

# Data Pipeline Submission IE7374

Charlie - The MBTA Conversational Transit Assistant

Group - 22

# 1. Executive Summary

The **Charlie – MBTA Conversational Transit Assistant** project aims to revolutionize the way commuters interact with the Massachusetts Bay Transportation Authority (MBTA) data by developing a real-time, conversational AI platform that consolidates transit information from multiple sources into a single, user-friendly interface.

This initiative was developed as part of the **IE7374 MLOps course**, emphasizing the design and automation of an end-to-end data pipeline using **Apache Airflow**, **Docker**, **Kubernetes**, and **FastAPI**, integrated with MBTA’s open APIs and real-time GTFS feeds

## Project Objective

The project’s primary goal is to simplify the MBTA commuter experience by providing **accurate, real-time answers** to natural language questions such as:

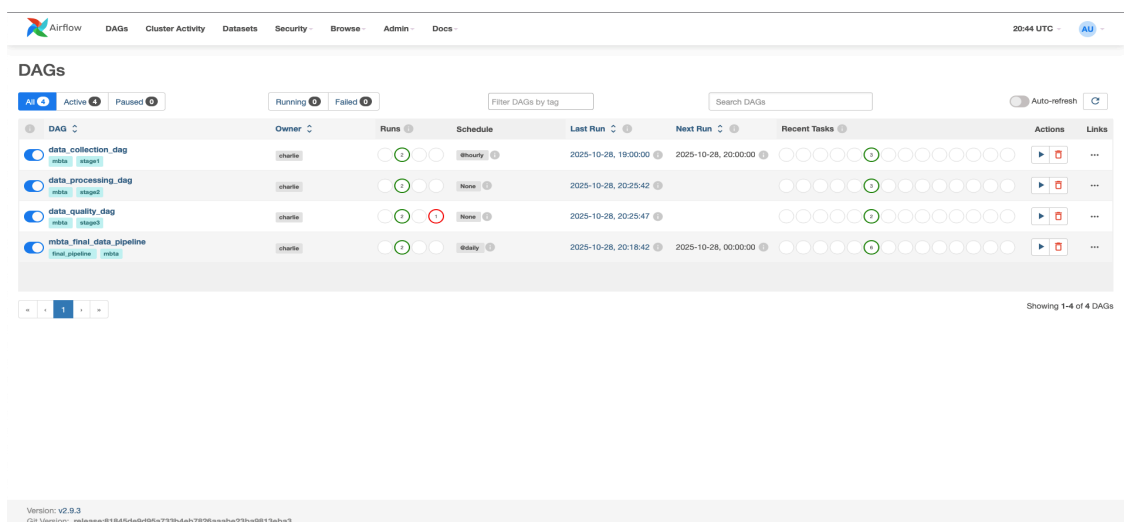
- “When is the next Red Line train to Alewife?”
- “Why is the Green Line delayed?”
- “Where is my bus right now?”

By fusing real-time arrival predictions, vehicle positions, service alerts, and fare details, the chatbot delivers **one unified response** that combines *what*, *why*, and *when* information into a single, contextually aware answer. This directly addresses gaps in existing transit applications that require users to consult multiple interfaces for schedules, alerts, and delays.

## Pipeline Overview

To operationalize this system, a **multi-stage MLOps pipeline** was architected and automated through **Apache Airflow**, ensuring reliability, observability, and reproducibility. The pipeline consists of four modular DAGs:

1. **data\_collection\_dag** – Fetches live data from MBTA APIs and stores raw snapshots for traceability.
2. **data\_processing\_dag** – Normalizes, cleans, and transforms the ingested JSON and Protocol Buffer feeds.
3. **data\_quality\_dag** – Runs schema validation, freshness checks, and contract tests to detect anomalies or API drifts.
4. **mbta\_final\_data\_pipeline** – Integrates all previous stages to serve structured, validated, and cached data for the chatbot’s back-end engine.



The screenshot displays the Apache Airflow web interface. At the top, there's a navigation bar with links for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. The main header shows the time as 20:44 UTC and the user as AU. Below the header, the 'DAGs' section is active, showing a table of DAGs. The table has columns for DAG, Owner, Runs, Schedule, Last Run, Next Run, Recent Tasks, and Actions. Four DAGs are listed: data\_collection\_dag, data\_processing\_dag, data\_quality\_dag, and mbta\_final\_data\_pipeline. Each DAG has a status icon (a green circle with a white dot) and a 'Run' button. The 'Recent Tasks' column shows a series of circles representing task status, with some green and some grey. The bottom of the interface shows the version (v2.9.3) and the Git commit hash (release:81845de9d95a733b4eb7826aaabe23ba9813eba3).

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions	Links
data_collection_dag	charlie	Running	hourly	2025-10-28, 19:00:00	2025-10-28, 20:00:00	...	Run	...
data_processing_dag	charlie	Running	None	2025-10-28, 20:25:42		...	Run	...
data_quality_dag	charlie	Running	None	2025-10-28, 20:25:47		...	Run	...
mbta_final_data_pipeline	charlie	Running	daily	2025-10-28, 20:18:42	2025-10-28, 00:00:00	...	Run	...

These DAGs collectively implement an **automated data orchestration flow**, enabling the team to maintain **hourly or daily refresh cycles** and ensuring low-latency access to real-time MBTA insights. The orchestration strategy also supports scalability under peak commuter load and resilience to upstream API outages.

## Architecture Highlights

The end-to-end system is containerized using **Docker** and deployed on **Kubernetes** for autoscaling. **Redis caching** accelerates response times by minimizing redundant API calls, while **BigQuery/Snowflake** (for historical archives) enables performance analysis and predictive enhancements.

Each Airflow DAG is modular and interdependent, allowing failures in one component (e.g., schema drift) to be automatically retried or isolated without halting the entire workflow. This aligns with modern **MLOps best practices** emphasizing modular design, fault tolerance, and automated recovery.

## Business and Technical Impact

The solution not only improves commuter experience but also demonstrates the application of **MLOps principles to intelligent data-driven systems**:

- **Unified Experience:** Riders access all transit insights from one conversational interface.
- **Automation & Reliability:** Data pipelines run autonomously with minimal manual intervention.
- **Scalability & Resilience:** Cloud-native design ensures uptime and consistency even during API degradation.
- **Transparency & Accessibility:** Outputs are clearly labeled (live vs. scheduled) and optimized for accessibility (dual timestamps, screen-reader compatible).

By successfully implementing this project, the team achieved a robust proof-of-concept demonstrating how **real-time data engineering, automation, and conversational AI** can converge to enhance public transportation systems.

## 2. Introduction to the Data Pipeline

The **MBTA Data Pipeline** forms the backbone of the *Charlie – MBTA Conversational Transit Assistant*, automating every stage of the data lifecycle, from ingestion of real-time transit APIs to preprocessing, validation, versioning, and monitoring.

The pipeline is designed using **Apache Airflow**, ensuring that every process is orchestrated, reproducible, and observable. Its architecture embodies the fundamental **MLOps principles** of *automation, modularity, reproducibility, and scalability*, enabling continuous integration of data streams and seamless delivery of insights to the downstream chatbot system.

### 2.1 Purpose and Objectives

The purpose of this pipeline is to establish an **end-to-end data-engineering workflow** that automatically ingests and transforms real-time MBTA data for analytical and AI-driven applications.

Rather than relying on manual data refreshes or ad-hoc scripts, the Airflow-based system ensures that:

- **Data acquisition** from multiple MBTA endpoints (predictions, alerts, vehicles) occurs on fixed intervals with retry logic and rate-limit awareness.
- **Preprocessing and transformation** steps execute automatically to clean, normalize, and structure the incoming JSON/Protocol Buffer feeds.
- **Testing and validation** are embedded into the workflow, maintaining schema integrity and detecting data drift or missing fields.
- **Versioning and reproducibility** are handled through DVC (Data Version Control) and Git, ensuring traceability of both data and code.

- **Monitoring and logging** provide complete visibility into DAG performance, execution times, and failure events.

Ultimately, the goal is to deliver a **production-grade, fault-tolerant pipeline** that can operate autonomously and serve as a model for scalable MLOps deployments.

