# Lab 05: Transferring large data volumes with Fabric Data Factory

## Introduction

In this lab, you will ingest data from a Azure SQL Database to the Managed Lakehouse Sales table. The objective is to understand how to connect to this external database and how the Fabric Data Factory Pipeline can scale when having to transfer 10M, 50M and 100M of rows.

## Objectives

After completing this lab, you will be better able to:

1. Connect to an Azure SQL Database and read the tables

2. Set up and run the Fabric Data Factory pipeline to transfer data from Azure SQL tables to a table in the managed Lakehouse.

## Estimated time to complete this lab

60 minutes

# Contents

## Lab Prerequisites

- Workspace: Fabric, Power Premium or Fabric trial
- Individual license: Power Pro or Premium Per User account
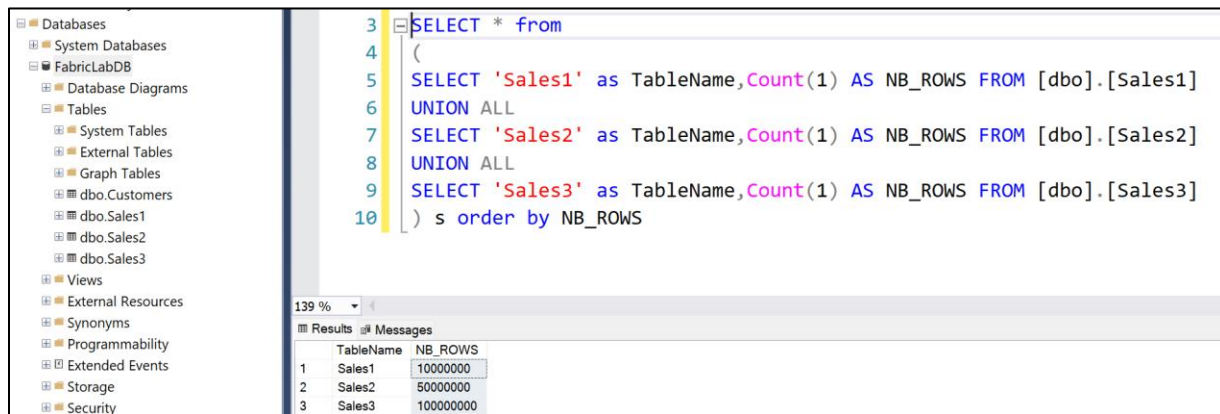
## Information provided by your training provider

- Trial tenant (if applicable): login & password, workspace to use for the lab.
- Azure SQL Server Database: server name, login and password.

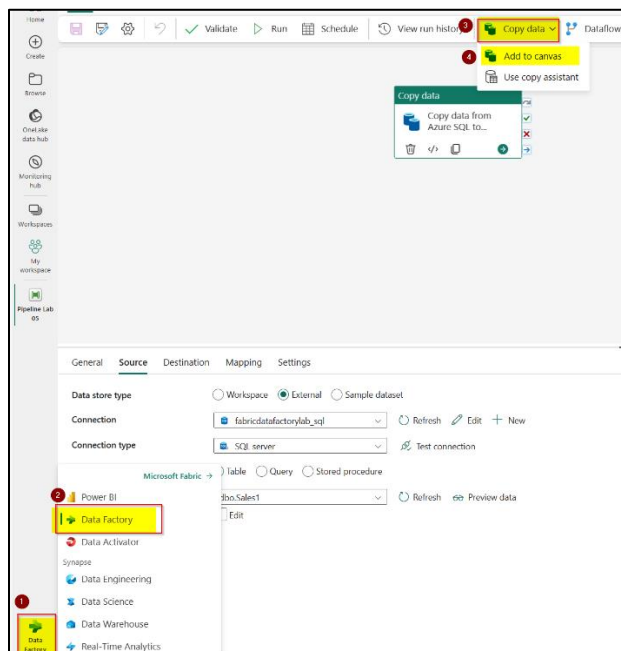# Task 1: Setup the connection to the Azure SQL Database

In this task, you will establish the Linked Service connection to an Azure SQL Database.

This database named FabricLabDB contains 3 tables with the same structure than the managed Sales table in the Lakehouse :

- Sales1 contains 10M rows
- Sales2 contains 50M rows
- Sales3 contains 100M rows



- On the MS Fabric menu, navigate to Data Factory and click on the Data pipeline. Provide the name Pipeline Lab 05 for this new Lab.
- On the newly created data pipeline, click on "Copy Data" then "Add to canvas"

- Configure the Data Copy Activity Task
    - Under the **General Tab**, provide the name for the Pipeline: ***Copy data from Azure SQL to Lakehouse***. The other settings remain unchanged.
    - 
- On the **Source Tab** choose the following options:
    - Data store type: External
    - Connection: Click on new and select Azure SQL Server, click on Continue.

