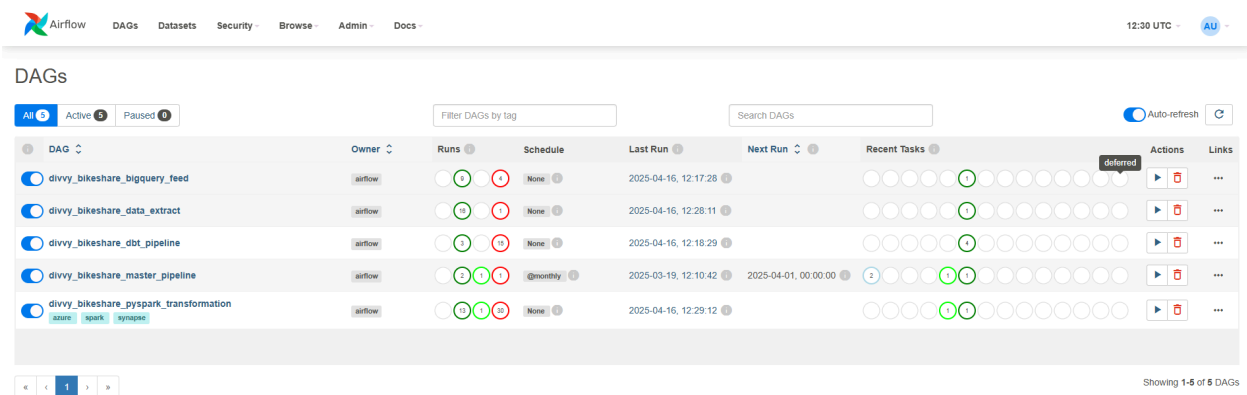# Divvy Bikeshare Analytical Project: Complete pipeline representation

We have five separate pipelines:
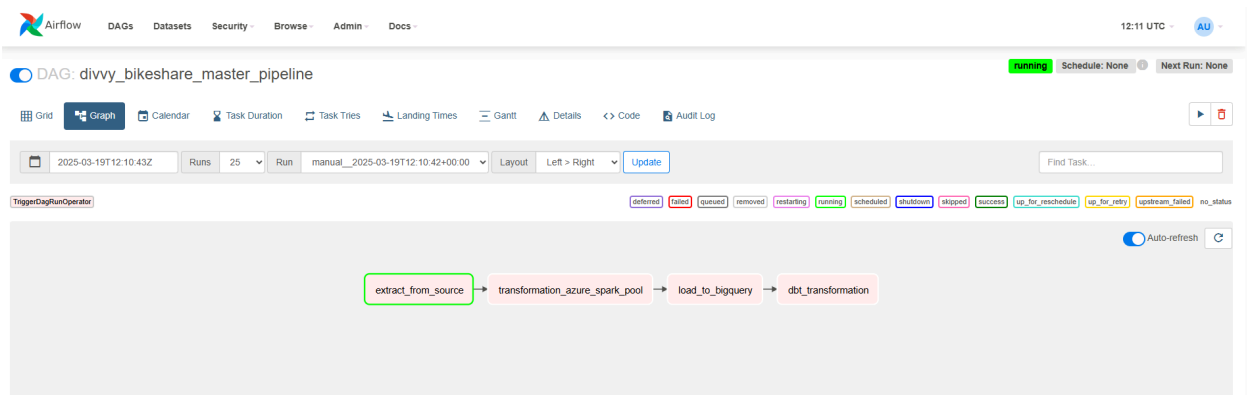
- divvy_bikeshare_data_extract
- divvy_bikeshare_pyspark_transformation
- divvy_bikeshare_bigquery_feed
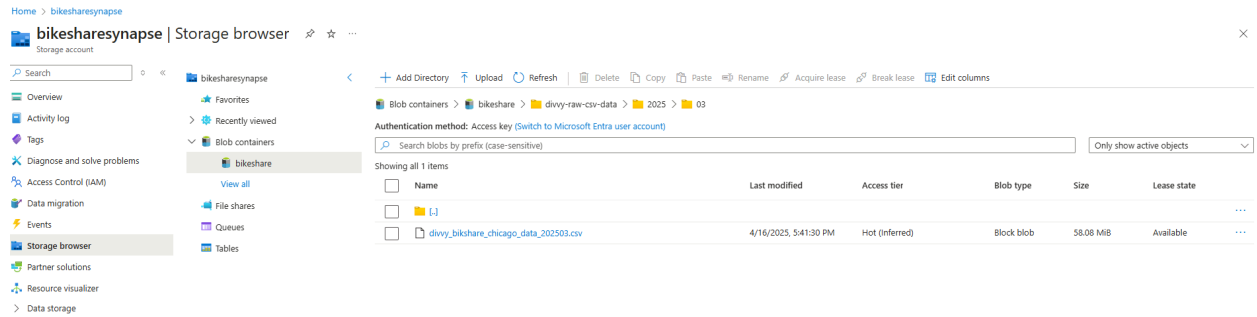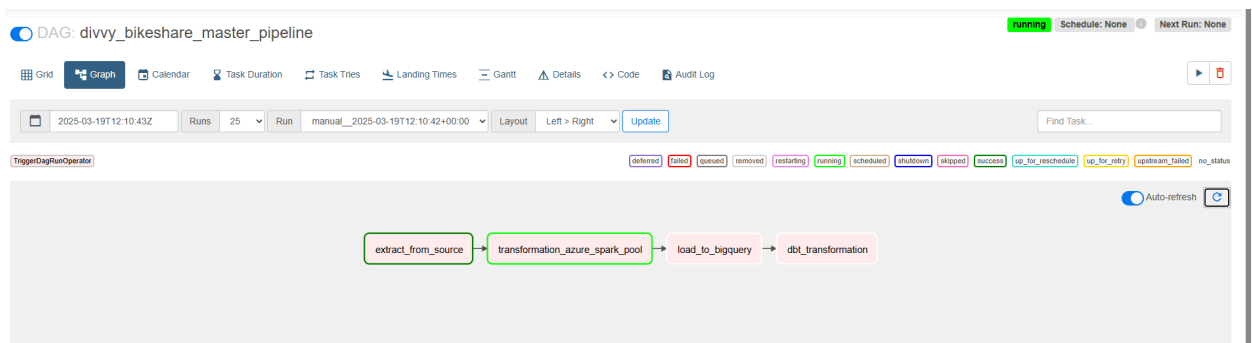- divvy_bikeshare_dbt_pipeline
- divvy_bikeshare_master_pipeline



The data pipeline is orchestrated by a master DAG  "divvy_bikeshare_master_pipeline" that coordinates the execution of several sub-DAGs in sequence.



divvy_bikeshare_data_extract : Downloads monthly Divvy Bikeshare data from S3, extracts ZIP files and standardizes naming  and uploads raw CSV data to Azure Blob Storage

divvy_bikeshare_pyspark_transformation is triggered by this and it submits a pyspark job in synapse analytics spark pool.



Executes PySpark jobs for data transformation

After the job is succeeded in synapse analytics, it saves data in Parquet format.



Divvy_bikeshare_bigquery_feed gets triggered it renames processed Parquet files with standardized naming and it loads transformed data into BigQuery

It triggers the final DAG "divvy_bikeshare_dbt_pipeline"

It runs dbt models with dynamic variables and applies business logic transformations and updates in bigquery.

It finally updates the dashboard hosted on looker: