

Breweries Case

VERSIONING

Versão	Data	Autor	Descrição
1.0	24/07/2024	Bhrenner Wilson Marcelino	Document Creation

Sumário

VERSIONING	2
Objective	4
Source.....	4
Solution.....	5
Orchestration	9
Monitoring	10
Alerting.....	10
How to execute	11

Objective

The goal of this test is consuming data from an API, transforming and persisting it into a data lake following the medallion architecture with three layers: raw data, curated data partitioned by location, and an analytical aggregated layer.

Source

The Source for this test is an api that can be consumed via the endpoint: api.openbrewerydb.org/breweries . In this api there is data about breweries of different types and different locations, the file is in json format and has the following schema:

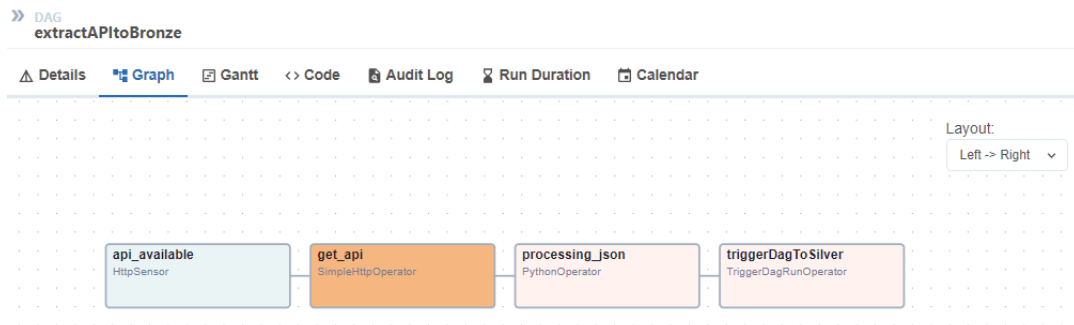
```
json_schema = {  
    "id": str,  
    "name": str,  
    "brewery_type": str,  
    "address_1": str,  
    "address_2": str,  
    "address_3": str,  
    "city": str,  
    "state_province": str,  
    "postal_code": str,  
    "country": str,  
    "longitude": str,  
    "latitude": str,  
    "phone": str,  
    "website_url": str,  
    "state": str,  
    "street": str  
}
```

Solution

To complete the case, the solution was developed in 3 different steps, each step is an airflow dag.

1. ExtracttoBronze

This step was designed to read data from the API using a connection and variables in Airflow, utilizing the HTTP sensor to validate and the HTTP operator to make requests. It transforms this data into a DataFrame using Pandas and writes it to the "bronze" layer of the data lake, preserving the original schema and JSON format properties.



1.1. Connection

Connections in Airflow are sets of configurations used to connect with other tools in the data ecosystem. This connection, called `base_api`, is used to validate and connect to the API. It works together with the variable `endpoint_api`, which is a parameter that defines the full path to access the data.

List Connection							
Search							
+ Actions -< Record Count: 1							
	Conn Id	Conn Type	Description	Host	Port	Is Encrypted	Is Extra Encrypted
<input type="checkbox"/>	base_api	http		https://api.openbrewerydb.org		False	False

1.2. Variable

Variables are a generic way to store and retrieve arbitrary content or settings as a simple key value store within Airflow.

1.2.1. End point

In this case, the variable is used to connect to the API to consume the data.

List Variable				
Search ▾				
<div> <div>+</div> <div>Actions ▾</div> <div>←</div> </div>				Record Count: 4
<input type="checkbox"/>	Key ▴ ▾	Val ▴ ▾	Description ▴ ▾	Is Encrypted ▴ ▾
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> endpoint_api	/breweries		False
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> path_bronze	datalake/bronze/brewe...		False
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> path_gold	datalake/gold/byType/...		False
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> path_silver	datalake/silver/breweri...		False

