# AG-PST DevOps-Driven Data Ingestion and Analysis

# Objective:

The objective of this assignment is to **design and implement a containerized data ingestion and analysis pipeline**. The solution demonstrates how to:

1. Select a complex dataset that maps naturally into a relational database schema.
2. Build a **Dockerized environment** with a database and ingestion service.
3. Ingest datasets into normalized tables using **Python and Pandas**.
4. Run **meaningful queries** (aggregations, joins, inserts) on demand.
5. Showcase **DevOps best practices** such as modularity, environment isolation, reproducibility, and idempotency.

# Dataset Selection:

We used the **TMDB 5000 Movies dataset** ( public dataset from Kaggle )
Url: https://www.kaggle.com/datasets/tmdb/tmdb-movie-metadata?select=tmdb_5000_credits.csv

This dataset is good, relational, and widely recognized. It was split into three CSVs for ingestion:
- **tmdb_5000_movies1.csv** renamed as **movie_metadata**
  Columns in this csv: id, budget, homepage, original_language, original_title
- **tmdb_5000_movies2.csv** renamed as **movie_details**
  Columns in this csv: id, popularity, release_date, revenue, runtime, status, vote_average, vote_count
- **tmdb_5000_movies3.csv** renamed as **movie_castcrew**
  Columns in this csv: movie_id, title, tagline

### Reason behind selecting this dataset:

Fits **multi-table relational modeling** (metadata, details, castcrew) The **TMDB 5000 Movies dataset** was chosen because it naturally lends itself to a **relational database model** while also being rich enough to demonstrate real-world data engineering practices. The dataset contains multiple aspects of movies such as metadata (budget, language, title), details (revenue, runtime, votes, ratings), and cast/crew information that can be cleanly separated into different tables and connected through foreign keys. This ensures that the project is not limited to flat file ingestion but instead highlights **database normalization, relationships, and joins**.Additionally, this dataset supports **meaningful analytical queries** such as finding the top movies by revenue, aggregating ratings with vote thresholds, or joining movie details with cast/crew information.

### Database Choice:

We chose **PostgreSQL 15** because:

- **Open-source, stable, production-grade** relational database widely used in industry
- Excellent support for **normalized schemas** with primary/foreign keys and constraints
- Strong **Python integration** via SQLAlchemy + psycopg2, making ingestion seamless
- Official lightweight **Docker image available**, making containerization simple and reliable
- Provides **robust query optimization** and indexing for analytical queries
- Ensures **data integrity and scalability**, suitable for real-world ingestion pipeline

Alternative DBs like MySQL or DuckDB could also work, but PostgreSQL offers the most **robust relational capabilities** for this project.

# Repository Setup:

```
├── data/
│   ├── tmdb_5000_movies1.csv      # Movie metadata (budget, language, title, etc.)
│   ├── tmdb_5000_movies2.csv      # Movie details (revenue, runtime, ratings, etc.)
│   └── tmdb_5000_movies3.csv      # Movie cast & crew info
│
├── ingest/
│   ├── Dockerfile                 # Python image setup for ingestion service
│   ├── main.py                    # Main ingestion & analysis logic
│   ├── requirements.txt           # Python dependencies
│
├── .env                           # Environment variables (DB credentials)
├── docker-compose.yml             # Defines database + ingestion services
└── README.md                      # Project documentation
```

The repository is structured to separate concerns clearly. Each folder and file serves a distinct purpose, making the project easy to understand, maintain, and to extend.

- **data/** → Contains the raw datasets (tmdb_5000_movies1.csv, tmdb_5000_movies2.csv, tmdb_5000_movies3.csv) used for ingestion. Keeping them in a dedicated folder ensures that data is organized and version-controlled separately from code.
- **ingest/** → Holds all the ingestion logic and supporting files:
  - **Dockerfile** → Builds a lightweight Python container for ingestion.
  - **main.py** → The core ingestion script that connects to PostgreSQL, loads CSVs into tables, and runs queries.
  - **requirements.txt** → Lists dependencies (pandas, sqlalchemy, psycopg2-binary) to guarantee reproducibility.
- **.env** → Stores sensitive configuration values (DB user, password, and database name). Using environment variables follows 12-factor app principles, keeps secrets out of code, and allows different environments (dev, test, prod) without code changes.