    **Server**: *<database server name>*.database.windows.net (information given by your trainer.

    **Database**: fabriclabdb

    **Connection**: Create new connection

    **Connection name:** azuresqldb

    **Authentication** kind: Basic

    **Username**: LabUser (to be confirmed by your trainer)

    **Password**: MicrosoftFabric@2023 (to be confirmed by your trainer)

    Thick Use **Encrypted Connection**

Make sure that the settings will look like this:

o Click on Test Connection to make sure that the connection is successful.



o Use query: Select Table
• **Table optional**: To find out how many records are on the dbo.Sales1, as a test you can select Query and paste this query:

```sql
SELECT COUNT(*) AS TotalRecords
FROM [dbo].[Sales1];
```



o As you can see, the table dbo.Sales1 contains 10M of rows.
o **Table**: select **dbo.Sales1**. Click on Preview Data to make sure that you can read the data from the Database

| General | **Source** | Destination | Mapping | Settings |
|---|---|---|---|---|

**Data store type**      ○ Workspace   ● External   ○ Sample dataset

**Connection**      azuresqldb User2   ⌄      ↻ Refresh   ✎ Edit   + New

**Connection type**      SQL server   ⌄      ✑ Test connection
                                              ✔ Connection successful

**Use query**      ● Table   ○ Query   ○ Stored procedure

**Table**      dbo.Sales1   ⌄      ↻ Refresh   👓 Preview data
              ☐ Edit

> Advanced

**Preview data**

| ⊞ | OnlineSalesKey | OrderDate | DeliveryDate | StoreKey | ProductKey | PromotionKey | CustomerKey | SalesOrderNumber | SalesOrd |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 19562610 | 2020-01-01T00:00:00 | 2020-01-04T00:00:00 | 307 | 293 | 5 | 19079 | 200701012CS425 | 2 |
| 2 | 19562623 | 2020-01-01T00:00:00 | 2020-01-06T00:00:00 | 307 | 672 | 5 | 19079 | 200701013CS425 | 2 |
| 3 | 19562636 | 2020-01-01T00:00:00 | 2020-01-03T00:00:00 | 307 | 1420 | 5 | 19079 | 200701015CS425 | 2 |
| 4 | 19562655 | 2020-01-01T00:00:00 | 2020-01-02T00:00:00 | 307 | 2167 | 5 | 19079 | 200701018CS425 | 2 |
| 5 | 19562863 | 2020-01-01T00:00:00 | 2020-01-10T00:00:00 | 307 | 47 | 5 | 19079 | 200701011CS425 | 2 |
| 6 | 19562869 | 2020-01-01T00:00:00 | 2020-01-06T00:00:00 | 307 | 1008 | 5 | 19079 | 200701014CS425 | 2 |
| 7 | 19563086 | 2020-01-01T00:00:00 | 2020-01-07T00:00:00 | 307 | 1688 | 5 | 18347 | 20070101729346 | 2 |
| 8 | 19563087 | 2020-01-01T00:00:00 | 2020-01-07T00:00:00 | 307 | 1688 | 5 | 2 | 200701017111001 | 2 |
| 9 | 19563088 | 2020-01-01T00:00:00 | 2020-01-07T00:00:00 | 307 | 1688 | 5 | 3 | 200701017111002 | 2 |
| 10 | 19563089 | 2020-01-01T00:00:00 | 2020-01-07T00:00:00 | 307 | 1688 | 5 | 4 | 200701017111003 | 2 |

- Under the **Destination** tab, select the following options:
  - Data store type: Workspace
  - Workspace data store type: Lakehouse
  - Lakehouse: Contoso (your managed Lakehouse created previously)
  - Root folder: Tables
  - Table name: Sales
  - Make sure that the settings will look like this:

- On the **Mapping** tab select Import schemas:



- On the **Settings** tab, no changes are required
- Validate and Save the Pipeline

- You can now **Run** your pipeline

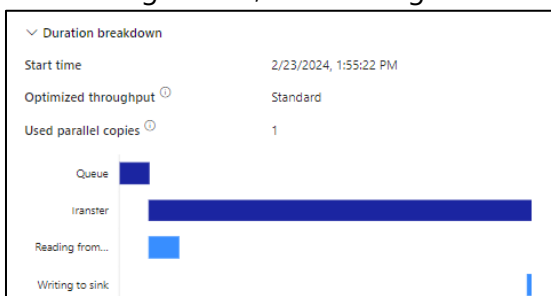- During the pipeline execution, you can monitor the progression by clicking on the Activity in the Output section



- Details should appear, and 10M rows should we be written.



