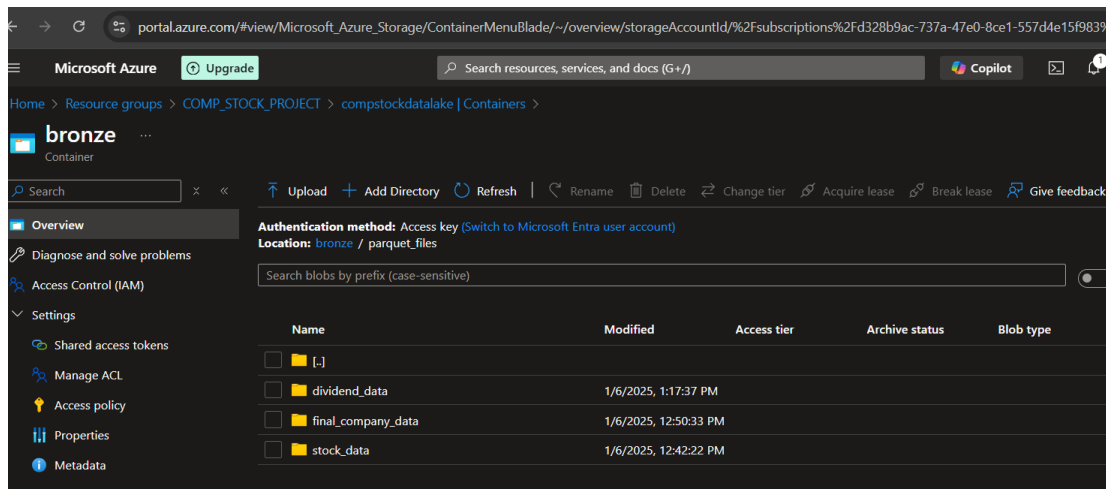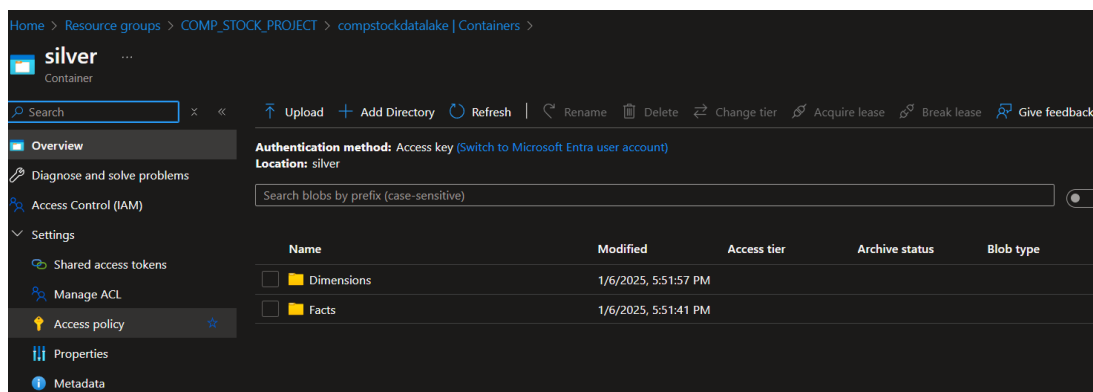# Financial Data Engineering in Azure- Guide
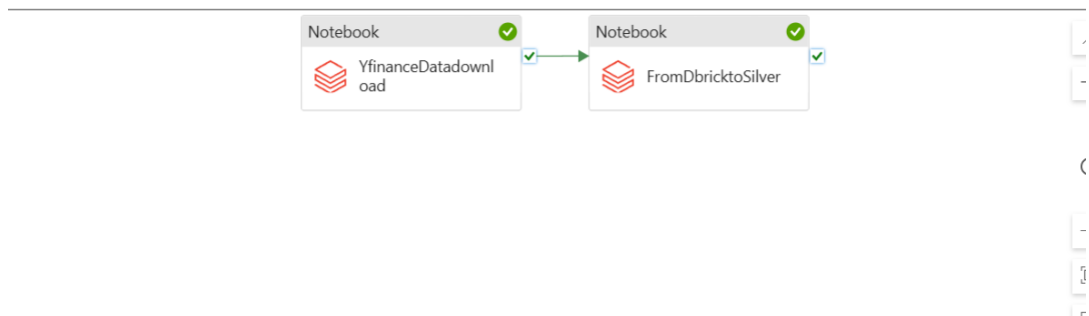
## Data extract and transformation

Initially, I intended to use an Azure Function to deploy my data fetching code, which retrieves data from yfinance. However, yfinance implemented a limitation on company data requests, so I had to introduce a time delay in the code, waiting one second between each company data retrieval. As a result, I exceeded the 5-minute limit allowed in the Azure Function package. Therefore, I switched to using a Databricks notebook to run the Python script instead. I created a service principal through app registration, which I assigned to my Data Lake to load into bronze container



Then, I performed transformations using a Databricks notebook, and loaded the results into the Silver container.

