Andy Cutler

**Independent Consultant & Contractor**

**Azure Data Platform & Power BI**

www.datahai.co.uk

https://twitter.com/MrAndyCutler

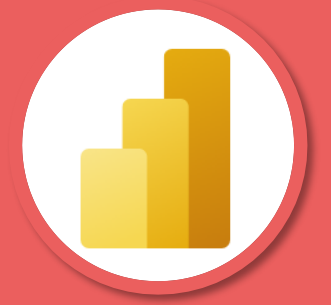https://www.linkedin.com/in/andycutler/

# Session Overview

- What are **Power BI Dataflows**?

- What is **Azure Synapse Analytics Serverless SQL Pools**?

- Why use **Serverless SQL Pools with Dataflows**?

- Serverless SQL Pools and Power BI Dataflows **Together**

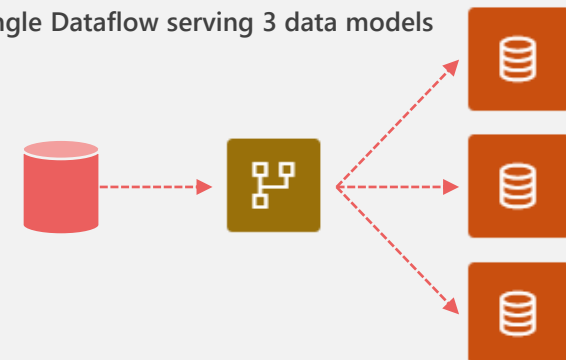- Considerations

# Power BI Dataflows

A **Power Query** based feature in the **Power BI Service** which enables:

- **Connecting** to a variety of data sources including SQL Databases and Data Lake Storage

- **Cleansing** and **Transforming** the data to suit requirements

- **Mapping** to common business entities using the **Common Data Model**

- Creating a **centralised repository** of data ready for using in Power BI Datasets



Single Dataflow serving 3 data models

# Azure Synapse Analytics Serverless SQL Pools

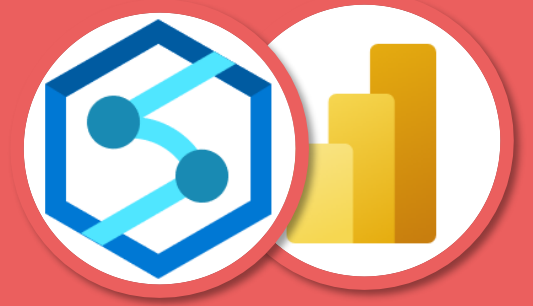**Cloud Analytics platform to process data in data lake storage**

- Ability to query file data "in place" using T-SQL without copying data to internal storage

- Ability to write data to external data lake storage

- Connect Business Intelligence tools using SQL endpoint

```sql
SELECT * FROM
OPENROWSET
(
    BULK 'conformed/factsalesorderheader/*/*/*/*.*',
    DATA_SOURCE = 'ExternalDataSourceDataWarehouse',
    FORMAT = 'parquet'
) as fctsl
```

```sql
SELECT * FROM
OPENROWSET
(
    BULK 'conformed/factsalesorderheader/*/*/*/*.*',
    DATA_SOURCE = 'ExternalDataSourceDataWarehouse',
    FORMAT = 'parquet'
) as fctsl
WHERE fctsl.filepath(1) = 2020
    AND fctsl.filepath(2) = 7
    AND fctsl.filepath(3) = 6
```
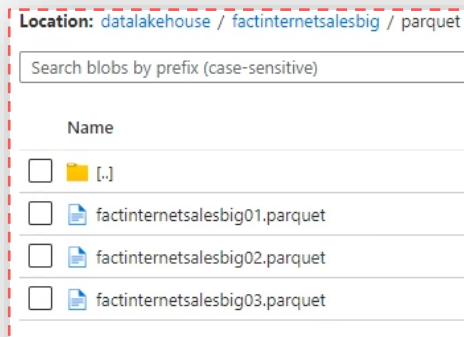
# Serverless SQL Pools and Dataflows Together