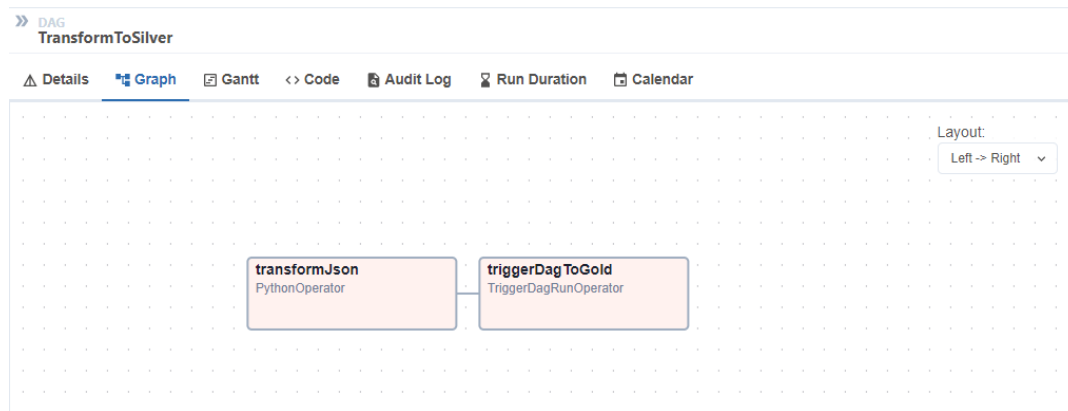
1.2.2. File path layer

For the next process to consume the generated file in the lake, a variable is set in Airflow using the same file path where the data was saved. This dynamically assists the process in consuming the file on the correct date.

List Variable				
Search ▾				
<div> <div>+</div> <div>Actions ▾</div> <div>←</div> </div>				Record Count: 4
<input type="checkbox"/>	Key ▴ ▾	Val ▴ ▾	Description ▴ ▾	Is Encrypted ▴ ▾
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> endpoint_api	/breweries		False
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> path_bronze	datalake/bronze/brewe...		False
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> path_gold	datalake/gold/byType/...		False
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> path_silver	datalake/silver/breweri...		False

2. TransformToSilver

After completing the previous step, the data is read from the bronze layer transformed into a dataframe by pandas, applying a fixed schema. For partitioning issues, a column called “process_date” was added indicating the date on which the data was processed. After these transformations, the dataframe is converted to parquet type and saved in the “silver” layer, partitioned respectively by: “process_date”, “country”, “state” and “city”.



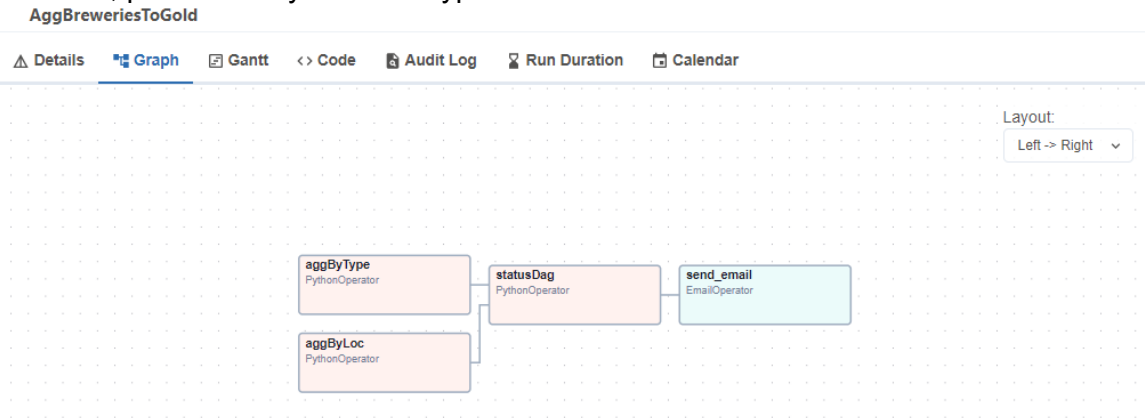
2.1. File path layer

For the next process to consume the generated file in the lake, a variable is set in Airflow using the same file path where the data was saved. This dynamically assists the process in consuming the file on the correct date.

List Variable				
Search				
+ Actions -				Record Count: 4
<input type="checkbox"/>	Key	Val	Description	Is Encrypted
<input type="checkbox"/>	endpoint_api	/breweries		False
<input type="checkbox"/>	path_bronze	datalake/bronze/brewe...		False
<input type="checkbox"/>	path_gold	datalake/gold/byType/...		False
<input type="checkbox"/>	path_silver	datalake/silver/brewen...		False