**Pipeline**



**Azure Synapse**

*Synapse Managed Identity*: Ensure the **Storage Blob Data Contributor** role is assigned to the Synapse workspace's managed identity in the storage account.

*User Identity:* Assign the **Storage Blob Data Contributor** role to your user account in the storage account.

**Why I Chose Synapse Serverless SQL**

I opted for the serverless SQL solution in Azure Synapse because it is significantly more cost-effective than storing data directly in Synapse's dedicated SQL pools. This approach allows me to query data stored in Azure Data Lake using an abstraction layer. The serverless SQL layer acts as a logical interface, enabling on-demand querying of files (e.g., Parquet, CSV) without the need for pre-loading or duplicating the data. This flexibility reduces storage and compute costs while maintaining scalability for ad-hoc analytics.
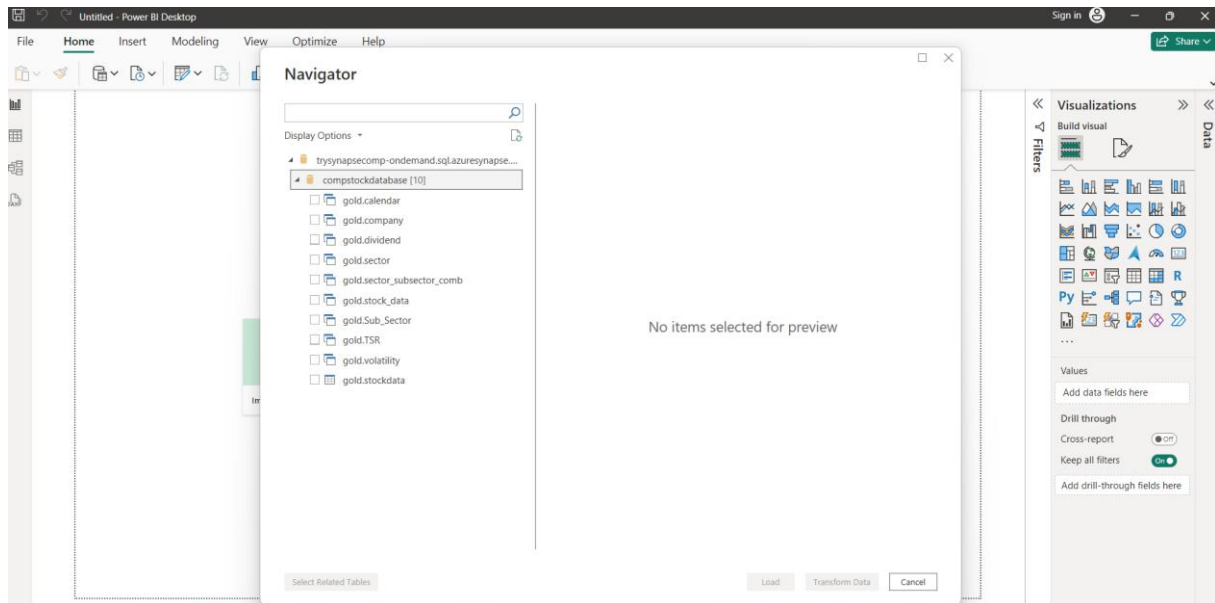
In **Synapse**, you created **views** and **external tables** to query the data directly from the Data Lake.