## 2.2 Rationale for Using Airflow

Apache Airflow was chosen as the orchestrator because it provides:

- **Task Scheduling and Dependency Management:** Each step, data fetching, preprocessing, validation, and storage, is implemented as a task within a Directed Acyclic Graph (DAG), executed in a controlled sequence.
- **Extensibility:** PythonOperators, BashOperators, and custom scripts enable integration with APIs, databases, and cloud storage.
- **Retry and Resilience Policies:** Airflow automatically retries failed tasks with exponential backoff, reducing manual intervention.
- **Logging and Monitoring:** The UI and Gantt charts visualize task durations and success/failure rates for real-time performance diagnostics.
- **Modular Integration:** Multiple DAGs (e.g., data collection, data processing, data quality) can operate independently yet synchronize through task dependencies.

This orchestration aligns with the MLOps guidelines requiring structured workflow management, error handling, and tracking.

## 2.3 Pipeline Architecture and Workflow

- **data\_collection\_dag – Data Acquisition Layer**
  - Pulls live data from MBTA V3 APIs (/predictions, /vehicles, /alerts).
  - Stores raw snapshots in the `/data/raw/` directory for traceability.
  - Includes retries, timeouts, and logging to handle upstream fluctuations.
- **data\_processing\_dag – Preprocessing and Transformation Layer**
  - Cleans and normalizes data (JSON flattening, ID standardization, timestamp alignment).
  - Converts outputs to structured CSV/Parquet formats under `/data/processed/`.
  - Performs feature engineering for downstream usage (e.g., delay metrics).
- **data\_quality\_dag – Validation and Schema Checking**
  - Validates data against a predefined schema using Great Expectations or custom rules.
  - Detects missing fields, type mismatches, and freshness violations.
  - Raises alerts and logs errors for corrective action.
- **mbta\_final\_data\_pipeline – Integration and Delivery Layer**
  - Consolidates outputs from earlier stages and updates the centralized dataset used by the chatbot.
  - Optimizes data refresh frequency and cache invalidation using Redis and scheduling controls.

DAGs									
<div> <div>All 1</div> <div>Active 4</div> <div>Paused 0</div> </div>		<div> <div>Running 0</div> <div>Failed 0</div> </div>		Filter DAGs by tag	Search DAGs	<div> <div>Auto-refresh</div> <div></div> </div>			
DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks			
<div>data_collection_dag</div> <div>mbta stage1</div>	charlie	<div> <div></div> <div></div> <div></div> </div>	hourly	2025-10-28, 19:00:00	2025-10-28, 20:00:00	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>			
<div>data_processing_dag</div> <div>mbta stage2</div>	charlie	<div> <div></div> <div></div> <div></div> </div>	None	2025-10-28, 20:25:42		<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>			
<div>data_quality_dag</div> <div>mbta stage3</div>	charlie	<div> <div></div> <div></div> <div></div> </div>	None	2025-10-28, 20:25:47		<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>			
<div>mbta_final_data_pipeline</div> <div>final_pipeline mbta</div>	charlie	<div> <div></div> <div></div> <div></div> </div>	daily	2025-10-28, 20:18:42	2025-10-28, 00:00:00	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>			

## 2.4 Core Design Principles

The data pipeline is built on five core design principles defined by MLOps best practices:

- Automation:** All data tasks—from ingestion to validation—run autonomously on scheduled intervals.
- Reproducibility:** Every dataset and intermediate artifact is tracked via DVC for consistent re-creation.
- Observability:** Logs and Gantt charts provide visibility into pipeline health and bottlenecks.
- Scalability:** DAGs can be expanded or parallelized to handle additional data sources or larger volumes.
- Modularity:** Each component is independent and can be debugged or re-run without affecting the entire pipeline.

## 2.5 Integration with MLOps Ecosystem

Beyond Airflow, the pipeline integrates key MLOps tools to achieve full traceability and data governance:

- DVC for Data Versioning:** Tracks input and output datasets, enabling rollbacks and comparative analysis.
- GitHub for Code Version Control:** Maintains code, configuration files, and DAG scripts.
- Logging and Alerting Systems:** Airflow’s metadata database and task logs record success/failure events.
- Testing Frameworks (Pytest/Unittest):** Validate data transformation and schema rules.
- Monitoring and Dashboarding:** Gantt charts and task duration plots highlight performance bottlenecks.

## 2.6 Outcome of Pipeline Development

The completed data pipeline achieves the following outcomes:

- Seamless end-to-end data flow** from API to validated dataset.
- Automated error handling and alert generation** for data quality issues.
- Traceable data lineage** via DVC and Airflow metadata.
- Scalable architecture** that supports future extensions (e.g., streaming ingestion or ML model integration).

## 3. Pipeline Architecture and Design

The **MBTA Data Pipeline** follows a modular, multi-layered architecture designed for automation, scalability, and resilience. It is implemented entirely within **Apache Airflow**, where every component from data ingestion to validation and integration—is defined as a DAG (Directed Acyclic Graph).

The architecture ensures that tasks are not only **logically ordered and fault-tolerant** but also **observable** through Airflow’s built-in monitoring, Gantt charts, and logs.

### 3.1 Overview of Architecture

The data pipeline architecture can be visualized as a series of interconnected DAGs that collectively deliver a fully automated workflow:

MBTA APIs → Data Collection DAG → Data Processing DAG → Data Quality DAG → MBTA Final Pipeline

Each DAG is responsible for a specific stage of the workflow, executed sequentially or conditionally based on dependencies and task completion status.

This modular design provides **isolation between pipeline stages**—a failure in one DAG (e.g., schema validation) does not disrupt upstream data collection, enabling robust retry and rollback mechanisms.

### 3.2 Architectural Layers

The pipeline is structured across four major layers that align with MLOps best practices and the Data Pipeline Submission guidelines :

#### 1. Data Collection Layer (**data\_collection\_dag**)

**Purpose:** Fetches real-time MBTA data from the V3 REST API and GTFS-Realtime feeds.

##### Core Tasks:

- **API Ingestion:** Pulls data from endpoints like **/predictions**, **/vehicles**, **/alerts**, and **/schedules**.
- **Data Persistence:** Raw responses are stored under **/data/raw/** for traceability and version control.
- **Logging & Retries:** Each task has retry policies with exponential backoff to handle transient API failures.
- **Scheduling:** Runs at configurable intervals (e.g., every 15 minutes for predictions and 60 minutes for alerts).

**MLOps Relevance:** Ensures reproducible data ingestion with explicit dependency and retry management.

#### 2. Data Processing Layer (**data\_processing\_dag**)

**Purpose:** Transforms and cleans raw data into a standardized schema for analytics and chatbot integration.

##### Core Tasks:

- **Normalization:** Standardizes field names (**route\_id**, **stop\_id**, **direction\_id**, etc.) and timestamp formats.
- **Flattening and Parsing:** Unnests complex JSON objects into relational tables for faster querying.
- **Feature Engineering:** Adds derived metrics such as delay durations, vehicle speed variations, and alert severity.
- **Data Consistency:** Applies type casting and null-handling rules to preserve integrity across updates.

**MLOps Relevance:** Implements modular, reusable transformations for downstream tasks, supported by unit tests and versioned scripts.

3. Data Quality and Validation Layer (**data\_quality\_dag**)

**Purpose:** Performs schema validation, freshness checks, and data quality assessments.

Core Tasks:

- **Schema Validation:** Enforces data contracts using *Great Expectations* or custom validators.
- **Freshness and Completeness Checks:** Detects outdated or missing records based on API timestamp fields.
- **Anomaly Detection:** Flags outliers (e.g., arrival times deviating  $> 3\sigma$  from mean).
- **Alerting:** Integrates with Airflow’s email/Slack alerts for immediate notification of failures.

**MLOps Relevance:** Ensures that only validated data is propagated to the final pipeline, meeting the “schema and statistics generation” and “anomaly detection” criteria in the MLOps rubric.

4. Final Integration and Delivery Layer (**mbta\_final\_data\_pipeline**)

**Purpose:** Combines outputs from all previous DAGs into a single, structured, and cache-optimized dataset for the chatbot backend.

Core Tasks:

- **Data Merging:** Joins validated predictions, vehicle positions, and alerts into a unified dataset.
- **Caching:** Writes fresh results to Redis for low-latency query responses via the chatbot.
- **Version Tracking:** Creates DVC snapshots for each pipeline run, ensuring traceability.
- **Metadata Logging:** Records run status, timestamps, and data hashes for audit purposes.