- Note the time spent to execute the pipeline to transfer 10M of rows



- As demonstrated in the lab 3, you can use the Duration Breakdown to understand where the execution time is mostly consumed – depending on the status of the Azure SQL DB used during the lab, the Reading from source (Light Blue) might be bigger.

- Using the Lakehouse explorer, you can also see the last Parquet file written during the last pipeline execution. This parquet file weights 50MB and contains 10M of records.



- It is also recommended to use the available analytics tools in Fabric to see the new ingested in the Sales Tables
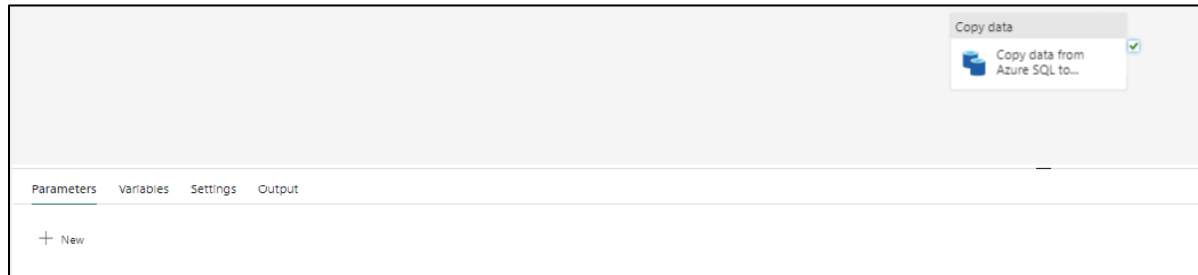  - **The SQL Analytics Endpoint:**



  - **The Power BI report**
    - You can add a simple visual card to count the number of rows in the Sales tables, and you don't have to refresh the report or the semantic model, thanks to the **Direct Lakehouse Mode**
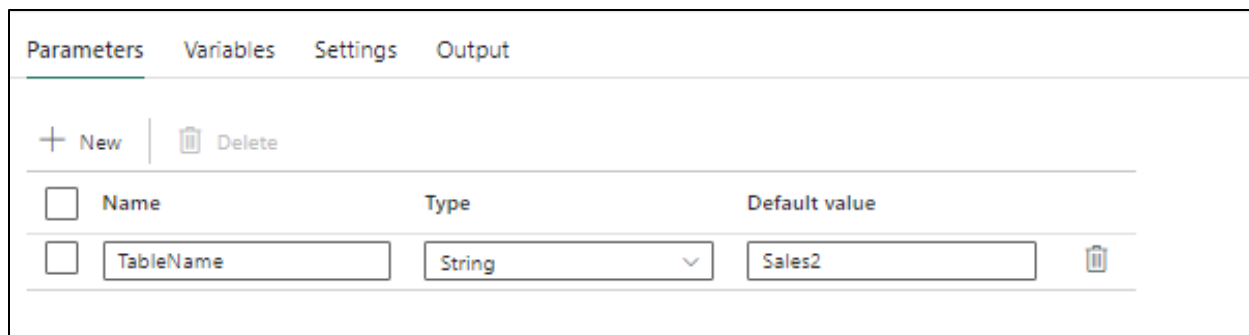
# Task 2: Use Data Pipeline Parameters to change the source table

In this task, you will declare a Parameter to be used in the Pipeline. Then you will load data from the second table (50M of records)

- Open the Fabric Data Factory Pipeline, and go to the **Parameters** section



- Define the following parameter:
  - Name : **TableName**
  - Type : **String**
  - DefaultValue : **Sales2**

- Select the **Copy Data Activity** in the Pipeline, go to the **Source**, and enable the **Enter manually** check box
- Click on the **Add dynamic content** link to access to the Pipeline expression builder



- From the Pipeline Expression Builder, drag and drop the TableName parameter, and the following expression should appear. Click on OK.



- The expression should be now visible in the table name section.