**External tables** allow you to query data stored in your Data Lake without having to move it into a database, providing efficient and scalable querying.
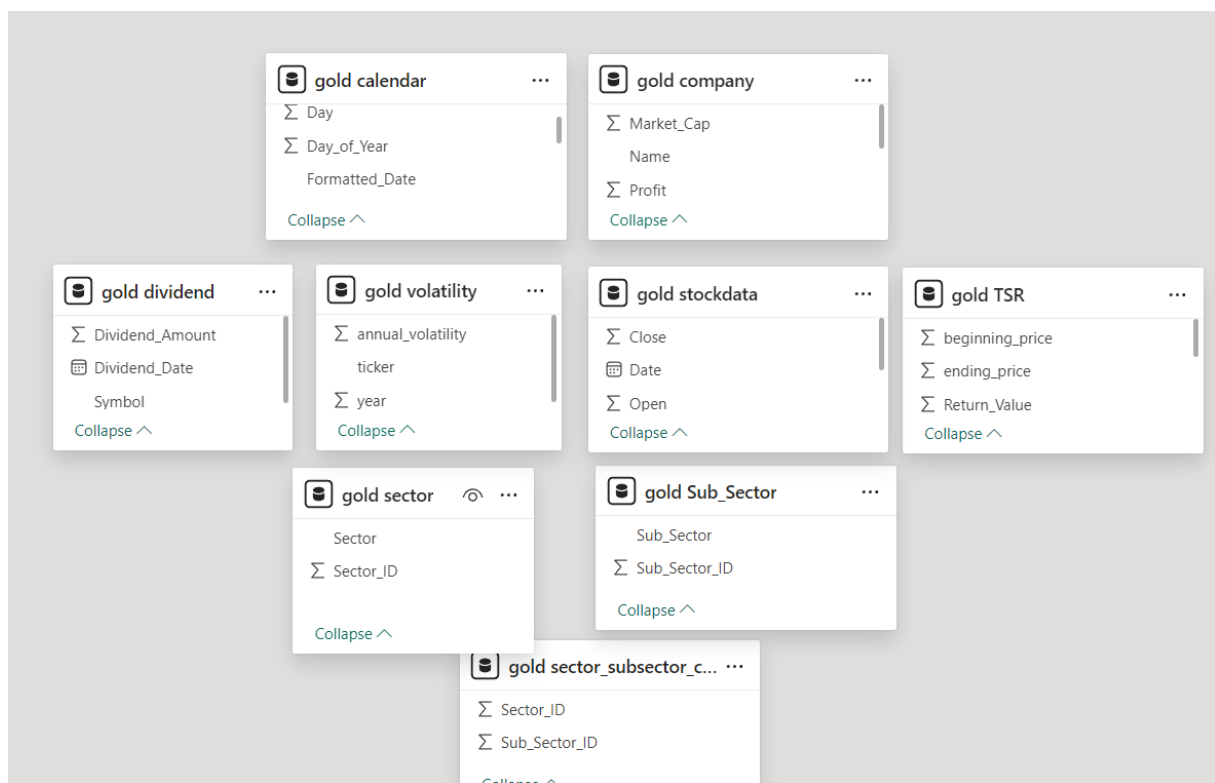
**Views** were created on top of these external tables to simplify access and provide a structured way to query the transformed data.

**Loading**

1. After creating the schema in Synapse, I proceeded with loading the the schema into Power BI desktop using SQL serverless endpoints.



Here is the schema:

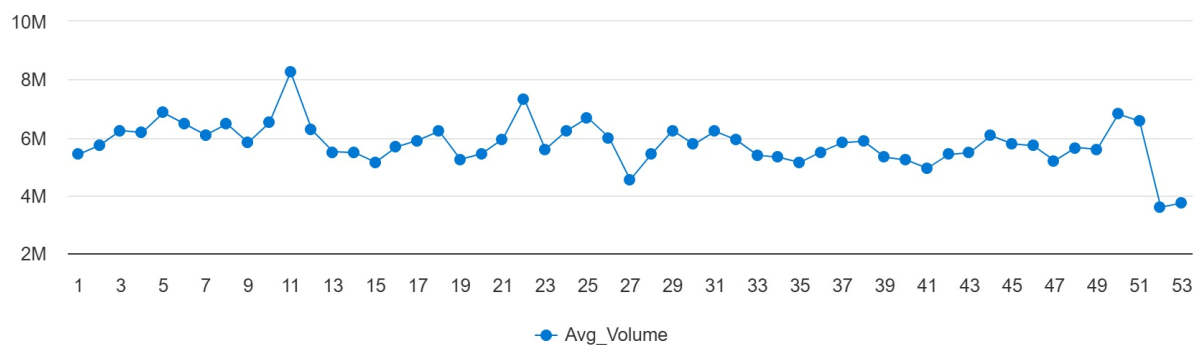2. I also got it loaded into the gold container



## Analytics

```sql
SELECT
    DATEPART(WEEK, Date) AS Week_Number,
    AVG([Close]) AS Avg_Close,
    AVG(Volume) AS Avg_Volume
FROM gold.stock_data_transf
GROUP BY DATEPART(WEEK, Date)
ORDER BY Week_Number;
```



```sql
--- Top 5 companies with the largest TSR
SELECT TOP 5
    ticker,
    year,
    TSR
```

```
FROM gold.TSR_annual
ORDER BY TSR DESC;
```



```
---TOP 5 companies with the best yearly TSR change --------
WITH TSR_Changes AS (
    SELECT
        ticker,
        year,
        TSR,
        LAG(TSR) OVER (PARTITION BY ticker ORDER BY year) AS Previous_TSR
    FROM gold.TSR_annual
)
SELECT TOP 5
    ticker,
    year,
    TSR,
    Previous_TSR,
    (TSR - Previous_TSR) AS TSR_Change
FROM TSR_Changes
WHERE Previous_TSR IS NOT NULL
ORDER BY TSR_Change DESC;
```