

IMDb Team Project - GROUP 3

Architecture Design

BRONZE LAYER (Raw Ingestion)

- Technology: Auto Loader (cloudFiles) + Delta Tables
- Schema evolution enabled
- Streaming ingestion from TSV files
- Metadata tracking: _source_file, _loaded_at
- Change Data Feed enabled
- Basic validation: ID format checks

SILVER LAYER (Cleansed Data)

- Technology: Delta Live Tables (Streaming)
- Data type conversions (String → Int/Double/Boolean)
- Data cleansing and standardization
- Range validations (years, ratings, runtime)
- Array transformations (split multi-valued fields)
- @dlt.expect and @dlt.expect_or_drop for quality
- Timestamp tracking: source_loaded_at, silver_processed_at

GOLD LAYER (Business-Ready Dimensional)

- Technology: Delta Live Tables (Streaming)
- Star schema: Dimensions + Facts + Bridge tables
- SCD Type 2 for DIM_PERSON
- Reference data enrichment (ISO codes)
- Business logic applied
- Optimized for analytics

QA LAYER (Quality Assurance)

- Load audit tracking (qa_load_audit)
- Row count history (qa_row_count_history)
- Data quality metrics (qa_data_quality)
- Schema drift detection
- Quality scoring system

The screenshot shows the Databricks workspace interface. The left sidebar contains navigation links like Home, Workspace, Catalog, Jobs & Pipelines, Marketplace, SQL, SQL Editor, Queries, Dashboards, Genie, Alerts, Query History, and SQL Warehouses. The main area has tabs for 'imdb' and 'imdb'. The 'imdb' tab is active, displaying a Pipeline configuration panel with 'Last runs' (4 runs, 0m 65s) and a Pipeline assets panel showing a new pipeline named 'New Pipeline 2025-12-0...'. Below these are sections for explorations (sample_exploration), transformations (bronze), and utilities (util.py). A file browser on the right lists files: gold.py, qa_loading.py, silver.py, util.py, and README.md. The central workspace shows a 'Pipeline graph' with a 'Maximized' view, displaying a complex network of nodes and edges representing the data pipeline. Below the graph is a 'Tables' section with 30 rows and a 'Performance' section with 29 rows. At the bottom, there's a 'Query performance' table.

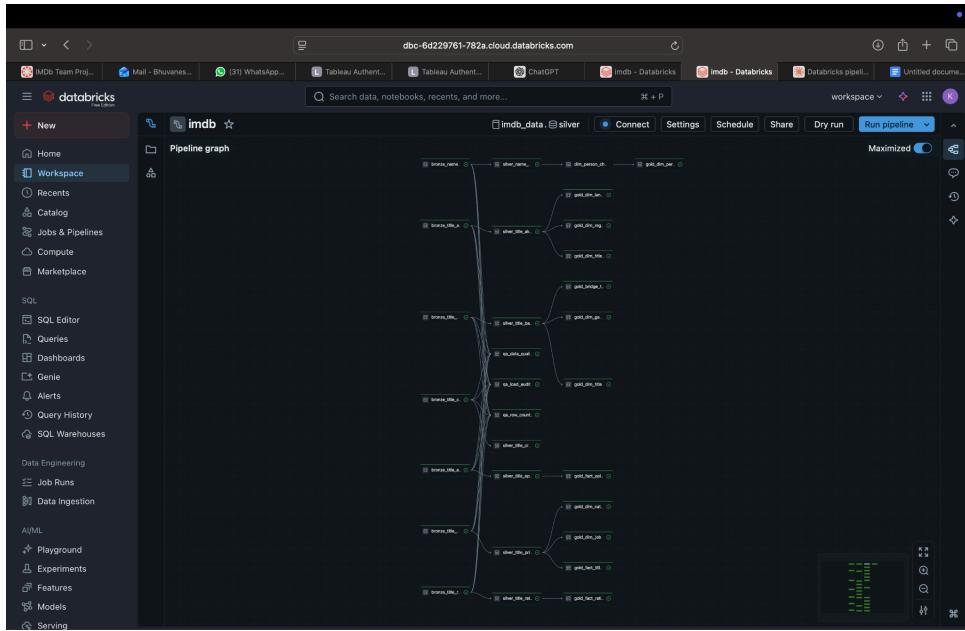
The screenshot shows the Databricks workspace interface. On the left, the sidebar includes links for Home, Workspace (selected), Recents, Catalog, Jobs & Pipelines, Compute, Marketplace, SQL, SQL Editor, Queries, Dashboards, Genie, Alerts, Query History, and SQL Warehouses. The main area displays a pipeline named 'imdb'. The pipeline configuration notebook contains the following code:

```
Lakeflow Pipelines Editor: ON

# Run file: /mnt/dbs/imdb/pipeline_imdb.py

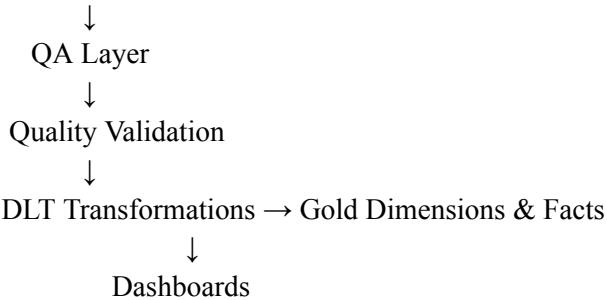
1 # InDB SILVER LAYER - DELTA LIVE TABLES (STREAMING)
2
3 # Pipeline Target Schema: silver
4 spark.sql("USE SCHEMA silver")
5 #
6
7 # Pipeline Target Schema: silver
8 # InDb SILVER LAYER - DELTA LIVE TABLES (STREAMING)
9 # Pipeline Target Schema: silver
10 # Cleaned, typed, and validated data from Bronze
11 #
12
13 import dtl
14 from pyspark.sql.functions import (
15     col, lit, when, trim, upper, lower, regexp_replace, split,
16     current_timestamp
17 )
18 from pyspark.sql.types import IntegerType, DoubleType
19
20
Tables: 30 Performance: 29
```