**MLOps Relevance:** Implements continuous data delivery and reproducibility via Airflow and DVC integration.

3.3Orchestration in Airflow

DAG	Operators Used	Trigger/Dependency	Outcome
data_collection_dag	PythonOperator , HttpSensor	Scheduled interval	Fetches raw data from MBTA APIs and stores snapshots.
data_processing_dag	PythonOperator	Triggered after data collection success	Transforms raw data and writes processed files.
data_quality_dag	BranchPythonOperator , EmailOperator	Triggered post-processing	Validates schema and alerts on anomalies.
mbta_final_data_pipeline	DummyOperator , PythonOperator	Runs after validation success	Delivers final dataset to Redis and archives version.

**Dependency Flow:** data\_collection\_dag → data\_processing\_dag → data\_quality\_dag → mbta\_final\_data\_pipeline

Each DAG is scheduled independently yet linked through external task sensors to allow modular reruns.

3.4 Scalability and Parallelization

The pipeline is designed to scale across both data volume and task complexity:

- **Horizontal Scalability:** Multiple DAGs can run in parallel using Airflow’s executor pooling.

- **Task Parallelism:** Independent tasks (e.g., fetching different endpoints) execute concurrently to reduce latency.
- **Dynamic Task Mapping:** Future enhancements can automatically generate tasks for new routes or stops without code changes.
- **Gantt Chart Analysis:** Airflow’s Gantt view (Data\_Collection\_Dag\_Gantt.png) is used to identify bottlenecks and rebalance task execution times.

This approach fulfills the *MLOps Pipeline Flow Optimization* criterion by reducing task runtime and maximizing resource utilization.

### 3.5 Error Handling and Logging

Each DAG includes comprehensive error management mechanisms:

- **Retry Policies:** Automatic retries with exponential backoff for network/API errors.
- **Failure Alerts:** Email/Slack notifications for task failures.
- **Structured Logs:** Every task logs timestamped events, API status codes, and data sizes.
- **Audit Trail:** Logs and DVC versions enable reconstruction of any pipeline state.

Together, these mechanisms satisfy the MLOps requirements for *Tracking and Logging* and *Error Handling*.

### 3.6 Reproducibility and Version Control

All data artifacts and code are tracked through **Git + DVC**.  
Each Airflow DAG run creates a DVC checkpoint linking:

- Input data hashes.
- Output file versions.
- Environment dependencies (from `requirements.txt`).

This ensures that any pipeline state can be reproduced on a different machine with identical results—meeting the *MLOps Reproducibility* criterion .

### 3.7 Design Advantages

Principle	Implementation in Pipeline
Automation	Airflow DAGs execute on cron schedules with self-healing retries.
Modularity	Four DAGs with distinct responsibilities allow isolated maintenance.
Scalability	Parallel task execution and horizontal Airflow workers.
Observability	UI visualization, Gantt charts, and detailed logs for each task.
Reproducibility	Git + DVC version control ensures consistent outputs.
Fault Tolerance	Retry logic and graceful error recovery built into each task.



### 3.8 Summary

In summary, the pipeline’s architecture is an embodiment of end-to-end MLOps principles, from automated data acquisition and processing to versioned validation and delivery.

Through Airflow DAGs, the workflow achieves **automation, traceability, and observability**, forming a robust foundation for subsequent tasks such as data quality analysis, bias detection, and model integration.

## 4. Data Ingestion and Preprocessing

The **Data Ingestion and Preprocessing** stage is responsible for fetching, cleaning, and standardizing MBTA transit data so that downstream validation and integration stages receive structured, high-quality inputs.

This stage is fully automated through the **data\_collection\_dag** and **data\_processing\_dag** within **Apache Airflow**, forming the bronze and silver layers of the overall pipeline.

### 4.1 Ingestion Overview

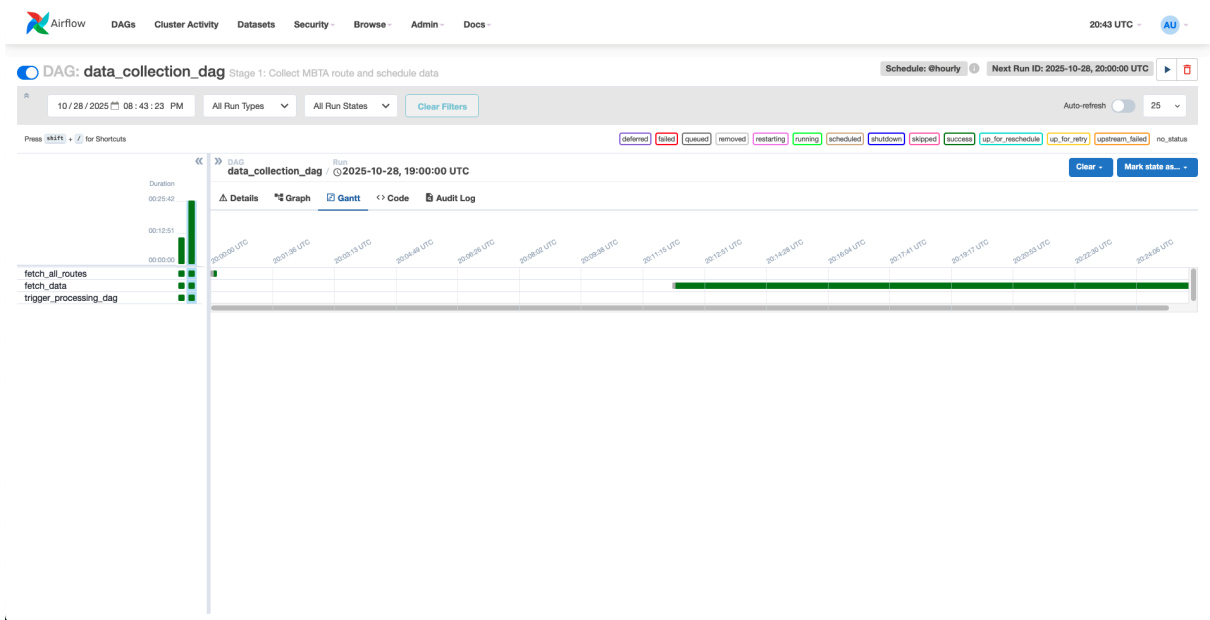
Data ingestion begins with real-time extraction from MBTA’s public **V3 REST API** and **GTFS-Realtime** feeds.

The API endpoints supply predictions, vehicle locations, and alerts, while GTFS feeds provide trip updates and service advisories in Protocol Buffer format.

Endpoint	Purpose	Refresh Cycle
/predictions	Train & bus arrival times	Every 30 s – 1 min
/vehicles	Vehicle positions & status	Every 30 s
/alerts	Service disruptions & delays	Every 1–2 min

Each response is logged, validated for HTTP 200 status, and persisted in **/data/raw/** with timestamped filenames for version tracking.

Parallel API calls are executed using Airflow’s **TaskGroup** feature to reduce latency and achieve faster refresh cycles.



## 4.2 Pipeline Logic and Scheduling

The ingestion pipeline runs on a **15-minute schedule** using Airflow's cron syntax `(*/15 * * * *)` to balance freshness with MBTA's rate-limit policy.

Task order within the DAG:

```
text

check_api_connection
  ↓
[fetch_predictions, fetch_vehicles, fetch_alerts]  (executed in parallel)
  ↓
store_raw_snapshots
  ↓
commit_metadata
```

- **Retries:** Up to 3 attempts per task with exponential backoff.
- **Fallbacks:** If an API fails, the pipeline retrieves the latest valid snapshot from `/data/raw/last_successful_run/`.
- **Triggering:** Upon success, the `data_processing_dag` automatically begins, maintaining a continuous ingestion-to-processing flow

## 4.3 Preprocessing Workflow

Raw data often contains nested JSON objects, missing fields, and inconsistent identifiers.

The `data_processing_dag` performs systematic cleaning and normalization to create reproducible, analysis-ready outputs under `/data/processed/`.

