

Airline Business Intelligence Database

End-to-End BI System Using SQL, ETL, PostgreSQL & Python

Grace Polito

Final Capstone Project - MSDS DTSC 691
Eastern University

PRESENTATION OVERVIEW

- 1 Intro & Problem Domain
- 2 Project Timeline & Phase Overview
- 3 Schema Design & Data Model
- 4 Data Sources & ETL Pipeline
- 5 SQL Analytics Layer
- 6 Python Integration & Visual Analytics
- 7 Key Business Insights
- 8 Challenges, Limitations & Lessons Learned
- 9 Future Enhancements
- 10 Conclusion & Next Steps



Introduction & Problem Domain

//

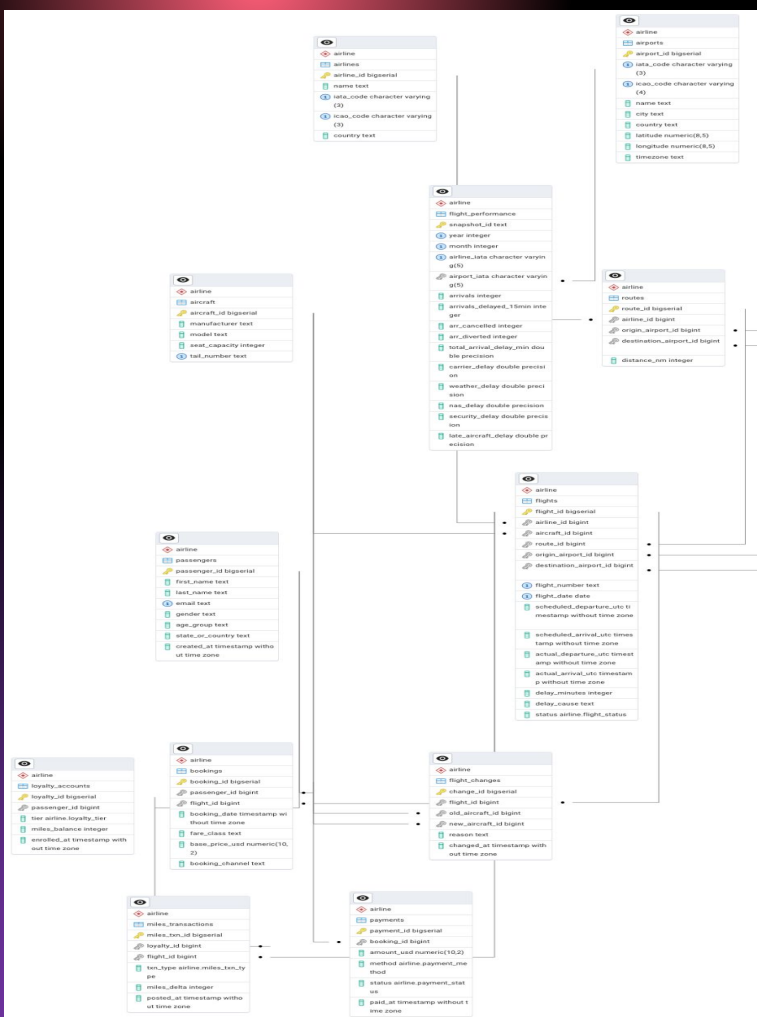
The core problem this project tackles is the **fragmentation** of airline **operational** and **commercial** data.

- Airlines generate complex operational, customer, and revenue data
- Spreadsheets & disconnected sources limit analytics quality
- Need for a unified, query-ready analytical database
- Project goal: build a BI-ready airline database + analytics layer







Project Timeline (Phase 1-6)













01	Design & Setup	Schema, ERD, environment, initial constraints
02	Data Collection & Insertion	OpenFlights/BTS import, synthetic generation
03	SQL CLeaning & Constraints	DML standardization, deduplication, indexing
04	Query Development	15+ analytical SQL queries + performance testing
05	Python Analytics	Engine connection, helper functions, charts
06	Final Deliverables	Overview PDF, exported code, 20-minute presentation



