence of null primary keys or out-of-range measure values. The `accepted_values` test on `review_score` is particularly important: it enforces the constraint that scores fall in the 1–5 range at the staging layer, preventing any corrupted API response from polluting the fact table with invalid data.

# Results and Visualisations

## 7.1 Summary of Materialised Data

After a complete pipeline run, the following table summarises everything built inside DuckDB and exported to Power BI:

**Table 7.1:** *Mart tables produced after a full pipeline run*

| Table | Type | Rows | Description |
|---|---|---|---|
| dim_apps | Dimension | 20 | One record per scraped application |
| dim_categories | Dimension | 12 | Unique Play Store category labels |
| dim_developers | Dimension | 17 | Unique developer entities |
| dim_date | Dimension | ∼100 | Calendar days covering the review period |
| dim_apps_scd | SCD Dimension | 21+ | Current and historical app versions |
| fact_reviews | Fact | 40 000+ | Individual review events at declared grain |

The fact table holds over 40 000 rows representing individual reviews posted within the last 100 days across the 20 monitored applications. WhatsApp Messenger alone—with 225 M total ratings on the Play Store—contributes the largest share of recent reviews to the dataset.

## 7.2 Power BI Dashboard

**Figure 7.1:** *Power BI dashboard: four analytical views over the Google Play Store review dataset*

The dashboard (Figure 7.1) presents four complementary analytical views.

**Average Review Score by App (bar chart, top-left).** Duolingo: Language Lessons leads the ranking with the highest average score among all 20 applications, followed by WhatsApp Messenger and Facebook. Discord and Amazon Shopping receive notably lower ratings. Netflix and YouTube, despite their massive install bases, hover around 3.9—consistent with the mixed reception that subscription-gated or ad-heavy platforms tend to attract from their most engaged users.

**Sum of Thumbs-Up Count by Category (pie chart, top-right).** The Social and Communication categories together account for the largest shares of community engagement on reviews. WhatsApp, Instagram, Facebook, Snapchat, and Discord all fall within these two categories and generate the highest absolute volumes of recent reviews. A meaningful contribution from the Business category (Microsoft Teams, LinkedIn) suggests that productivity users engage actively with each other's feedback.

**Review Count Over Time (line chart, bottom-left).** The time series covers November 2025 through February 2026 and reveals a sharp spike in review activity in mid-February 2026. This kind of signal—immediately obvious in a properly modelled time-series chart—would be nearly invisible in a flat CSV export. The date dimension makes

21

it trivial to filter by app and drill into the spike to identify its cause.

**Average Review Score vs. Total Ratings (scatter chart, bottom-right).** Applications with very large rating counts (above 100 M: YouTube, WhatsApp, Instagram) do not necessarily achieve the highest average scores. Duolingo and Microsoft Teams, with comparatively modest rating counts, reach averages above 4.5—suggesting that niche and productivity-focused apps attract more considered and generous reviews than mass-market entertainment platforms.

## 7.3   Power BI Relationship Configuration

After importing the CSVs, four relationships were configured in Power BI's model view to reproduce the snowflake schema:

**Table 7.2:** *Power BI relationships matching the snowflake schema design*

| From table / column | To table / column | Cardinality |
|---|---|---|
| dim_apps.app_sk | fact_reviews.app_sk | 1 : many |
| dim_date.date_key | fact_reviews.date_key | 1 : many |
| dim_categories.category_sk | dim_apps.category_sk | 1 : many |
| dim_developers.developer_sk | dim_apps.developer_sk | 1 : many |

# Guideline Compliance Summary

The table below cross-references every requirement from the lab guidelines with its implementation evidence:

**Table 8.1:** *Compliance matrix against lab guidelines*

| Guideline Requirement | Status | Evidence |
|---|---|---|
| Install DuckDB, dbt-core, dbt-duckdb adapter | ✓ Done | `requirements.txt` |
| Use Python for data ingestion (from Lab 1) | ✓ Done | `ingest.py` |
| Use JSONL as immutable raw landing files | ✓ Done | `data/raw/apps.jsonl`, `reviews.jsonl` |
| Initialise dbt project with DuckDB adapter | ✓ Done | `dbt_project.yml`, `profiles.yml` |
| Enforce staging as views, marts as tables | ✓ Done | `dbt_project.yml` materialisation config |
| Staging: rename, cast types, surrogate keys, null filter | ✓ Done | `stg_playstore_apps.sql`, `stg_playstore_reviews.sql` |
| Apply Kimball 4-step process and Bus Matrix | ✓ Done | Chapter 3 of this report |
| Build `dim_apps`, `dim_categories`, `dim_developers` | ✓ Done | `models/marts/dimensions/` |
| Build conformed `dim_date` (YYYYMMDD integer key) | ✓ Done | `dim_date.sql`, Figure 4.3 |
| Build `fact_reviews` at declared grain | ✓ Done | `fact_reviews.sql`, Figure 4.4 |
| Apply dbt schema tests (not_null, unique, relationships) | ✓ Done | `schema.yml` files, Table 6.1 |
| Serve analytics-ready data to a BI tool | ✓ Done | `export_to_powerbi.py`, Figure 7.1 |

*Continued. . .*

*(Table 8.1 continued)*

| Guideline Requirement | Status | Evidence |
| --- | --- | --- |
| Map foreign keys correctly in BI model view | ✓ Done | Figure 3.1, Table 7.2 |
| Implement SCD Type 2 with dbt snapshots | ✓ Done | `snapshots/`, `dim_apps_scd.sql`, Figure 5.1 |
| Implement incremental loading on `fact_reviews` | ✓ Done | Incremental materialisation, Chapter 5 |
| Create data visualisations (at least 3 charts) | ✓ Done | Four charts, Figure 7.1 |

*(Table 8.1 continued)*