# Final Project Documentation

## Data Analyst Internship 2024

Stefan Ristikj

This document contains the overview and detailed documentation of the final project for the Data Analyst internship. The goal of the project is to assess my knowledge gained through the past two months, and it encompasses the Microsoft-provided database WideWorldImporters and its detailed analysis. Firstly, the database needs to be thoroughly checked for any inconsistencies, missing values and continuity in SQL Server Management Studio, then according to my area of analysis, I need to take whatever data I think is necessary and upload it to Azure services, whose environment I need to setup before that, then create views over the data that I finally import into PowerBI and create a report for a fictitious customer using various visualizations to make the report story-like and appealing for non-professional BI users.

First off, I checked the data for any anomalies. After making sure none were left (there weren't any significant changes I needed to make), I started preparing the Azure environment for importing the data. I started by creating an Azure resource group (figure 1) to keep all created resources in a centralized environment.
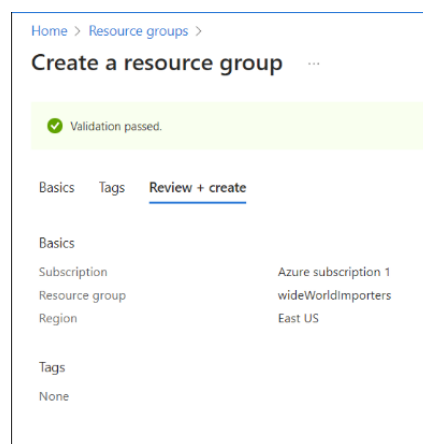


*Figure 1*

Then I proceeded to create a Synapse workspace (figure 2) which during initialization created a built-in serverless SQL pool automatically, to be able to work on the project. During the process of creating a workspace, I was prompted to create an Azure Data Lake Storage Gen2 account a file system name and credentials (a username, which I left to the default value of sqladminuser, and a password which I would use later to connect to my local machine).

*Figure 2*

After that, the next few steps were to open Synapse Studio to work in the workspace I created previously, then go to the Manage tab and Linked Services, where besides the two automatically created services to the Synapse Analytics workspace and the ADLS Gen2 account, I needed to create a linked service to the SQL server which hosts the WideWorldImporters database, so I did as in figure 3.



*Figure 3*

Next up, I needed to create a SQL database on the serverless pool, to host the external tables I will create in the next steps (figure 4).
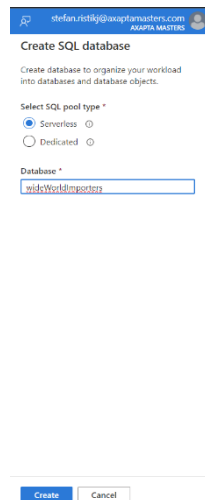
*Figure 4*

After that, I decided to create a preliminary schema (figure 5) of the two main schemas in the database part of my project, Warehouse and Sales. This schema contains every table from both schemas plus Application.People because it was needed, contains every column from the tables and my assumption on whether the tables and columns will be necessary for my final model and the fact table.
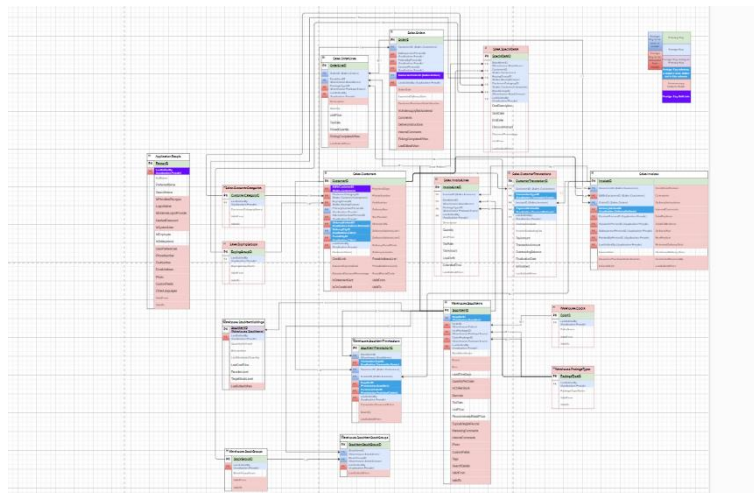