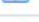





Schema Design & Data Model Overview




























Core Operational Entities


 airline
 airlines
 airline_id bigserial
 name text
 iata_code character varying (3)
 icao_code character varying (3)
 country text













 airline
 airports
 airport_id bigserial
 iata_code character varying (3)
 icao_code character varying (4)
 name text
 city text
 country text
 latitude numeric(8,5)
 longitude numeric(8,5)
 timezone text












 airline
 routes
 route_id bigserial
 airline_id bigint
 origin_airport_id bigint
 destination_airport_id bigint
 distance_nm integer











 airline
 aircraft
 aircraft_id bigserial
 manufacturer text
 model text
 seat_capacity integer
 1 tail_number text











 airline
 flights
 flight_id bigserial
 airline_id bigint
 aircraft_id bigint
 route_id bigint
 origin_airport_id bigint
 destination_airport_id bigint
 1 flight_number text
 1 flight_date date
 scheduled_departure_utc timestamp without time zone
 scheduled_arrival_utc timestamp without time zone
 actual_departure_utc timestamp without time zone
 actual_arrival_utc timestamp without time zone
 delay_minutes integer
 delay_cause text
status airline.flight_status









Commercial Entities: Passengers, Bookings & Payments


 airline
 passengers
 passenger_id bigserial
 first_name text
 last_name text
 1 email text
 gender text
 age_group text
 state_or_country text
 created_at timestamp without time zone


 airline
 bookings
 booking_id bigserial
 passenger_id bigint
 flight_id bigint
 booking_date timestamp without time zone
 fare_class text
 base_price_usd numeric(10, 2)
 booking_channel text


 airline
 payments
 payment_id bigserial
 booking_id bigint
 amount_usd numeric(10,2)
 method airline.payment_method
 status airline.payment_status
 paid_at timestamp without time zone


 airline
 miles_transactions
 miles_txn_id bigserial
 loyalty_id bigint
 flight_id bigint
 txn_type airline.miles_transaction_type
 miles_delta integer
 posted_at timestamp without time zone


 airline
 loyalty_accounts
 loyalty_id bigserial
 passenger_id bigint
 tier airline.loyalty_tier
 miles_balance integer
 enrolled_at timestamp without time zone

Analytical Fact Tables: Performance & Change Tracking

airline
flight_changes
change_id bigserial
flight_id bigint
old_aircraft_id bigint
new_aircraft_id bigint
reason text
changed_at timestamp with out time zone

airline
flight_performance
snapshot_id text
year integer
month integer
airline_iata character varyin g(5)
airport_iata character varyin g(5)
arrivals integer
arrivals_delayed_15min inte ger
arr_cancelled integer
arr_diverted integer
total_arrival_delay_min dou ble precision
carrier_delay double preci on
weather_delay double preci sion
nas_delay double precision
security_delay double precis ion
late_aircraft_delay double pr ecision

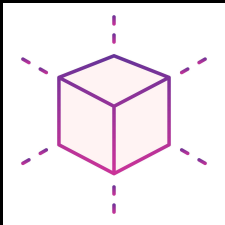
airline
miles_transactions
miles_txn_id bigserial
loyalty_id bigint
flight_id bigint
txn_type airline.miles_txn_ty pe
miles_delta integer
posted_at timestamp witho ut time zone

WHY 3NF?

Schema Design Principles

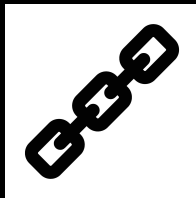
ATOMIC TABLES

- Each table models one concept
- Eliminates duplicate fields



REFERENTIAL INTEGRITY

- Strong PK/FK relationships
- Prevents orphaned or inconsistent records



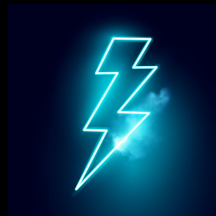
CONTROLLED BUSINESS RULES

- ENUMs for statuses & tiers
- CHECK constraints protect data quality



PERFORMANCE OPTIMIZATION

- Indexed join keys
- Supports fast analytical queries



DATA SOURCES

OpenFlights real-world operational data

provides **real-world** reference data for airports and airlines

- Airports & airlines reference tables
- Global identifiers (IATA/ICAO)
- Used for route + flight generation
- Real Data