The Pipeline graph on the right shows the data flow betweenbronze, silver, and gold layers, with various transformations and triggers.



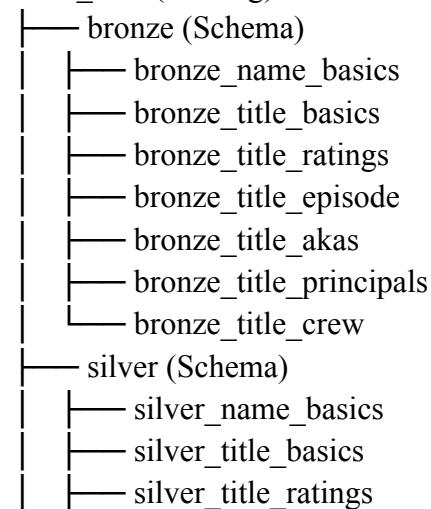
DATAFLOW DIAGRAM

TSV Files → Auto Loader → Bronze Tables → DLT Streaming → Silver Tables



Schema Architecture

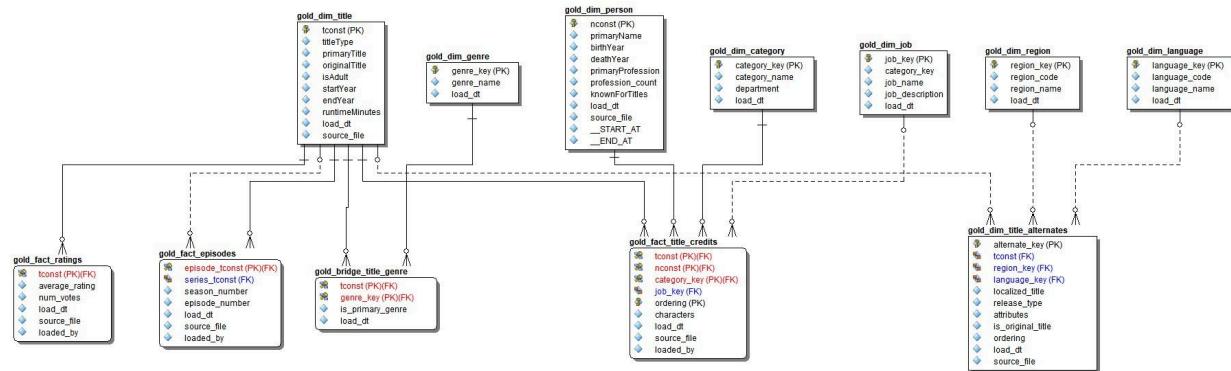
imdb_data (Catalog)



```

silver_title_episode
silver_title_akas
silver_title_principals
silver_title_crew
gold (Schema)
gold_dim_person (SCD Type 2)
gold_dim_title
gold_dim_genre
gold_dim_category
gold_dim_job
gold_dim_region
gold_dim_language
gold_dim_title_alternates
gold_bridge_title_genre
gold_fact_ratings
gold_fact_episodes
gold_fact_title_credits
qa (Schema)
qa_load_audit
qa_row_count_history
qa_data_quality

```



Standard metadata columns-

Column Name	Data Type	Purpose	Layer
source_file	STRING	Original file path	Bronze
_loaded_at	TIMESTAMP	Bronze ingestion time	Bronze
source_loaded_at	TIMESTAMP	Bronze timestamp reference	Silver/Gold
silver_processed_at	TIMESTAMP	Silver processing time	Silver
load_dt	TIMESTAMP	Gold load time	Gold
Source_file	STRING	Source identifier	Gold
loaded_by	STRING	ETL job identifier	Gold