*Figure 5*

The next step is to create a fact table and some views on the data, which would eliminate the need for importing too many tables. In figures 6, 7 and 8 are the SSMS queries for creating views on the Sales.Customers and the Warehouse.StockItems tables, and for creating the fact table. For the Sales.Customers table, I needed information for the customers' location state-wise, so I joined the Cities and StateProvinces tables from the Application schema instead of taking them as whole tables (I also excluded a column with type geography due to a problem with these columns in Synapse). Same for the StockItems table, and I created the fact table as a joint table between the two tables that I previously identified as tables that contain facts, Sales.Invoices and Warehouse.StockItemTransactions.

*Figure 6*



*Figure 7*



*Figure 8*

Next up, I imported the necessary data in Synapse through the Ingest option on the home screen (figure 9). The option is a copy data tool which creates a pipeline automatically, recognizes the tables from the database through the linked service to SQL Server, creates source and destination integration datasets and executes the pipeline. The other options pertaining to the project are file format (parquet, as needed) and the location in which these parquet files should reside. After verifying everything is in order, I executed the pipeline, creating parquet files for every necessary table.



*Figure 9*

Following this, I proceeded to create the external data source and file format, which will then be used to create the external tables. In figure 10, I give the code for creating the data source and file format, as well as an external table on the ApplicationPeople.parquet file. Every table is created using the same pattern, selecting the columns I need and using the already created data source and file format.
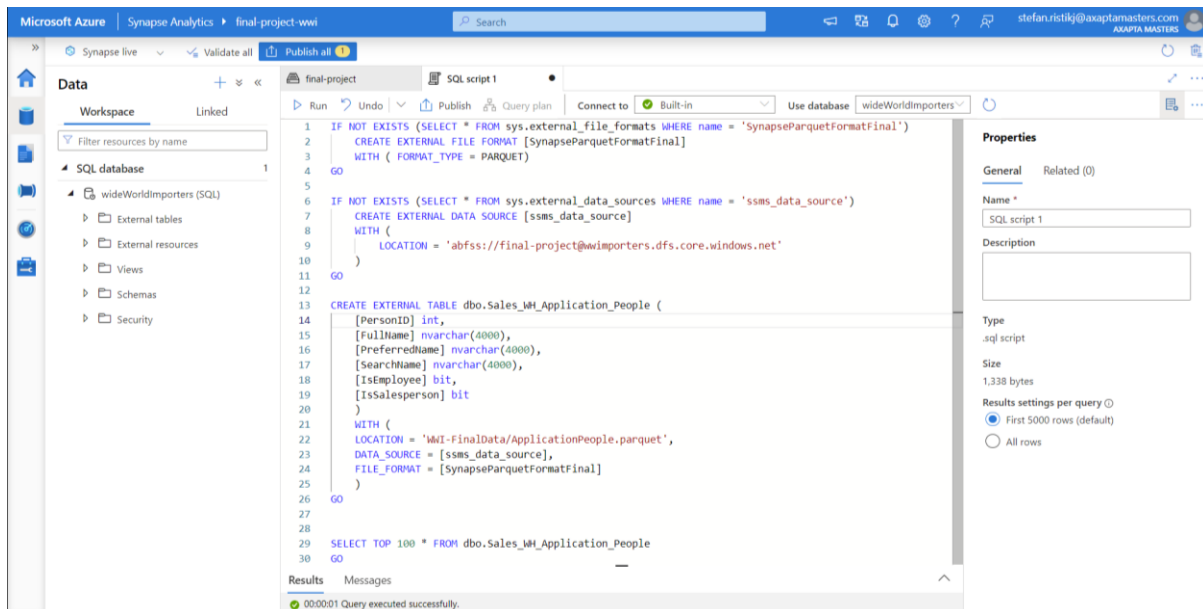