BTS On-Time Performance real-world delay data

contains monthly airline **delay** metrics

- Monthly delay performance metrics
- Carrier, weather, security, NAS delays
- Used to enrich operational analysis
- Real Data

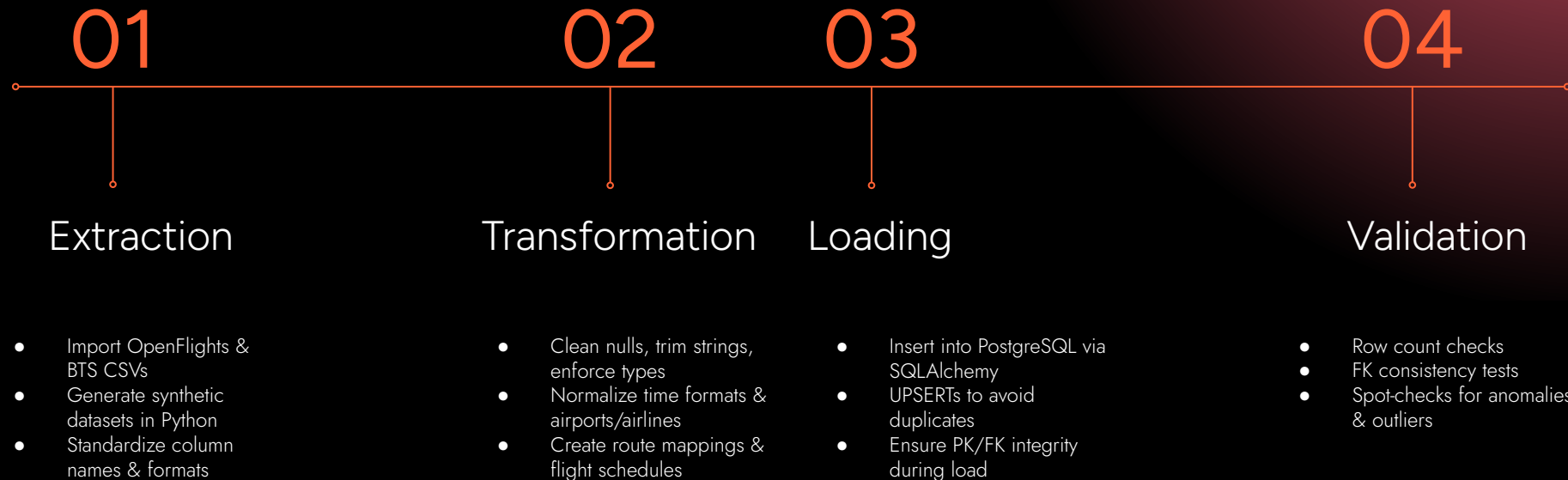
Synthetic Data flights, passengers, bookings, payments, loyalty

generated to complete the **commercial side** of the model

- Flights, passengers, bookings, payments
- Loyalty accounts & miles transactions
- Enables complete operational + commercial coverage
- Generated in Python with faker

ETL PIPELINE

The ETL pipeline follows a structured Extract—Transform—Load workflow.



SQL ANALYTICS LAYER

15 analytical SQL queries built for operations, revenue, network, and loyalty insights.

Window Functions

Ranking airlines by average delay

```
SELECT airline_name,  
       iata_code,  
       AVG(delay_minutes) AS  
       avg_delay_minutes,  
       DENSE_RANK() OVER (ORDER BY  
       AVG(delay_minutes) DESC) AS  
       delay_rank  
FROM flights  
GROUP BY airline_name, iata_code;
```