- **docker-compose.yml** → Defines and orchestrates multiple services:
  - postgres_db (database service)
  - ingest_service (Python ingestion service)
    Both run inside a shared Docker network, ensuring secure communication.
- **README.md** → Comprehensive documentation with setup instructions, project explanation, screenshots, and usage details.

# Dockerized Solution:

We used **Docker Compose** with two services mainly db (PostgreSQL) and ingest_and_query (Python)

## Service 1: db (PostgreSQL)
- Based on postgres 15 image.
- Credentials + DB name injected via .env.
- Exposes port 5432.
- Persists data in Docker volume db_data.
- Uses Docker volume (db_data) for persistence

## Service 2: ingest_and_query (Python)
- Built with python:3.10-slim.
- Installs required libraries ( pandas, sqlalchemy, psycopg2 )
- Runs main.py.
  - Waits for DB readiness
  - Loads CSV data into tables (movie_metadata, movie_details, movie_castcrew)
  - Runs queries (aggregations, join, insert)
- Runs inside the same Docker network as postgres_db
- Executes sample SQL queries.

Both services run on the same **Docker network** for secure communication.

**Docker Compose Build & Startup**

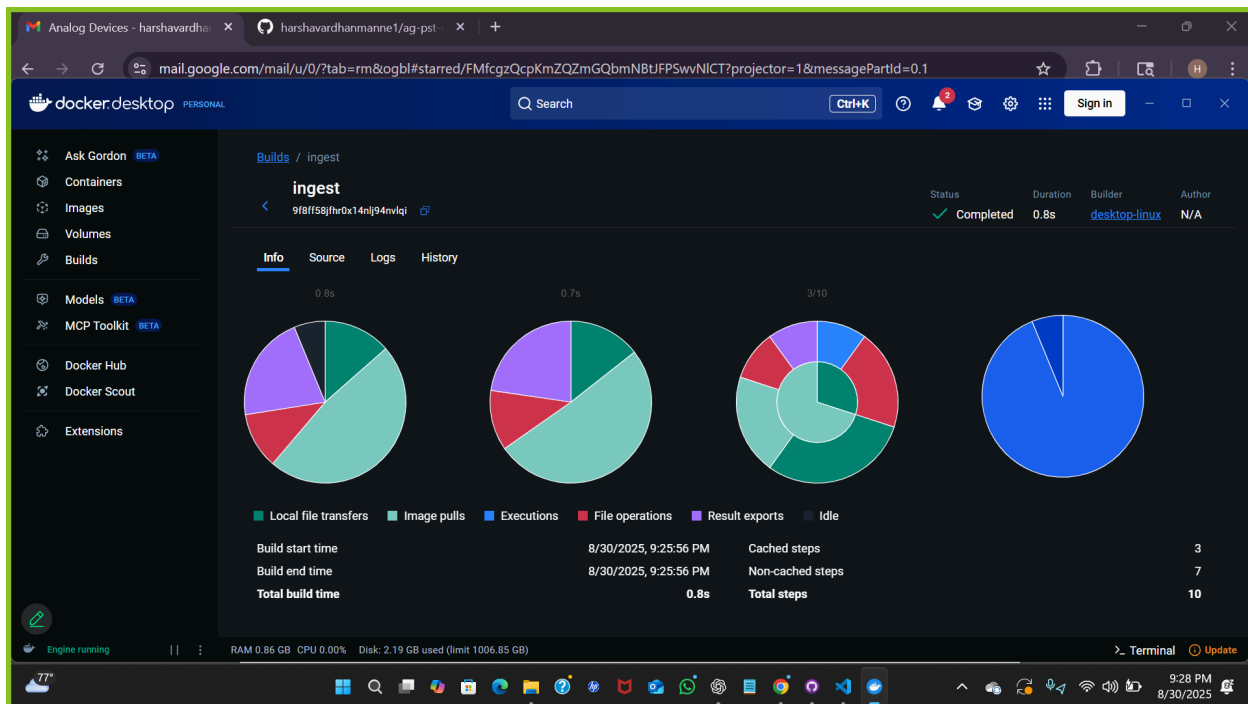This screenshot shows the execution of docker-compose up --build, where:

- The ingest_and_query service image is built successfully.
- A dedicated Docker network (ag-pst-devops-driven-data-ingestion_app-network) is created for secure communication between services.
- A persistent volume (ag-pst-devops-driven-data-ingestion_db_data) is created to ensure PostgreSQL data survives container restarts.
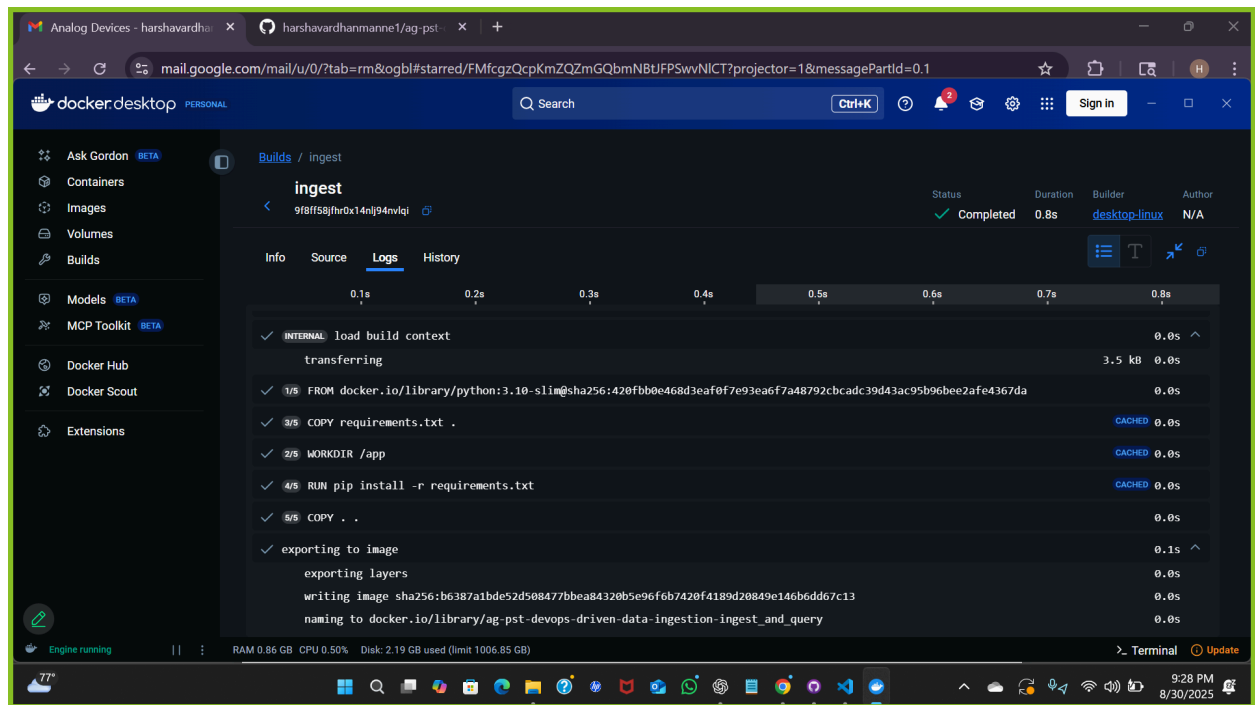- Both containers (postgres_db and ingest_service) are started successfully.

**Docker Image Build**

This screenshot shows the custom Python ingestion image successfully built using the ingest/Dockerfile. It validates that dependencies (pandas, sqlalchemy, psycopg2) were installed and the container is ready to run ingestion logic.

**Docker Build Logs:** These logs show the step-by-step build process for the Python service.



This screenshot shows build stats & timing of Docker layers.

This screenshot shows detailed installation logs, confirming all dependencies were installed correctly.

## Runtime Phase

**Running Containers**



This screenshot from Docker Desktop shows both services running:

- postgres_db → PostgreSQL container
- ingest_service → Python ingestion service

It proves that docker-compose.yml correctly orchestrated the services.

**Docker Volumes**

This screenshot shows the Docker volume db_data created for PostgreSQL persistence. It ensures that even if the container restarts, the data is not lost.