**Power BI can connect to the Serverless SQL endpoint just like any SQL database**

Dedicated SQL endpoint      : synapsedemodh.sql.azuresynapse.net
Serverless SQL endpoint     : synapsedemodh-ondemand.sql.azuresynapse.net
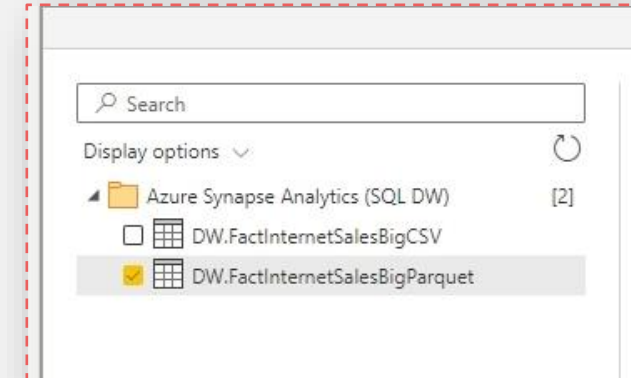Development endpoint         : https://synapsedemodh.dev.azuresynapse.net

**Source data in Data Lake**

Location: datalakehouse / factinternetsalesbig / parquet

Search blobs by prefix (case-sensitive)

Name

☐ 📁 [..]

☐ 📄 factinternetsalesbig01.parquet

☐ 📄 factinternetsalesbig02.parquet

☐ 📄 factinternetsalesbig03.parquet

**Create table to read data from Data Lake**

```
CREATE EXTERNAL TABLE DW.FactInternetSalesBigParquet (
    ProductKey INT,
    OrderDateKey INT,
    DueDateKey INT,
    ShipDateKey INT,
    CustomerKey INT,
    PromotionKey INT,
    CurrencyKey INT,
    SalesTerritoryKey INT,
    ....(all other columns)
WITH (
        LOCATION = '/factinternetsalesbig/parquet',
        DATA_SOURCE = ExternalDataSourceDataLake,
        FILE_FORMAT = SynapseParquetFormat
```

**Connect to table within Power BI Dataflow**

🔍 Search

Display options ⌄                                    ↻

▲ 📁 Azure Synapse Analytics (SQL DW)          [2]

  ☐ ▦ DW.FactInternetSalesBigCSV

  ☑ ▦ DW.FactInternetSalesBigParquet

- If you have **CSV, Parquet or JSON** files in Storage, let SQL Serverless do the data processing

- Can create a Dataflow and use **Serverless SQL Pool Database and Tables** as a data source

- SQL is pushed down to SQL Serverless due to Power Query's **Query Folding** feature

# Power BI Query Folding

When transformations are applied to the dataset and can be "folded", the logic to perform the transformations is passed to the data source. In this scenarios, Serverless SQL Pool will receive a SQL statement.