CTEs & Aggregations

Multi-hop airport connectivity

```
WITH RECURSIVE connections AS (  
    SELECT origin_iata, dest_iata, 1 AS hops,  
    ARRAY[origin_iata, dest_iata] AS path  
    FROM routes  
    WHERE origin_iata = 'YCK'  
    UNION ALL  
    SELECT c.origin_iata, r.dest_iata, c.hops + 1,  
    path || r.dest_iata  
    FROM connections c  
    JOIN routes r ON c.dest_iata = r.origin_iata  
    WHERE c.hops < 3)  
SELECT * FROM connections;
```

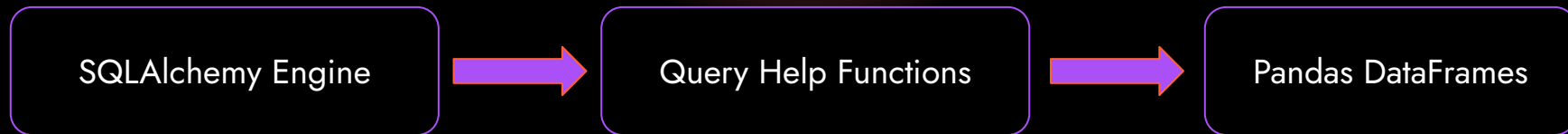
Advanced Analytical Queries

**Revenue, fare class, payment success,
CLV-style revenue**

```
SELECT fare_class,  
       COUNT(*) AS num_bookings,  
       SUM(amount_usd) AS total_revenue  
FROM bookings  
JOIN payments USING (booking_id)  
GROUP BY fare_class  
ORDER BY total_revenue DESC;
```

Python Integration & Analytics Layer

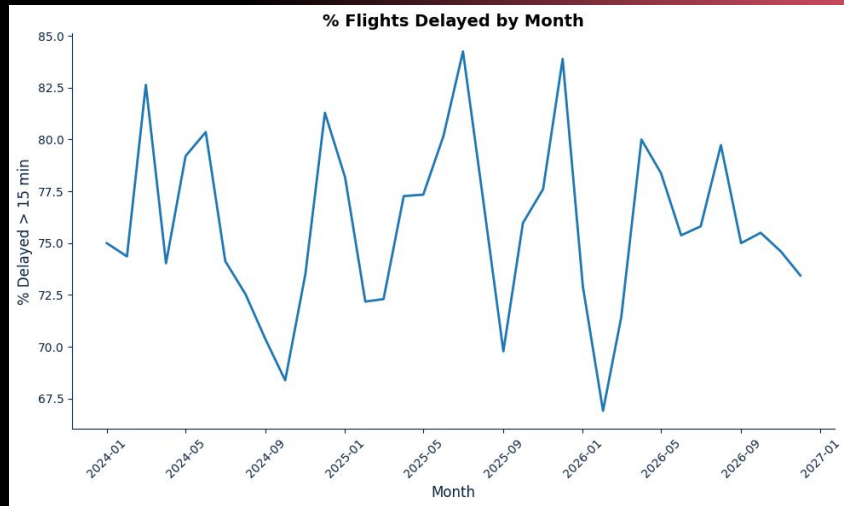
Python, SQLAlchemy, Pandas, Matplotlib, and Plotly power the BI analysis.



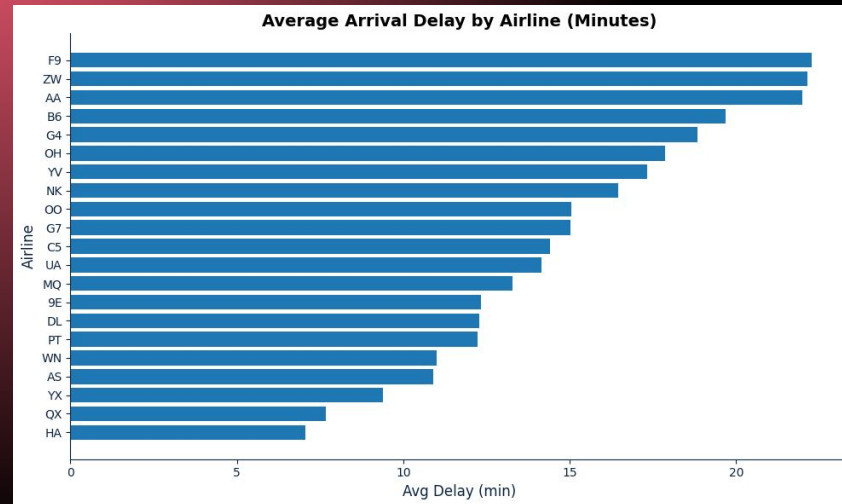
```
def get_df(sql: str, params=None):  
    engine = get_engine()  
    with engine.connect() as conn:  
        return pd.read_sql(text(sql), conn, params=params)
```