# Database Schema:

```
CREATE TABLE movie_metadata (
    id INT PRIMARY KEY,
    budget BIGINT,
    homepage VARCHAR,
    original_language VARCHAR,
    original_title VARCHAR
);

CREATE TABLE movie_details (
    id INT PRIMARY KEY,
    popularity FLOAT,
    release_date DATE,
    revenue BIGINT,
    runtime INT,
    status VARCHAR,
    vote_average FLOAT,
    vote_count INT,
    FOREIGN KEY (id) REFERENCES movie_metadata(id)
);
```

```
CREATE TABLE movie_castcrew (
    movie_id INT,
    title VARCHAR,
    tagline TEXT,
    FOREIGN KEY (movie_id) REFERENCES movie_metadata(id)
);
```

The schema was designed to be **normalized, relational, and extensible**,

- Different aspects of the dataset (metadata, details, cast/crew) are stored in separate tables, reducing redundancy.
- Foreign key relationships (e.g., movie_details.id → movie_metadata.id) ensure consistency between tables.
- New attributes or additional datasets can be added without disrupting existing tables now or later.
- Splitting into multiple tables allows targeted queries (ratings, revenue, or tagline lookups) without scanning unnecessary data.
- Normalization prevents data duplication, making storage efficient and queries faster as the dataset grows.
- Mimics how relational databases are used in production systems for movie catalogs, analytical dashboards, or recommendation engines

**PostgreSQL Logs**



This screenshot shows logs from the PostgreSQL container, confirming successful initialization and that

the database is ready to accept connections. It validates the schema and DB service are running as expected.

# Ingestion Workflow ([main.py](main.py)):
## Steps performed:

1. **Wait for Database Readiness**
   - The ingestion service waits until the postgres_db container is fully up.
   - Prevents connection failures by retrying until PostgreSQL accepts connections.
2. **Read CSV Files with Pandas**
   - Loads each dataset ( tmdb_5000_movies1.csv, tmdb_5000_movies2.csv, tmdb_5000_movies3.csv ) into Pandas DataFrames.
   - Ensures the data is clean and structured before loading.
3. **Write to PostgreSQL using SQLAlchemy**
   - Uses df.to_sql() to insert data into corresponding tables (movie_metadata, movie_details, movie_castcrew).
   - Ensures type mapping and efficient bulk insert.
   - Uses append/replace mode to avoid duplicate schema creation.
4. **Run SQL Queries via SQLAlchemy**
   - Executes predefined SQL queries for analysis.
   - Results are printed in container logs for verification.
5. **Insert a Test Movie**
   - Demonstrates that the pipeline supports inserts in addition to reads.
   - Ensures idempotency by using ON CONFLICT DO NOTHING.
6. **Print results**
   - visible in Docker logs

# Queries Implemented:

The project includes **four categories of queries** to demonstrate practical use cases:

- Top 5 movies by revenue
- Top 5 movies by average rating (votes > 1000)
- Insert a test movie record
- Join metadata + details + castcrew

**Aggregation Queries:**

*Top 5 Movies by Revenue*

```sql
SELECT m.original_title, d.revenue
FROM movie_details d
JOIN movie_metadata m ON d.id = m.id
ORDER BY d.revenue DESC
LIMIT 5;
```

*Top 5 Movies by Average Rating (votes > 1000)*

```sql
SELECT m.original_title, d.vote_average, d.vote_count
FROM movie_details d
JOIN movie_metadata m ON d.id = m.id
WHERE d.vote_count > 1000
ORDER BY d.vote_average DESC
LIMIT 5;
```

**Insert Query**

*Add a Test Movie into Metadata*

```sql
INSERT INTO movie_metadata (id, budget, homepage, original_language,
VALUES (9999999, 100000, 'http://test.com', 'en', 'Test Movie')
ON CONFLICT (id) DO NOTHING;
```

**Join Query**

*Combine Metadata, Details, and CastCrew*

```python
print("\n Example 3-table join (metadata + details + castcrew):")
result = conn.execute(text("""
    SELECT m.original_title, d.release_date, d.vote_average, c.tagline
    FROM movie_metadata m
    JOIN movie_details d ON m.id = d.id
    JOIN movie_castcrew c ON m.id = c.movie_id
    LIMIT 5;
"""))
```

# Query Outputs :

This screenshot shows the output of the query **Top 5 movies by revenue**. It validates that the dataset was ingested correctly and queries are running successfully inside the container.

**Query Outputs (Rating, Insert, Join)**



This screenshot shows the results of multiple queries average rating, inserting test movie query and Join query (metadata + details + castcrew).