Figure 10

Then, I imported the external tables into PowerBI using a connection to Azure Synapse Analytics (figure 11) and reviewed the tables to make sure everything is in order, as shown in figure 12.
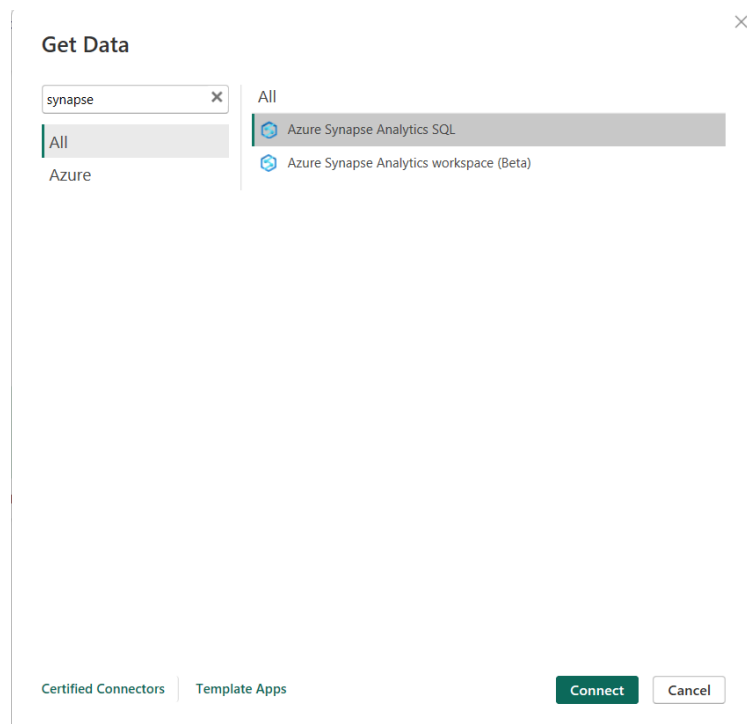


Figure 11

*Figure 12*

Last but not least, I reevaluated the schema and created a final one (figure 13) that most closely represents the data model I used to create the report, as shown in my presentation.
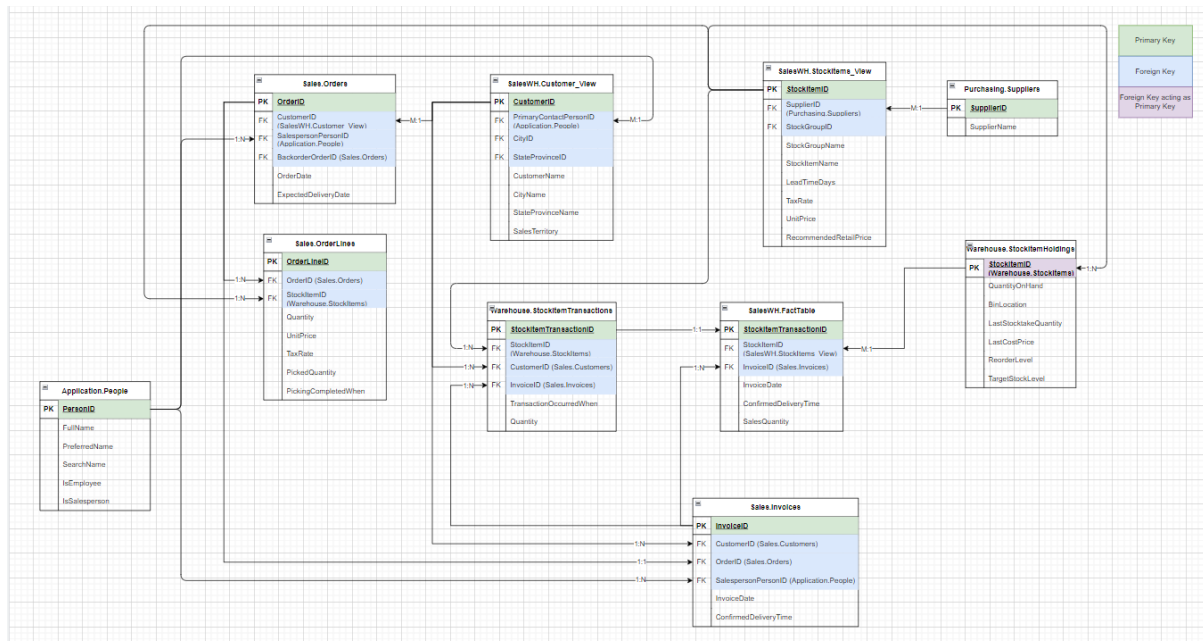


*Figure 13*