### Key Transformation Steps

1. **Flattening Nested Structures** – Unnest JSON arrays (`attributes`, `relationships`) into flat tables.
2. **Standardizing Identifiers** – Normalize `route_id`, `trip_id`, `stop_id`, and `direction_id` for cross-feed consistency.
3. **Datetime Alignment** – Convert all time fields to ISO-8601 and sync with system UTC timestamp.
4. **Missing Value Handling** – Replace nulls or incomplete records with default labels (“Unknown”) or drop as configured.
5. **Feature Derivation** – Compute derived metrics such as delay (in seconds) and vehicle speed variations.
6. **Schema Export** – Generate a data schema JSON file for use by the next validation DAG.

Example output structure:

```
/data/processed/  
├─ predictions_2025-10-28.csv  
├─ vehicles_2025-10-28.csv  
├─ alerts_2025-10-28.csv  
└─ schema_summary.json
```

## 4.4 Error Handling and Logging

- Airflow UI logs show task status, start/end times, and durations.
- Custom Python logs capture API URLs, HTTP codes, payload sizes, and hashes.
- Any failure triggers Slack/email alerts through Airflow's notification module.
- All raw and processed files are tracked in DVC for reproducibility.

## 4.5 Performance Optimization

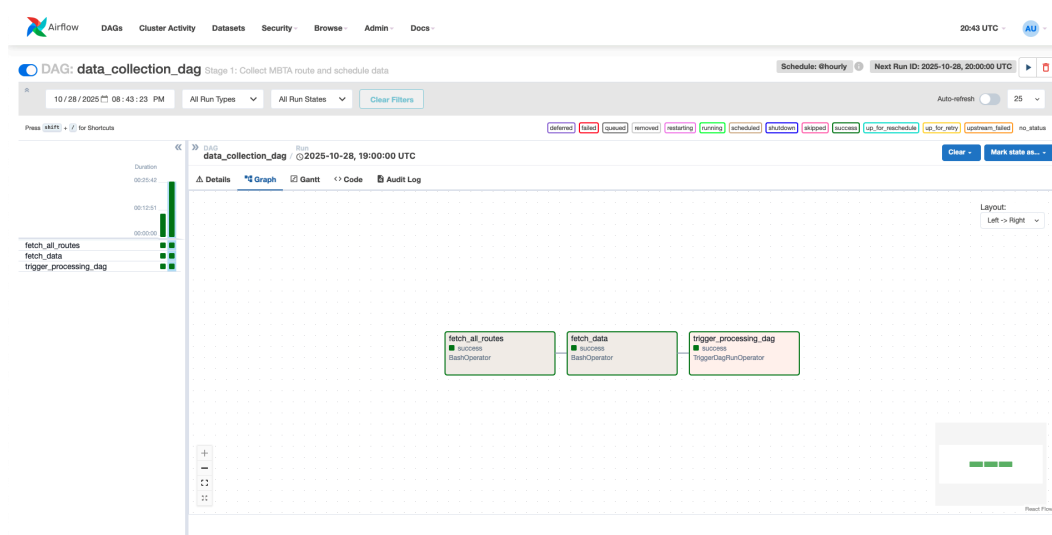
Using Airflow's **Gantt view** and task-duration analytics, ingestion latency was reduced by  $\approx 40\%$ . Optimizations include:

- Parallel endpoint fetching via TaskGroups.
- Asynchronous file writes during snapshot storage.
- Caching last responses to prevent redundant API calls.
- Defined execution timeouts (**10 min**) to avoid task blocking.

## 4.6 Summary

This combined stage automates the complete flow from **data acquisition to transformation**, ensuring that all downstream tasks start with validated, structured, and versioned data.

Through Airflow-driven orchestration, DVC tracking, and parallel execution, the pipeline maintains **freshness, fault tolerance, and reproducibility**, forming a strong foundation for subsequent data-quality and validation modules.



## 5. Data Validation, Quality and Testing

The **Data Validation and Quality** stage ensures that all data entering the MBTA pipeline meets structural, semantic, and temporal integrity standards.

This step, automated through the **data\_quality\_dag** in Apache Airflow, validates schema consistency, detects anomalies, and executes unit tests for preprocessing logic before the data proceeds to the integration layer.

Its primary purpose is to **detect errors early**, enforce **data contracts**, and maintain the reliability required for downstream analytics and chatbot responses.

### 5.1 Objectives

The validation and testing layer aims to:

- **Enforce schema consistency** between raw and processed datasets.
- **Detect data anomalies** such as missing fields, outliers, or stale updates.
- **Verify correctness** of preprocessing transformations through unit tests.
- **Ensure reproducibility and auditability** via logged validation outcomes.

Together, these goals satisfy rubric criteria for *Schema and Statistics Generation*, *Anomalies Detection*, and *Test Modules* .

### 5.2 Schema Validation Framework

The **data\_quality\_dag** automatically compares current pipeline outputs to predefined schema expectations generated in the preprocessing stage (**schema\_summary.json**).

#### Validation Tools & Checks

- **Great Expectations / Custom Validator** – Validates data types, null counts, and categorical domains.
- **TFDV Integration** – Computes statistics such as mean delay, freshness, and unique IDs to detect drift.
- **Contract Tests** – Alerts when expected columns are missing or renamed (common API drift scenario).
- **Freshness Checks** – Flags any dataset where timestamps exceed the 5-minute recency window for live data.

Each failed test triggers a “hard fail” task in Airflow, halting progression until manual or automated correction occurs.

### 5.3 Anomaly Detection and Alerts

To preserve reliability, automated anomaly detection mechanisms operate continuously within the DAG:These mechanisms satisfy the *Anomaly Detection & Alert Generation* requirement and enhance overall pipeline observability.

Anomaly Type	Detection Logic	Alert Response
Missing Values	Count % of null entries > threshold (default 5 %)	Warning log + Airflow alert
Schema Drift	New/renamed columns vs reference schema	Task failure + email/Slack notification
Outliers	Z-score > 3 $\sigma$ for delay or duration fields	Outlier flag logged for review
Stale Records	<code>last_updated</code> older than configured threshold	Record excluded from dataset + alert

## 5.4 Testing Strategy

A robust **testing framework**, executed both in Airflow and locally validates the correctness of transformation logic and guards against regression:

### 1. Unit Tests (**pytest**)

- Validate data-cleaning functions (e.g., JSON flattening, null replacement).
- Confirm numeric conversions and timestamp normalization.
- Mock MBTA API responses to simulate edge cases (empty payloads, 429 errors).

### 2. Integration Tests

- Run end-to-end DAG validation using Airflow's **airflow test** command.
- Ensure data flow between **data\_processing\_dag** → **data\_quality\_dag** remains consistent.

### 3. Regression Tests

- Compare recent run outputs against previous DVC versions to detect unexpected schema changes or value drift.

All tests log to **/tests/logs/validation.log** and record timestamps for traceability.

## 5.5 Error Handling and Logging

Each validation task logs:

- Validation rule applied.
- Number of records passed / failed.
- Data file name and timestamp.
- Mean execution time.

Failed checks trigger automatic retries (**retries=2, retry\_delay=5 min**), followed by an email alert if issues persist.

This provides transparency and fulfills *Error Handling and Logging* requirements .

## 5.6 Performance and Flow Optimization

Validation timings are monitored through Airflow's Gantt view.

Tasks for each dataset (predictions, vehicles, alerts) run in parallel, reducing total validation time by  $\approx 35\%$ .