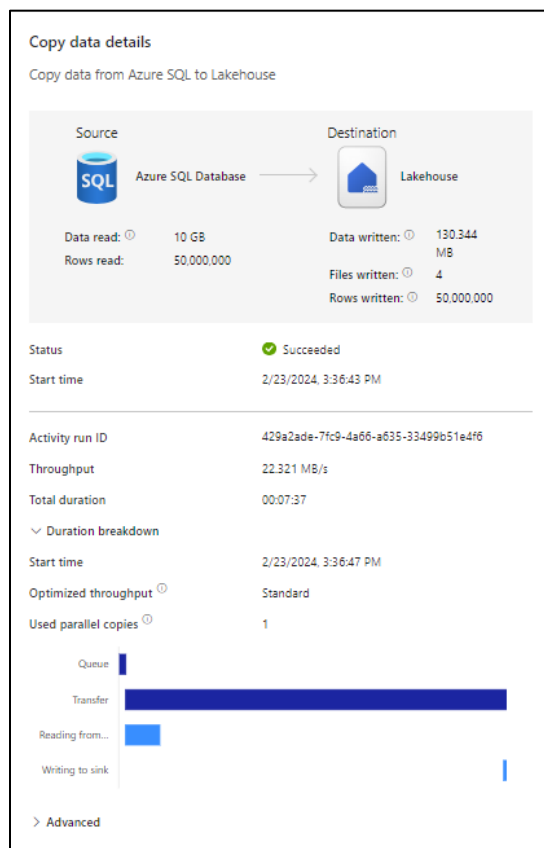- The expression should be now visible in the table name section.

- Save the pipeline and click on Run.
- As the Pipeline now contains a parameter, you will have the choice to keep the default parameter value, or change it before the execution. Keep the default parameter value for this execution to get data from the table Sales2.
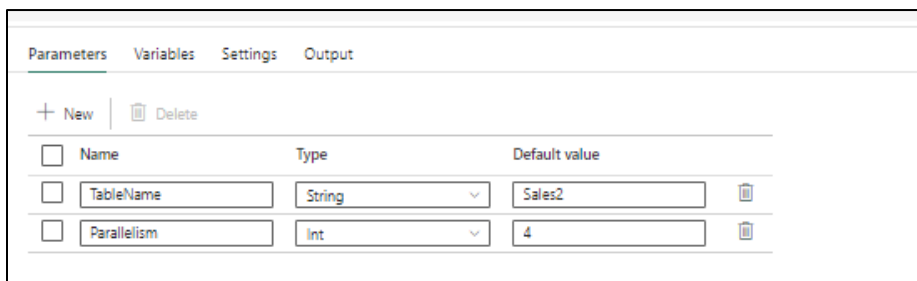


- At the end of the pipeline execution, you can observe that :
  - 50M of records were loaded
  - 4 files were written
  - The parallelism was 1 (in the next lab you will configure the parallelism using a parameter).

# Task 3: Optimize Data Pipeline performance with parallelism and partitions

In this task, you will declare a second Parameter to be used as the level of parallelism in order to accelerate the data transfer and optimize performance by reducing the execution duration.
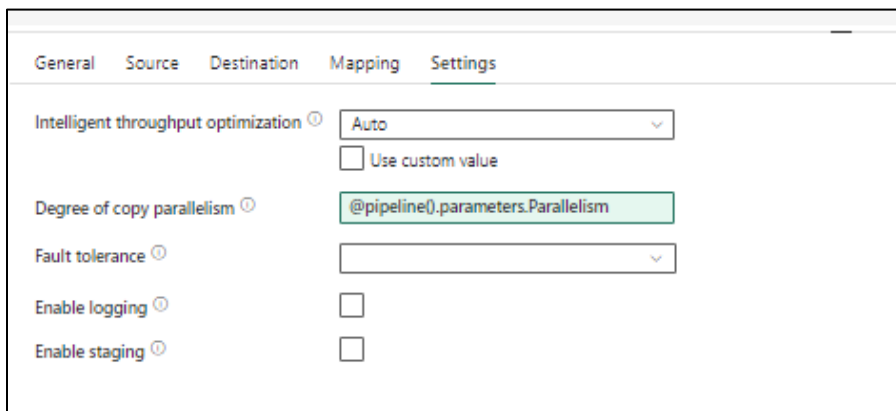
- Create a second parameter with the following settings:
  - Name: Parallelism
  - Type: Int
  - Default Value: 4



- Use this parameter to dynamically configure the Degree of copy parallelism (in the settings section of the Copy activity)

- On the Source section, enable the partitioning by selecting **Physical partitions** of table in the **Partition option** item.



- Save the pipeline and run it with the following configuration:
  - TableName : Sales3 (100M of rows)
  - Parralelism : 4

- At the end of the pipeline execution you can observe :
  - A better throughput when transferring data
  - Good performance (only 2mn more when moving from 50M to 100M of rows)