- Connects Python ↔ PostgreSQL
- Reusable query wrappers (e.g., `get_revenue_by_fare_class`)
- Supports operational, revenue, and loyalty analytics
- Drives visualizations (Matplotlib + Plotly)

Operational Performance Visuals

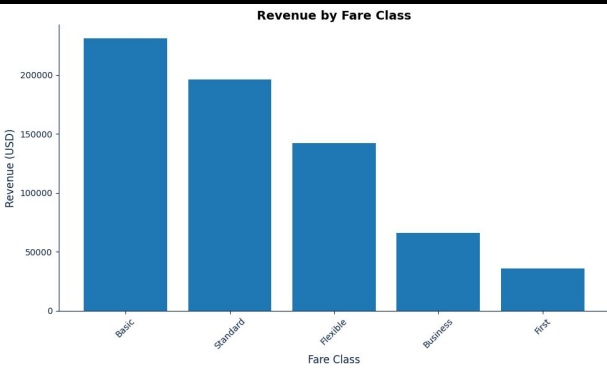


Delay rates fluctuate **seasonally**, with spring and early winter showing the **highest disruption** levels. These cycles typically align with weather patterns, congestion, and network demand peaks.



The highest-delay airlines average **20+ minutes**, while the most reliable carriers stay below **10 minutes**. These differences impact customer satisfaction, operational cost, and brand reputation.

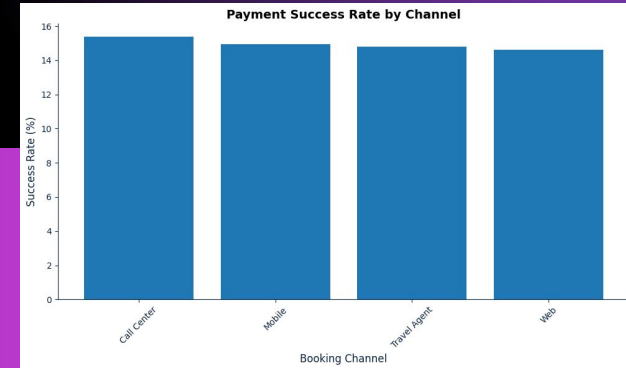
Commercial Performance: Revenue, Payments & CLV



Basic and Standard fares drive the majority of revenue volume, reflecting price-sensitive demand in the synthetic dataset.



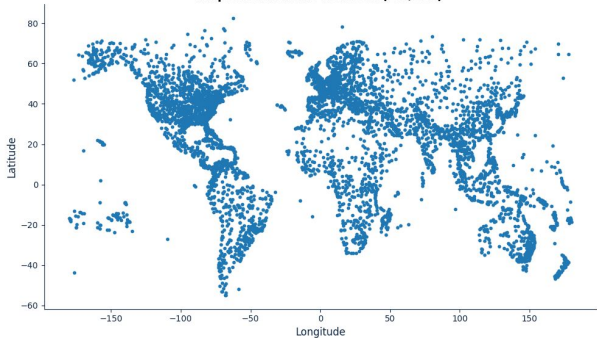
Revenue trends illustrate cyclical demand, confirming peak travel periods and slower off-season months.



All channels exhibit similar success rates, with Call Center slightly outperforming digital channels.

Network & Geographic Visualizations

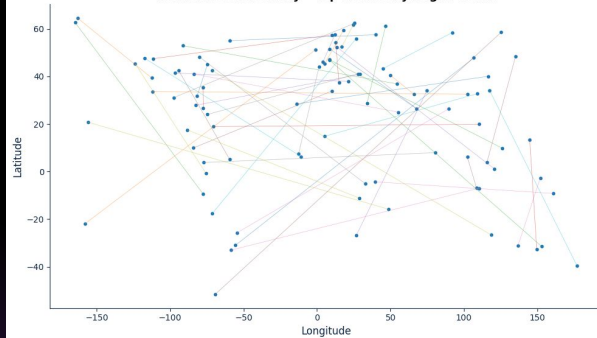
Airports in Airline Network (Lat/Lon)



Busiest Origin-Destination Pairs (Flights)



Network Connectivity - Top Routes by Flight Count



Airport coordinates plotted globally to reveal the geographic footprint of the modeled network.

High-volume OD pairs—such as ALN→BFI and BUR→FET—highlight the strongest traffic flows in the network.

Most-frequent routes drawn as direct coordinate links, illustrating the core structure of the airline's route system.

KEY BUSINESS INSIGHTS

Operational Insights

- ★ Delays show clear **seasonal** patterns
- ★ Airline **reliability** varies significantly

Revenue Insights

- ★ **Basic** and **Standard** fares drive most revenue volume
- ★ Revenue exhibits cyclical trends

Payment & Booking Insights

- ★ Payment success rates are **consistent** across channels
- ★ Synthetic booking patterns still **reflect realistic** conversion behavior

Customer & Loyalty Insights

- ★ Top **5%** of customers generate **~13%** of total revenue
- ★ Loyalty tiers show balanced distribution

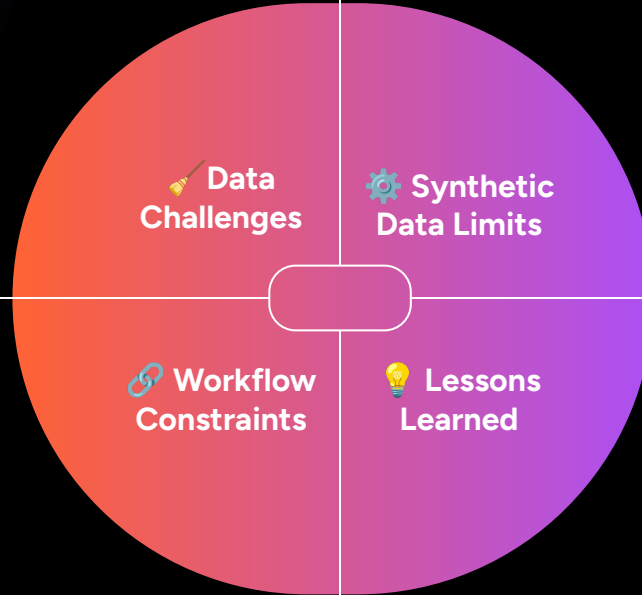
Network Insights

- ★ Top **OD pairs** form key network corridors
- ★ Connectivity highlights major route clusters

CHALLENGES, LIMITATIONS & LESSONS LEARNED

- Required heavy **cleaning** & **standardization**
- Multiple sources needed ID/time alignment

- PK/FK integrity required **iterative cleanup**
- Recursive CTEs & large joins needed tuning



- Some airports unmapped → limited geography
- Future dates used to model **full-year trends**

- Clean schema drives **better analytics**
- Reusable pipelines reduce rework

FUTURE ENHANCEMENTS



Data Expansion

Integrate **real airline** schedules, fares, and ancillary data for **richer modeling**



Advanced Analytics

Develop CLV forecasting, demand prediction, and **anomaly detection** models



Enhanced Visualizations

Build **interactive dashboards** and basemap-backed route maps using **Tableau**



Production Readiness

Containerize the pipeline and add **automated testing** for reproducible workflows

Conclusion & Next Steps

Project Summary



Built a unified, BI-ready airline database integrating operational, booking, loyalty, and payment data.

Key Outcomes

Delivered end-to-end analytics uncovering operational reliability, revenue drivers, and customer value.

Next Steps

Expand with real schedules, predictive models, and interactive dashboards.

THANK YOU

Grace Polito

Airline Business Intelligence Database