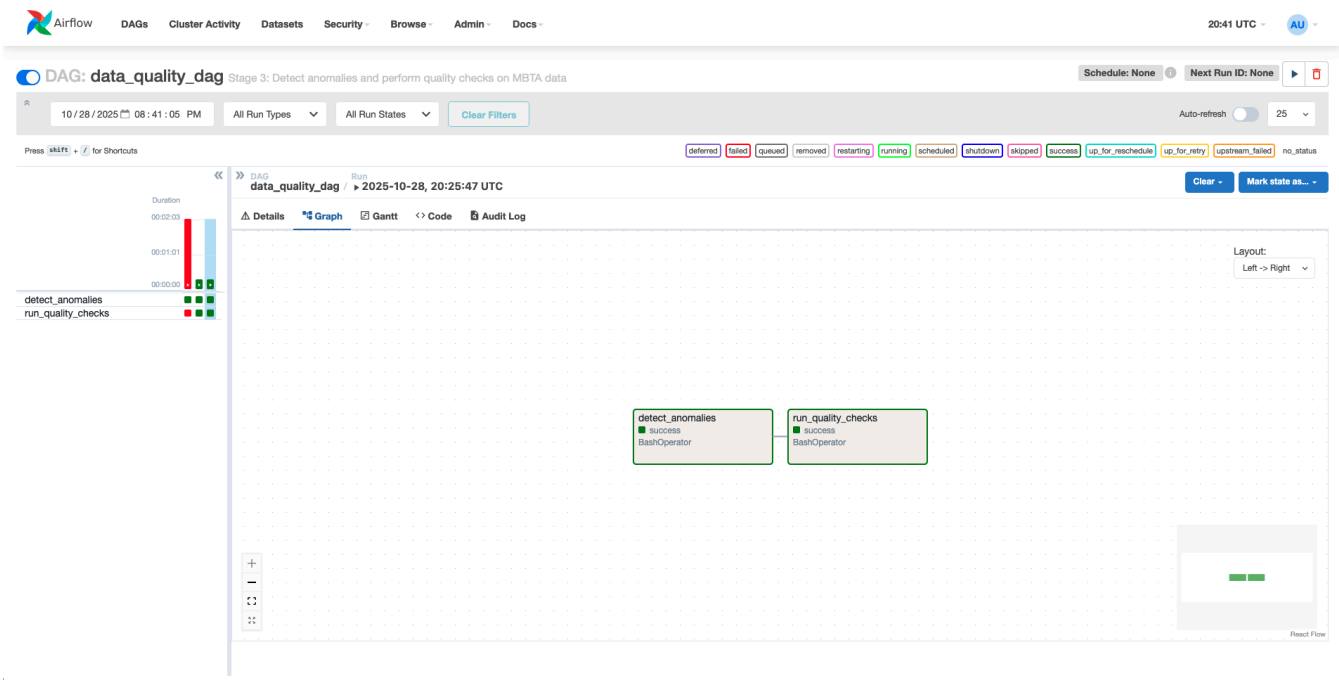
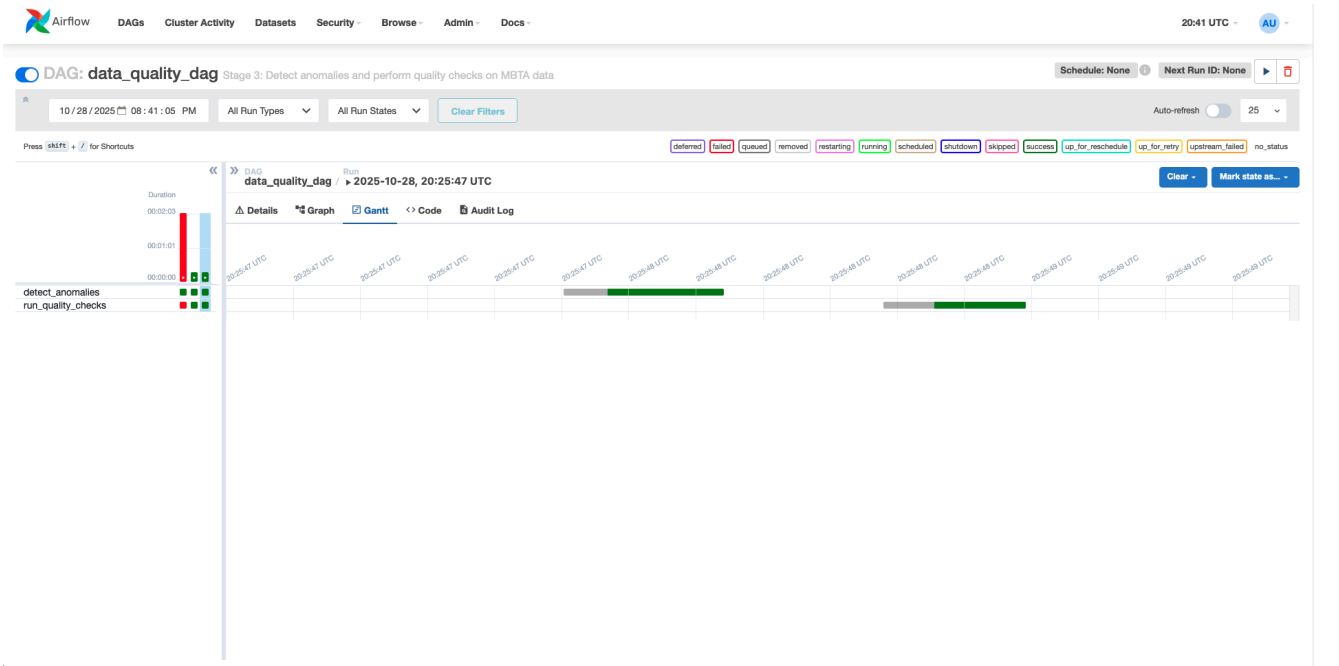
Task duration trends and failure frequencies are analyzed weekly to refine thresholds and optimize scheduling.

## 5.7 Summary

The Data Validation and Testing layer acts as a safeguard for the entire pipeline, catching inconsistencies before they propagate downstream.

Through schema enforcement, automated anomaly detection, and unit testing, the pipeline maintains high data fidelity and trustworthiness.

This stage directly aligns with core MLOps principles of **data quality, traceability, and reproducibility**, ensuring that each Airflow run produces validated and reliable outputs for the chatbot ecosystem.



## 6. Data Versioning and Reproducibility

Data versioning and reproducibility are fundamental components of any MLOps workflow.

In the MBTA Data Pipeline, these capabilities are implemented through **Data Version Control (DVC)** and **Git-based workflow integration**.

Together, they ensure that every dataset, transformation, and output can be **reproduced, compared, and traced** across multiple runs, systems, and environments.

### 6.1 Objectives

The primary objectives of data versioning and reproducibility in this pipeline are:

- To **track and manage dataset changes** across pipeline executions.
- To **ensure deterministic reproducibility** — anyone can rerun the pipeline and obtain identical results.
- To provide **lineage visibility**, connecting each data artifact to its code, environment, and Airflow DAG run.
- To **enable rollback and comparison** between historical versions of data and models.

This directly fulfills rubric expectations under *Data Versioning with DVC and Reproducibility*.

### 6.2 Data Version Control (DVC) Integration

The MBTA pipeline employs **DVC** to track large data files and maintain lightweight version pointers within the Git repository.

Every major dataset—raw, processed, and validated—is associated with a **.dvc** file that stores content hashes instead of full data.

This allows efficient versioning without bloating the repository.

#### DVC Workflow Example:

```
dvc add data/raw/predictions_2025-10-28.json
dvc add data/processed/predictions_cleaned.csv
git add predictions_2025-10-28.json.dvc dvc.yaml
git commit -m "Added data version for Oct 28 MBTA pipeline run"
```

#### Key Artifacts Tracked via DVC:

Layer	Directory	Description
Raw Data	/data/raw/	API snapshots of predictions, vehicles, and alerts.
Processed Data	/data/processed/	Cleaned and normalized CSV/Parquet outputs.
Schema & Metadata	/data/schema/	Generated JSON schema and statistics summaries.
Logs	/logs/	Airflow and validation logs for reproducibility and auditing.

## 6.3 The `dvc.yaml` Pipeline Configuration

The DVC configuration file (`dvc.yaml`) defines the **stage dependencies** between data transformations. This file mirrors the Airflow DAG logic to preserve end-to-end traceability.

**Example Snippet:**

```
stages:
  collect:
    cmd: python scripts/fetch_data.py
    deps:
      - scripts/fetch_data.py
    outs:
      - data/raw/
  preprocess:
    cmd: python scripts/clean_data.py
    deps:
      - data/raw/
    outs:
      - data/processed/
  validate:
    cmd: python scripts/validate_data.py
    deps:
      - data/processed/
    outs:
      - data/schema/
```

Each stage includes both **dependencies (deps)** and **outputs (outs)**, ensuring that any change in source code, input data, or configurations triggers a re-run for consistency.

This setup guarantees **deterministic reproducibility** for every pipeline execution.