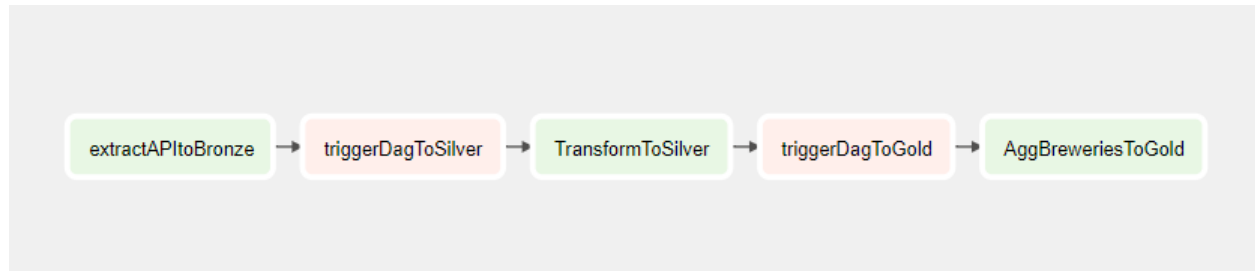
3. AggBreweriesToGold

In this last step, the data in parquet format is read from the silver layer by pandas and the view construction process begins. The first view called “agg_by_type” shows the number of breweries aggregated by type. The second is called “agg_by_location”, it shows the number of breweries aggregated by Country, state and city respectively. After the transformations, the data is saved in parquet format in the “gold” layer of the datalake, partitioned by date and type of view.



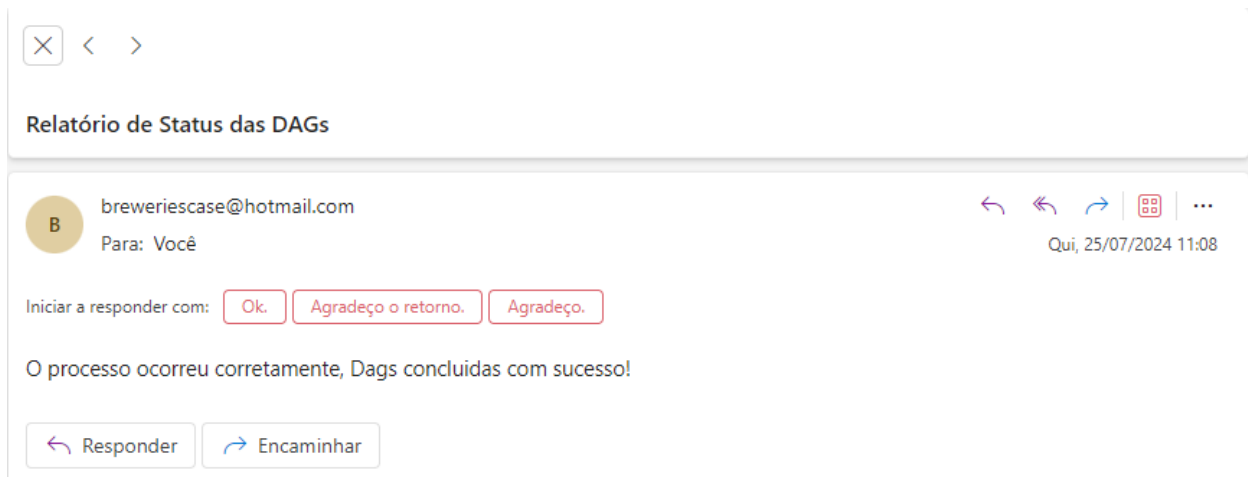
Orchestration

The process is scheduled to run once a day. The first process starts and when completed the second process starts and so on.



Monitoring

At the end of all Airflow processes, there is a task that checks the status of each DAG to log and send emails in case of failure or success.

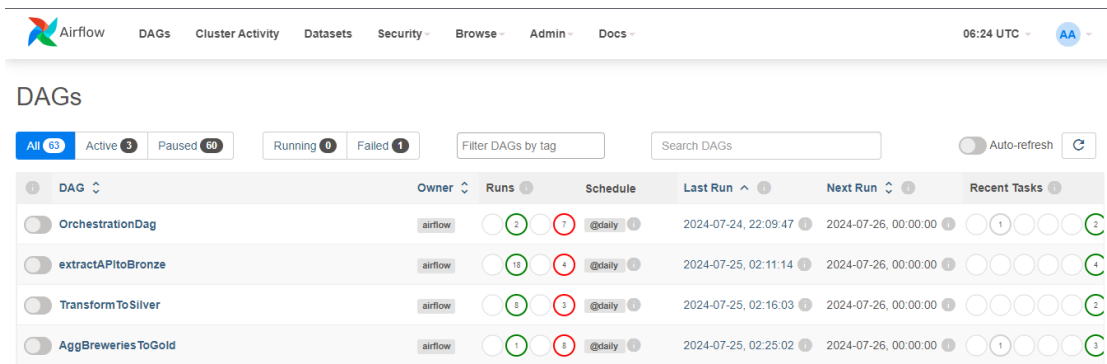


Alerting

In a future would be to make a report on days when there is a problem with the API and we were unable to read its data, when the day's folder in datalake has no records, it would send an alert to the developers.