ETL Pipeline Implementation

Pipeline Architecture

Pipeline Components:

1. **explorations** - Setup and initial data quality checks
2. **transformations/bronze.py** - Raw data ingestion
3. **transformations/silver.py** - Data cleansing and validation
4. **transformations/gold.py** - Dimensional modeling
5. **transformations/qa_loading.py** - Quality assurance
6. **utilities/utils.py** - Helper functions

Bronze Layer Implementation

Key Features:

- Streaming ingestion from Volume paths
- Schema hints for all columns
- Automatic schema evolution (addNewColumns)
- Change Data Feed enabled
- Checkpoint management for incremental loads
- IMDb null handling ('\\N' → NULL)

Bronze Tables Created:

1. bronze_name_basics
2. bronze_title_basics
3. bronze_title_ratings
4. bronze_title_episode
5. bronze_title_akas
6. bronze_title_principals
7. bronze_title_crew

Validations:

- ID format checks (tt% for titles, nm% for persons)
- NOT NULL constraints on primary keys
- Metadata tracking for lineage

Silver Layer Implementation

Key Transformations:

Data Type Conversions

String → Integer

```
when(col("birthYear").rlike("[0-9]{1,4}"), col("birthYear").cast(IntegerType()))
```

String → Double

```
col("averageRating").cast(DoubleType())
```

String → Boolean

```
when(col("isAdult") == "1", lit(True)).otherwise(lit(False))
```

Data Standardization

```
# Region codes: UPPER case
```

```
upper(trim(col("region")))
```

```
# Language codes: lower case
```

```
lower(trim(col("language")))
```

```
# Category standardization
```

```
lower(trim(col("category")))
```

Array Transformations

```
# Split comma-separated values into arrays
```

```
when(col("primaryProfession").isNotNull(),
```

```
    split(trim(col("primaryProfession")), ","))
```

```
# Split genres
```

```
when(col("genres").isNotNull(),
```

```
    split(trim(col("genres")), ","))
```

Data Cleansing

```
# Remove JSON brackets from characters field
```

```
regexp_replace(regexp_replace(trim(col("characters")), "^\\"[\"", "\"), "\""\\]$\", "")
```

Silver Tables Created:

1. silver_name_basics
2. silver_title_basics
3. silver_title_ratings
4. silver_title_episode
5. silver_title_akas
6. silver_title_principals
7. Silver_title_crew

Gold Layer Implementation

Dimension Tables

DIM_PERSON (SCD Type 2)

- Tracks historical changes in: primaryName, deathYear, primaryProfession
- Uses dlt.apply_changes() for automatic SCD Type 2 management
- Includes __start_at, __end_at, __current columns

gold_dim_title - Title/movie details

gold_dim_genre - Genre lookup (exploded from arrays)

gold_dim_category - Role category with department classification

gold_dim_job - Job titles by category

gold_dim_region - Region codes enriched with ISO country names

gold_dim_language - Language codes enriched with ISO language names

gold_dim_title_alternates - Alternative titles by region/language

Reference Data Integration:

Bridge Tables

gold_bridge_title_genre

- Handles many-to-many relationship between titles and genres
- Explodes genre arrays into individual records
- Includes is_primary_genre flag

Fact Tables

Gold_fact_ratings

- Granularity: One record per title
- Measures: average_rating, num_votes
- Foreign Keys: tconst → dim_title

Gold_fact_episodes

- Granularity: One record per episode
- Links episodes to their parent series
- Includes season_number, episode_number

gold_fact_title_credits

- Granularity: One record per person per title per role
- Links cast/crew to titles
- Includes ordering, job, characters
- Foreign Keys: tconst → dim_title, nconst → dim_person

QA Layer Implementation

Data Quality & Validation

QA Layer Implementation

- Purpose: Automated quality assurance, monitoring, and validation framework
- Technology: Delta Live Tables (Batch mode)
- Schema: imdb_data.qa
- Trigger: Executes after Bronze layer completes

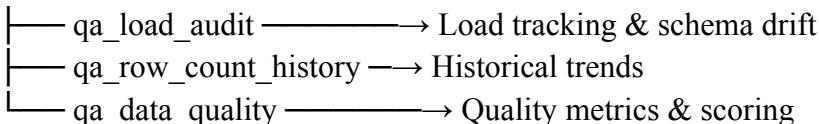
Architecture Overview

The QA layer provides automated data quality monitoring through three complementary tables:

Bronze Tables (7 tables)



QA Layer (Reads Bronze)



Monitoring & Alerts

Key Features:

- Dynamic configuration-driven processing
- Schema drift detection
- Automated quality scoring
- Historical trend tracking

- Duplicate detection
- Null value analysis
- No hardcoded table references

DASHBOARD