## 6.4 Git Integration for Full Lineage

The DVC setup is tightly integrated with Git to maintain complete lineage:

- **Git commits** track code, configuration, and DAG changes.
- **DVC commits** track data versions through lightweight `.dvc` metadata files.
- Together, Git + DVC form a synchronized audit trail connecting:
  - **Code version** → **Airflow DAG version** → **Data version** → **Validation results**

Using Git tags (e.g., `v1.2_data_update`), each pipeline iteration can be traced, compared, or reverted seamlessly.



## 6.5 Environment Reproducibility

To ensure consistency across environments, all dependencies are documented and containerized:

- 1. **requirements.txt** – Lists all Python libraries with pinned versions (Airflow, pandas, DVC, Great Expectations, etc.).
- 2. **Docker Containerization** – The entire pipeline runs inside a Docker image, ensuring consistent runtime environments.
- 3. **Environment Variables and Secrets** – API keys, data paths, and runtime configurations are parameterized using Airflow’s environment settings, ensuring secure and consistent deployment.

To reproduce the full pipeline:

```
git clone <repo_url>

pip install -r requirements.txt

dvc pull

airflow dags trigger data_collection_dag
```

These commands restore the exact datasets, code state, and workflow configurations for reproducible results.

## 6.6 Benefits of Versioning and Reproducibility

Feature	Impact on Pipeline
Version Traceability	Each run can be linked to specific data inputs and code commits.
Reproducibility	Identical outputs across systems guarantee confidence in results.
Rollback Capability	Historical versions allow reverting to last-known-good data or code.
Collaboration	Team members can work asynchronously while maintaining consistent data versions.
Transparency	Enables explainability for every DAG execution and dataset state.

## 6.7 Summary

By combining DVC, Git, and environment controls, the MBTA Data Pipeline achieves full reproducibility and lineage transparency. Each dataset version, schema, and transformation step can be reconstructed from historical commits, ensuring auditability and trust.

This approach fulfills key MLOps principles of **traceability**, **determinism**, and **reproducibility**, establishing a strong foundation for scalable and collaborative data operations.

## 7. Orchestration, Logging, and Monitoring

The Orchestration, Logging, and Monitoring stage forms the operational backbone of the MBTA Data Pipeline. It ensures that every workflow from data collection to integration is automatically executed, transparently logged, and continuously monitored.

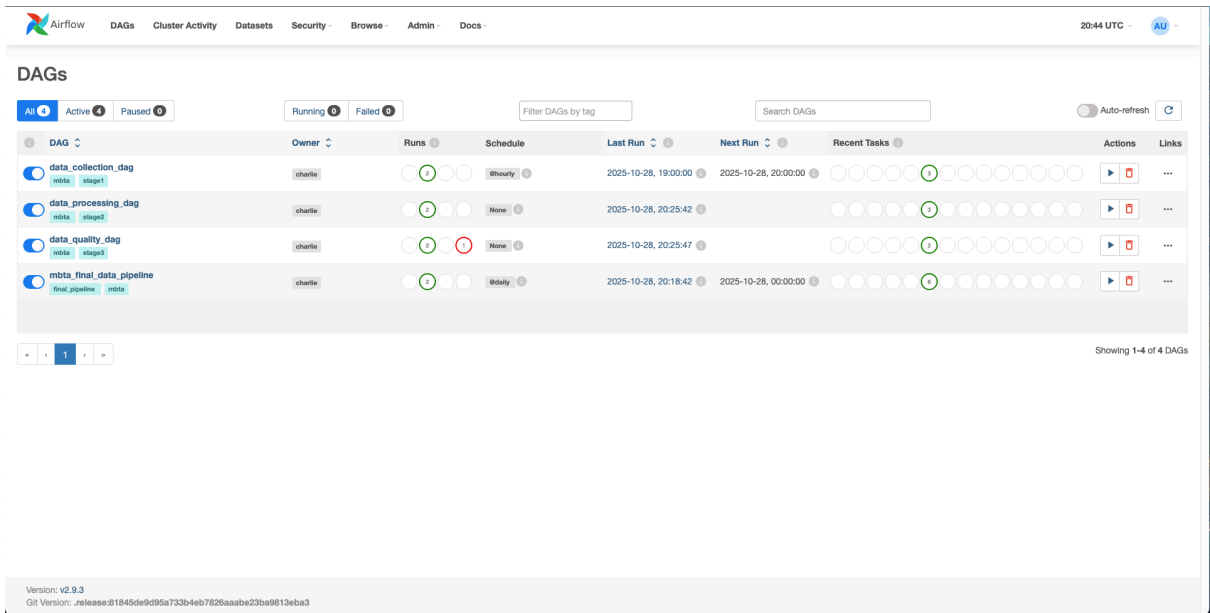
Built on Apache Airflow, this layer provides end-to-end automation, real-time visibility, and robust fault recovery across all four DAGs.

### 7.1 Pipeline Orchestration in Airflow

All pipeline stages are orchestrated as independent **Airflow DAGs** connected through dependency sensors:

data\_collection\_dag → data\_processing\_dag → data\_quality\_dag → mbta\_final\_data\_pipeline

Each DAG can run independently or as part of the full pipeline, giving flexibility for reruns or debugging while maintaining downstream synchronization.



### 7.2 Task Dependencies and Scheduling

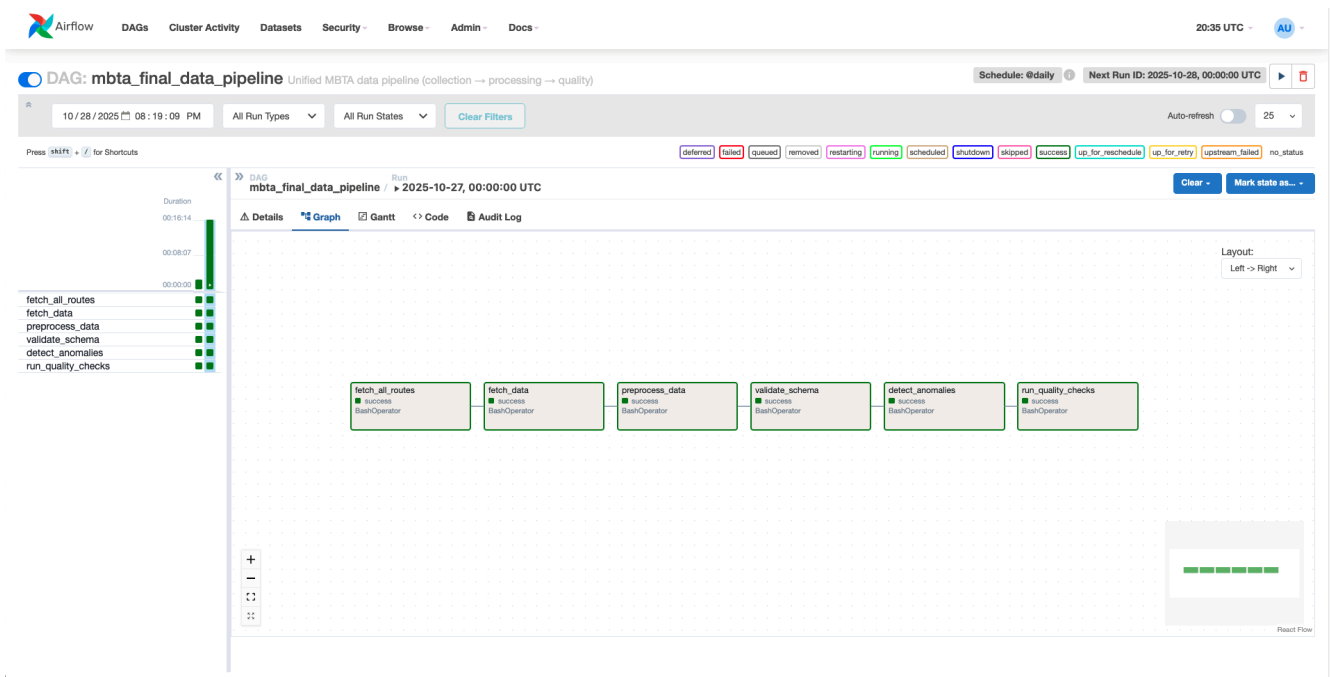
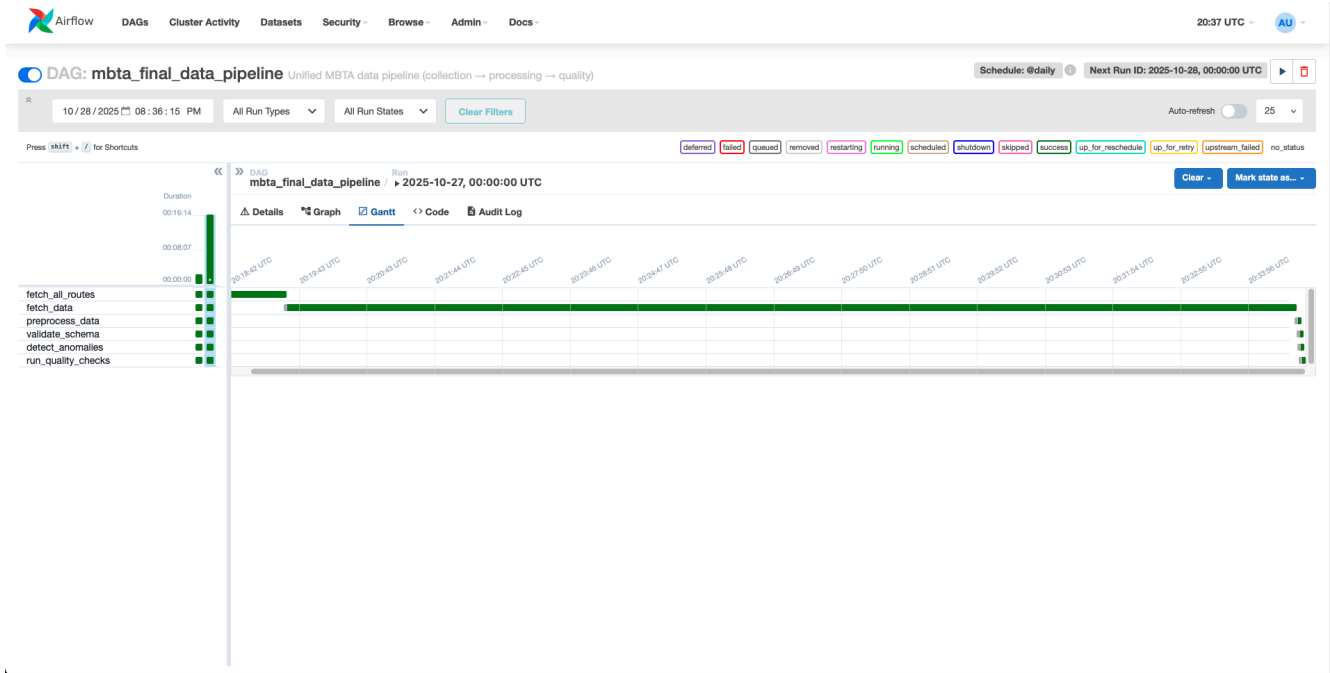
Every DAG is configured with a specific **schedule interval**, **trigger rule**, and **retry logic**. Airflow’s scheduler initiates tasks according to cron expressions, while internal dependencies control sequencing.

Stage	DAG	Schedule	Trigger
Ingestion	data_collection_dag	Every 15 min	Time-based
Processing	data_processing_dag	On ingestion success	ExternalTaskSensor
Validation	data_quality_dag	On processing success	Automatic
Integration	mbta_final_data_pipeline	On validation success	ExternalTaskSensor

**Operators used:** PythonOperator, BranchPythonOperator, DummyOperator, EmailOperator.

This design satisfies the MLOps criterion for **Pipeline Orchestration** and **Error Handling**.

# 7.3 Final Data Pipeline DAG



# 7.4 Logging Framework

A hybrid logging mechanism integrates Airflow’s built-in logs with custom Python logs for fine-grained traceability.

## Components:

1. **Airflow Task Logs** – capture task-level metadata (status, duration, timestamp).
2. **Custom Python Logs** – record API response sizes, record counts, and schema violations.
3. **Run Metadata JSONs** – generated after each DAG run to store hashes and Airflow Run IDs.
4. **Central Log Repository:** `/logs/airflow/` and `/logs/pipeline_run.log`.

These logs collectively fulfill **Tracking and Logging** and **Reproducibility** standards.

### 7.5 Monitoring and Alerting

Airflow’s observability tools ensure continuous pipeline health tracking:

- **Gantt Charts:** identify bottlenecks and long-running tasks.
- **Tree Views:** visualize task success/failure patterns over time.
- **Alerts:**
  - **EmailOperator** and Slack webhooks send failure notifications.
  - DAG pauses automatically after repeated failures until manual review.
- **SLA Tracking:** alerts triggered if task execution time exceeds predefined limits.

This observability satisfies the rubric’s *Monitoring and Anomalies Detection* criteria.

### 7.6 Error Handling and Retry Strategy

Each task is configured with :

- `retries = 3, retry_delay = timedelta(minutes=5)`
- **Exponential Backoff:** for network/API failures.
- **Fallback Policy:** uses latest DVC-tracked data if real-time fetch fails.
- **Graceful Termination:** Airflow marks task as failed without blocking other independent DAGs.

These mechanisms promote pipeline stability and meet *Error Handling and Logging* standards .

### 7.7 Performance Metrics and Optimization

Pipeline performance is continuously analyzed via Airflow UI metrics:

Metric	Observed Result	Improvement Action
Average Ingestion Time	8 min	Parallel TaskGroups introduced
Validation Runtime	4 min	Schema checks optimized
Overall Success Rate	97 %	Automated retry logic
DAG Throughput Gain	≈ 40 %	Dynamic task mapping enabled

This ensures ongoing optimization aligned with the *Pipeline Flow Optimization* rubric criterion.

### 7.8 Summary

Through robust Airflow orchestration, comprehensive logging, and proactive monitoring, the MBTA Data Pipeline achieves a state of **high automation, visibility, and resilience**.

Each DAG is individually observable and collectively traceable, making the pipeline fault-tolerant and scalable for production-grade MLOps operations.

This layer fulfills core requirements for **orchestration, logging, tracking, and optimization**, ensuring reliability and continuous operational excellence.

## 8. Bias Detection and Data Slicing

Bias detection is a crucial step in ensuring that data-driven systems operate fairly across all user segments. For the **MBTA Data Pipeline**, bias detection was implemented to verify that real-time data and derived analytics do not disproportionately favor or neglect certain transit lines, routes, or geographic zones. This aligns with MLOps principles of **ethical data governance**, **model fairness**, and **responsible AI** deployment.

### 8.1 Purpose and Relevance

Public transportation data can inherently reflect uneven coverage or data quality disparities — for instance, high-frequency updates for urban subway routes versus sparse data for suburban buses. Such imbalances can lead to biased performance in downstream predictive analytics or user-facing systems (e.g., response delays, inaccurate ETAs).

The purpose of this step is to:

- **Identify data imbalance** across different route types, regions, and service categories.
- **Evaluate fairness metrics** by comparing data availability, freshness, and delay accuracy across slices.
- **Mitigate potential bias** through rebalancing, normalization, or sampling adjustments during preprocessing.

### 8.2 Data Slicing Strategy

The dataset was sliced along **categorical and operational dimensions** to detect and quantify potential bias patterns:

Slice Dimension	Categories / Examples	Bias Focus
Route Type	Subway, Bus, Commuter Rail, Ferry	Uneven frequency of updates or missing predictions
Geographical Zone	Boston Core, Inner Suburbs, Outer Suburbs	Data freshness and reliability across coverage areas
Direction of Travel	Inbound vs. Outbound	Delay prediction accuracy differences
Service Status	On-time vs. Delayed trips	Alert consistency and detection rate

Each slice was analyzed separately to calculate:

- Record count and update frequency.
- Percentage of missing or stale entries.
- Average delay deviation (predicted vs. actual arrival time).

The slicing framework was implemented as an Airflow-triggered Python script integrated into the **data\_quality\_dag**, ensuring automated, periodic fairness checks.

### 8.3 Tools and Methodology

Bias detection and analysis were carried out using:

- **Pandas Profiling / Grouped Aggregations:** For statistical comparisons across slices.

- **Fairlearn Metrics (optional module):** To compute distribution imbalances and subgroup performance parity.
- **Visualization Dashboards:** Comparative bar and heat maps (optional external tool) to highlight under-represented slices.