The GROUP BY transformation in Dataflows

The SQL generated by the transformation which will be sent to the Serverless SQL Pool (Query Folding)
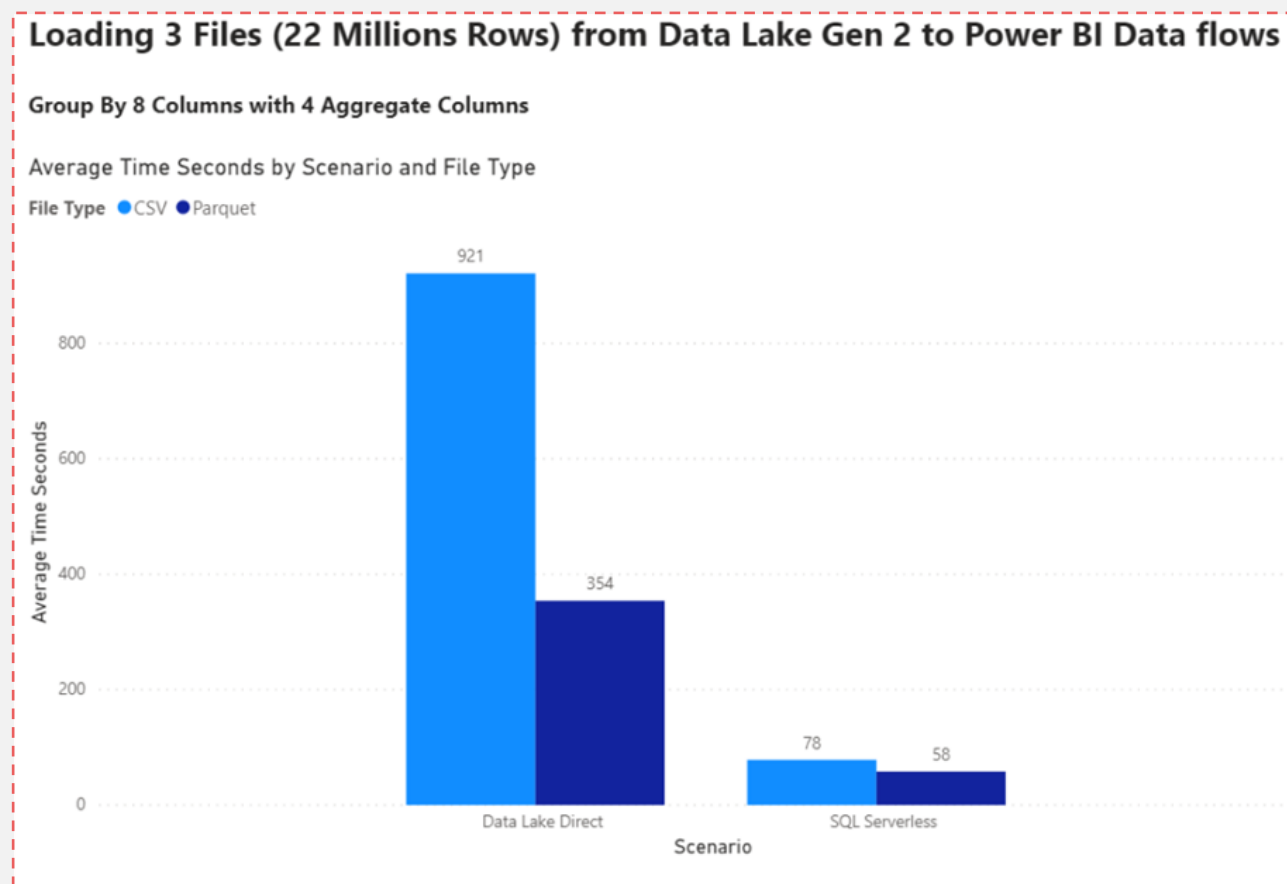
# Serverless SQL Pools Monitoring



The SQL generated by Power BI can be seen in the Monitor area of Synapse Analytics Studio.

The statistics include the query time and also the Data Processed amount, which is a vital statistic.

# Data Loading Performance Analysis

Comparing the performance of Serverless SQL Pools and the Power BI native Data Lake connector
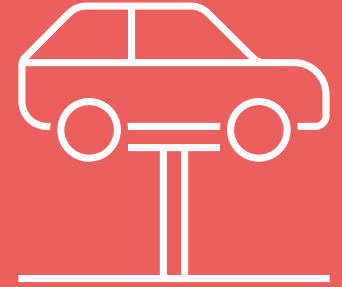


Loading 3 Files (22 Millions Rows) from Data Lake Gen 2 to Power BI Data flows

Group By 8 Columns with 4 Aggregate Columns

Average Time Seconds by Scenario and File Type

File Type ● CSV ● Parquet

The data loading tests were carried out on **CSV** and **Parquet** data.
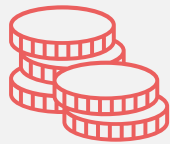
Total files size together: **4.5 GB**

Pushing the data transformations down to the Serverless SQL Pool has reduced the time taken to load the Power BI Dataflow.

**Your results may vary!**
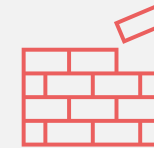
# Considerations

## Cost

SQL Serverless does cost money!
£3.727 per 1 Terabyte (1TB) of Data Processed

When developing/testing, use a smaller file or set of smaller files

Data at rest does not necessarily translate directly into data processed

## Infrastructure

Adds another service into a data architecture which will need managing

However, you can use Synapse Analytics SQL Serverless as a processing engine without any data warehousing

# Reference

[https://www.datahai.co.uk/power-bi/harnessing-azure-synapse-analytics-sql-serverless-in-power-bi-dataflows/](https://www.datahai.co.uk/power-bi/harnessing-azure-synapse-analytics-sql-serverless-in-power-bi-dataflows/)