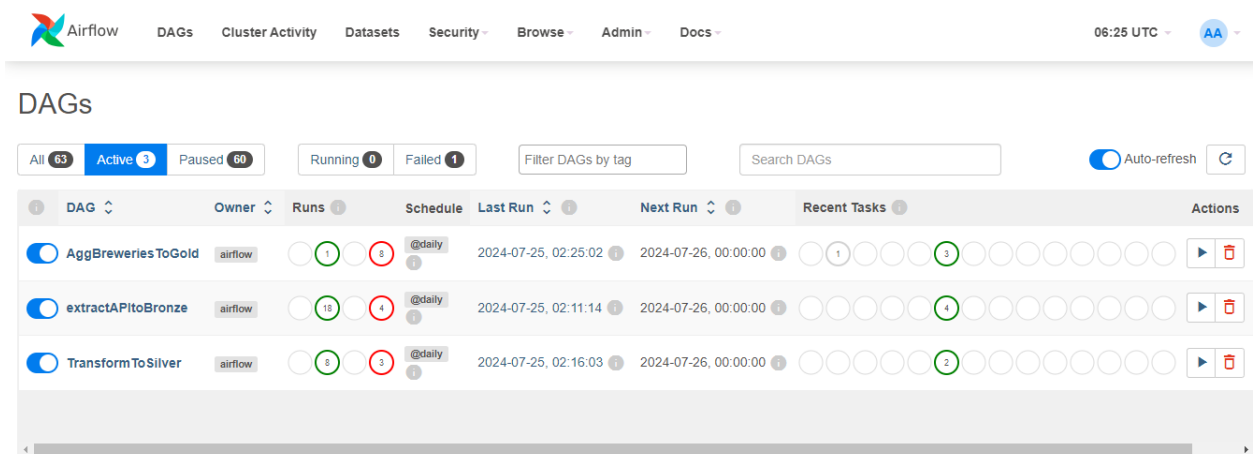
How to execute

1. First of all it's necessary complete the steps in: **docs/ Setup.docx**
2. After starting the environment, there will be a panel called “Dags”, you must enable the following dags: **ExtracttoBronze, TransformToSilver, AggBreweriesToGold.**



DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
OrchestrationDag	airflow	3	@daily	2024-07-24, 22:09:47	2024-07-26, 00:00:00	1 2
extractAPtoBronze	airflow	10	@daily	2024-07-25, 02:11:14	2024-07-26, 00:00:00	4
TransformToSilver	airflow	8	@daily	2024-07-25, 02:16:03	2024-07-26, 00:00:00	2
AggBreweriesToGold	airflow	1	@daily	2024-07-25, 02:25:02	2024-07-26, 00:00:00	1 3

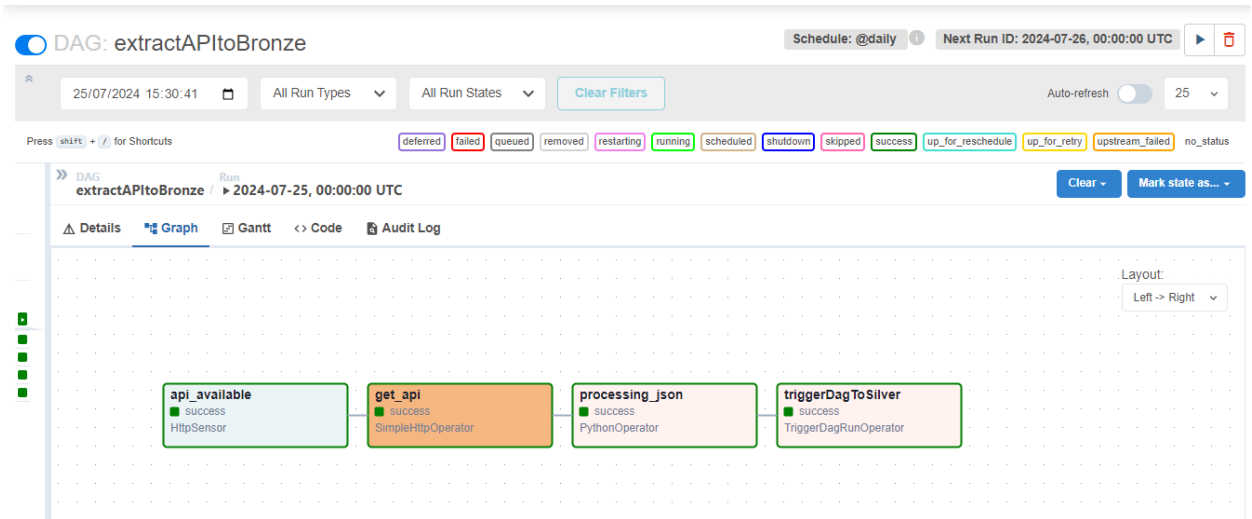
With **just these three dags** active, your “active” panel should look like this:



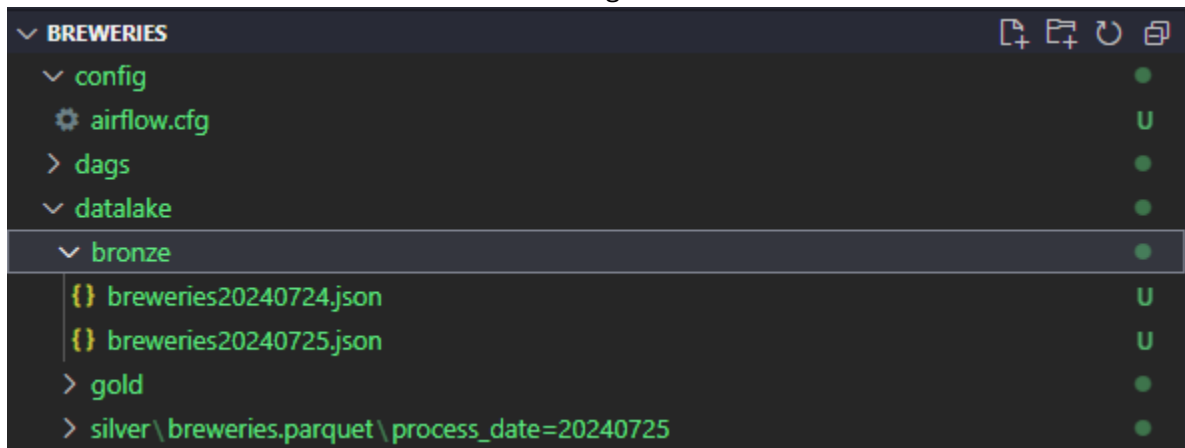
DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions
AggBreweriesToGold	airflow	1	@daily	2024-07-25, 02:25:02	2024-07-26, 00:00:00	1 3	▶ 🗑
extractAPtoBronze	airflow	10	@daily	2024-07-25, 02:11:14	2024-07-26, 00:00:00	4	▶ 🗑
TransformToSilver	airflow	8	@daily	2024-07-25, 02:16:03	2024-07-26, 00:00:00	2	▶ 🗑

3. After activation, you need to set the connections parameters and endpoint variables to connect with the API and the you can wait for the scheduled time in airflow or run the dag manually, just click on the **ExtracttoBronze** Run icon. This dag when executed will automatically execute the next dagPartition.

Clicking on **ExtracttoBronze**, you will be able to view a panel with the status of each task, as shown below:



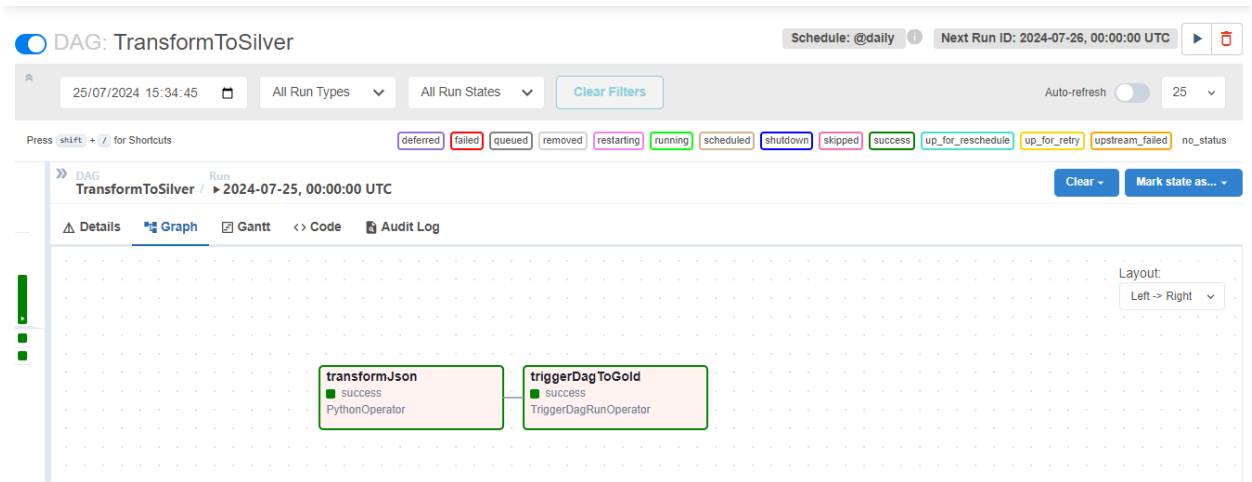
The result in the datalake should be the following in case of success:



- After the first step is completed, just check if the other dags executed automatically after the first.

The results should be as follows:

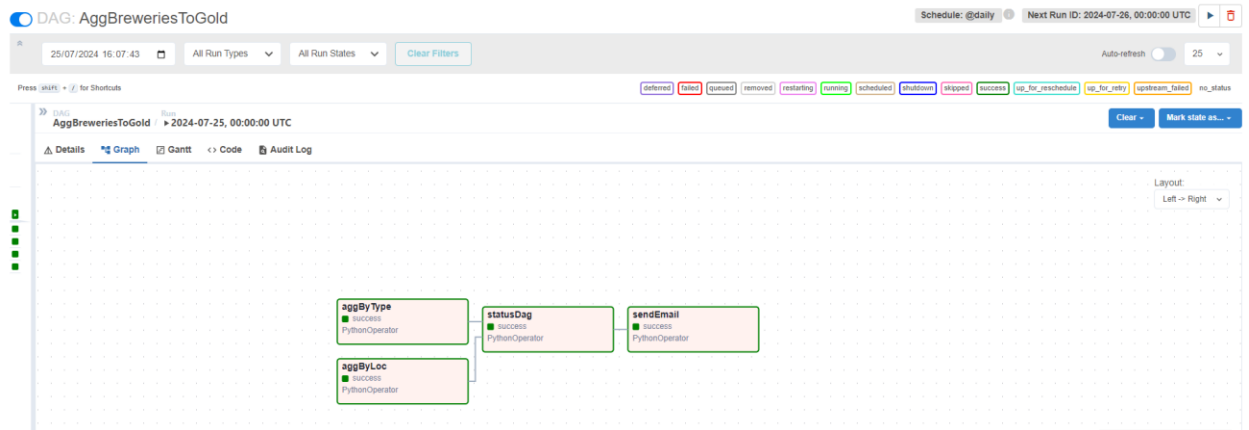
Success on **TransformToSilver**



On datalake

BREWERIES		
> dags		●
> datalake		●
> bronze		●
> gold		●
> silver \ breweries.parquet \ process_date=20240725		●
> country=Ireland \ state=Laois \ city=Killleshin		●
≡ 2f30fd43dc264b1ea35735a1d32a0a48-0.parquet		U
≡ 0970c6ccc1574f6eb6f70f13ad9bf77d-0.parquet		U
> country=United%20States		●
> state=Arizona		●
> city=Gilbert		●
≡ 2f30fd43dc264b1ea35735a1d32a0a48-0.parquet		U
≡ 0970c6ccc1574f6eb6f70f13ad9bf77d-0.parquet		U
> city=Mesa		
> city=Tucson		
> state=California		●
> city=Mariposa		●
> city=Petaluma		
> city=San%20Diego		
> city=Westlake%20Village		
> state=Colorado		
> city=Castle%20Rock		
> city=Denver		
> city=Louisville		
> state=Delaware		

Sucess on AggBreweriesToGold:



On datalake:

BREWERIES			
config			
airflow.cfg			U
dags			
datalake			
bronze			
gold			
byLocation \ view_breweries_by_location.parquet \ process_date=20240725			
bb0cbea2a433485d9beb4b2ebfef861b-0.parquet			U
byType \ view_breweries_by_type.parquet \ process_date=20240725			
f32045315ce4434e95065b9c302abe1e-0.parquet			U
silver \ breweries.parquet \ process_date=20240725			