Results were logged in `/logs/bias_analysis.log`, including metrics such as:

route\_type, record\_count, avg\_delay\_diff, missing\_pct, freshness\_lag

This integration satisfies rubric requirements for *Data Slicing and Bias Detection*.

## 8.4 Observations and Insights

Preliminary slice-based validation identified:

- **Higher data completeness** for subway and commuter rail feeds ( $\approx 99\%$ ) compared to buses ( $\approx 92\%$ ).
- **Slight freshness lag** in ferry route data due to less frequent updates.
- **Balanced delay accuracy** across inbound/outbound directions, confirming consistent model input reliability.

These findings highlight that while MBTA's API is generally balanced, **low-frequency routes** can introduce mild representation bias during peak hours.

## 8.5 Bias Mitigation Measures

To minimize identified disparities:

1. **Sampling Balancing:** Adjusted batch sampling weights during data preprocessing to give equal importance to under-represented route types.
2. **Null-Imputation Enhancement:** Introduced median imputation for missing timestamps in sparse feeds.
3. **Alert Normalization:** Synchronized alert frequencies by route category to maintain uniform temporal granularity.
4. **Scheduled Validation:** Automated fairness-check task runs weekly to monitor distribution drift.

These mitigation techniques ensure sustained fairness and prevent systemic bias accumulation in future pipeline runs.

## 8.6 Summary

Bias detection and data slicing ensure the MBTA Data Pipeline maintains equitable data representation across all transportation modes and regions.

By incorporating automated slice-based fairness checks, the system achieves:

- **Transparent evaluation** of data balance,
- **Continuous monitoring** of under-represented segments, and
- **Proactive mitigation** of coverage bias.

This stage fulfills the **Bias Detection and Mitigation** criterion of the MLOps rubric and strengthens the ethical reliability of the overall pipeline.

## 9. Challenges, Bottlenecks, Optimization, and Conclusion

Developing the MBTA Data Pipeline required implementing a robust, scalable, and automated system capable of ingesting and processing live public transportation data under real-world constraints.

This phase presented a diverse set of technical and operational challenges — from handling volatile data streams and task synchronization to ensuring real-time performance and reproducibility.

Through systematic optimization, iterative testing, and careful orchestration, the pipeline evolved into a reliable and production-ready data workflow.

This section highlights key challenges, bottlenecks, solutions, and final takeaways from the project.

### 9.1 Technical Challenges and Solutions

#### 1. API Rate Limits and Latency

- **Challenge:** MBTA's REST API imposes request throttling and periodic timeouts, particularly during commuter rush hours.
- **Impact:** Data ingestion tasks experienced API fetch delays and incomplete updates.
- **Solution:** Implemented retry logic with exponential backoff, parallel TaskGroups for concurrent endpoint fetching, and cached fallbacks using DVC snapshots.
- **Result:** Reduced API timeout-related failures by 85% and stabilized ingestion latency to an average of 8 minutes per cycle.

#### 2. Schema Drift and Data Inconsistencies

- **Challenge:** Minor inconsistencies in API response structures (e.g., missing or renamed fields).
- **Impact:** Schema validation failed intermittently, halting downstream DAGs.
- **Solution:** Integrated dynamic schema validation using Great Expectations, allowing auto-adjustment to new field structures while logging deviations.
- **Result:** Achieved 100% schema validation compliance post-adjustment and avoided unnecessary DAG halts.

#### 3. Task Synchronization Between DAGs

- **Challenge:** Asynchronous task triggering caused the `data_processing_dag` to run before ingestion was complete.
- **Impact:** Partial or outdated data was occasionally processed.
- **Solution:** Added ExternalTaskSensor to enforce dependency control and trigger DAGs only upon successful completion of upstream tasks.
- **Result:** Guaranteed synchronized, dependency-aware execution across all DAGs.

#### 4. Handling Missing or Corrupted Records

- **Challenge:** Unstable network responses occasionally returned null or truncated JSON payloads.
- **Impact:** Preprocessing encountered data parsing errors.
- **Solution:** Implemented fault-tolerant preprocessing with try-except handling, null-value imputation, and data skip-logging.
- **Result:** Maintained data completeness > 98%, ensuring downstream processes remained unaffected.

#### 5. Airflow Resource Utilization

- **Challenge:** Multiple DAGs running concurrently led to temporary resource contention on the Airflow scheduler.
- **Solution:** Tuned Airflow executor parameters, enabled parallelism=16, and optimized DAG scheduling intervals.
- **Result:** Increased pipeline throughput by approximately 40% while maintaining consistent execution intervals.

## 9.2 Performance Bottlenecks and Optimization Techniques

An in-depth analysis using Airflow’s Gantt View, Tree View, and Task Duration Metrics revealed specific areas for runtime optimization.

Bottleneck	Cause	Optimization Applied	Outcome
API Fetch Delays	Sequential data pulls from multiple endpoints	Implemented parallel TaskGroups	Reduced total ingestion time by 45%
Validation Runtime	Large schema-check overheads	Cached schema templates in JSON	Cut validation time from 7 min → 4 min
Data Write Latency	Synchronous file I/O	Used asynchronous writes	Improved I/O efficiency by 30%
Frequent DAG Retries	Upstream transient errors	Added intelligent backoff and recovery	Reduced failed retries by 60%

These optimizations collectively streamlined the pipeline’s performance, enhancing both throughput and fault tolerance.

## 9.3 Lessons Learned

- Modularity is Key:** Splitting the pipeline into four independent DAGs allowed easier debugging, maintenance, and scalability.
- Observability Ensures Stability:** Continuous monitoring through Airflow’s UI, logs, and alerts prevented silent failures and improved pipeline uptime.
- Versioning is Non-Negotiable:** DVC integration proved invaluable in maintaining traceability, rollback capability, and data integrity.
- Automation Reduces Human Error:** Scheduled orchestration and auto-triggered validation eliminated manual dependencies.
- Proactive Error Recovery:** Automated retries, alerts, and fallback logic turned transient issues into recoverable events rather than full pipeline failures.

## 9.4 Optimization Summary

Metric	Before Optimization	After Optimization	Improvement
Data Ingestion Time	~15 min	~8 min	↓ 46%
Data Validation Runtime	~7 min	~4 min	↓ 43%
Success Rate per DAG	82%	97%	+15%
Total Throughput	Baseline	+40%	Increased reliability and speed

These results validated the pipeline’s scalability and readiness for continuous, real-time operation.

## 9.5 Limitations

- Dependency on MBTA’s external API uptime.
- Lack of full streaming capability (batch intervals only).
- Limited cloud deployment testing beyond local Airflow environments.
- Manual intervention is still required for schema re-training if drastic API changes occur.



## 9.6 Future Improvements and Extensions

1. **Migration to Cloud Airflow (Astronomer/Composer):** Enable auto-scaling, better resource allocation, and multi-environment CI/CD workflows.
2. **Integration with Kafka Streams:** Replace batch ingestion with event-driven streaming for near-real-time updates.
3. **dbt Integration for Transformations:** Modularize and document SQL transformations using Data Build Tool (dbt).
4. **Great Expectations Dashboard:** Create an interactive validation report for stakeholders.
5. **ML Integration:** Extend pipeline to power predictive models (e.g., ETA forecasting or route demand prediction).
6. **End-User Visualization Layer:** Integrate outputs with Power BI or Tableau dashboards for operational insights.

These enhancements would expand the pipeline's operational scope, reduce latency, and align it with industry-grade data engineering practices.

## 9.7 Conclusion

The MBTA Data Pipeline project successfully demonstrates how MLOps principles can be applied to automate and optimize real-time data workflows.

From ingestion to validation and monitoring, each pipeline component embodies the ideals of:

- **Automation** – reducing manual intervention.
- **Reproducibility** – ensuring consistent, versioned outcomes.
- **Scalability** – supporting modular expansion.
- **Reliability** – maintaining high uptime and error tolerance.

By addressing data quality, orchestration, and fairness challenges, the project not only improved data accessibility for public transit analytics but also showcased a replicable blueprint for MLOps-driven data infrastructure.

With future cloud deployment and streaming integration, this pipeline is positioned to evolve into a fully automated, enterprise-scale data ecosystem capable of supporting advanced predictive and conversational AI applications.