# Niloofar Moosavi

## Big Data and Cloud - DE23

## Architecture

The key components include:

1. **Azure Subscription**: A subscription was created to manage all the resources.

2. **Azure Resource Group**: A resource group was established to organize all related services.

3. **Azure Storage Account**: A storage account was created with three containers:

   - **StructuredData**: Contains a CSV file with sales data from Amazon, sourced from Kaggle.
   - **UnstructuredData**: Contains an XML log file that was generated in Python.
   - **Curated**: Transformed data is stored here in Parquet format.

   **Creating Azure Data Lake Storage**:



4. **Azure Data Factory (ADF)**: ADF is utilized to orchestrate data movement and transformation. Data from the storage account is extracted, transformed, and loaded into the curated folder in Parquet format.

   **Creating Azure Data Factory**:

   - The screenshot below demonstrates the initial setup process for ADF, including selecting the appropriate resource group and naming the data factory.

5. **Azure Synapse Analytics**: Synapse is used for data exploration and querying. It allows running Spark SQL to create tables on top of the Parquet files and analyze the transformed data.

   **Creating Synapse**:

   ○ The screenshot illustrates the setup of Synapse, including the selection of performance tiers and configurations for serverless resources.

   

   **Data Exploration**:

   I used this sql query to create a table on top of the data

   ```sql
   %%sql
   CREATE TABLE logs USING PARQUET
   LOCATION
   'abfss://curated@storageaccountbigdata23.dfs.core.windows.net/logs'
   ```

   Here are some example results from the log data:

```
1    select count(*), referrer from logs group by referrer
     ✓ - Command executed in 2 sec 862 ms on 8:57:15 AM, 10/27/24
```

View  [ Table ]  [ Chart ]        ↦ Export results ∨

| count(1) | referrer |
|----------|----------|
| 450949 | null |
| 107728 | google.com |
| 107288 | direct |
| 107885 | facebook.com |

```
1    SELECT count(*), action from logs group by action
     ✓ - Command executed in 1 sec 115 ms on 8:57:27 AM, 10/27/24
```

View  [ Table ]  [ Chart ]        ↦ Export results ∨

| count(1) | action |
|----------|--------|
| 32201 | logout |
| 322901 | visit_page |
| 192899 | view_product |
| 128975 | purchase |
| 96874 | add_to_cart |

And some example results from the sales data:

```
1    select count(*), fulfilment from sales group by fulfilment
     ✓ - Command executed in 1 sec 115 ms on 9:01:26 AM, 10/27/24
```

View  [ Table ]  [ Chart ]        ↦ Export results ∨

| count(1) | fulfilment |
|----------|-----------|
| 31589 | Merchant |
| 82112 | Amazon |

```
1    select sum(TotalAmount), productCategory from sales group by productCategory
     ✓ - Command executed in 1 sec 143 ms on 9:01:31 AM, 10/27/24
```

View  [ Table ]  [ Chart ]        ↦ Export results ∨

| sum(TotalAmount) | productCategory |
|------------------|-----------------|
| 762949 | Ethnic Dress |
| 5242931 | Top |
| 37934434 | Set |
| 915 | Dupatta |
| 125767 | Saree |
| 142870 | Bottom |
| 10707932 | Western Dress |
| 441259 | Blouse |
| 20675349 | kurta |

```
1    SELECT sum(TotalAmount), orderDateTransformed as date from sales group by orderDateTransfor
     ✓ - Command executed in 1 sec 125 ms on 9:01:39 AM, 10/27/24
```

View  [ Table ]  [ Chart ]        ↦ Export results ∨

| sum(TotalAmount) | date |
|------------------|------|
| 717201 | 2022-05-17 |
| 1170547 | 2022-05-04 |
| 902319 | 2022-06-01 |
| 858268 | 2022-04-04 |
| 943102 | 2022-06-07 |
| 716058 | 2022-05-16 |
| 719271 | 2022-06-03 |
| 745190 | 2022-05-19 |
| 755686 | 2022-06-15 |
| 903676 | 2022-05-07 |
| 888252 | 2022-06-04 |
| 745213 | 2022-05-24 |
| 851288 | 2022-04-12 |

6. **Power BI**: Power BI is employed to create dashboards and visualizations based on the transformed data stored in Azure Synapse.

# Data Management

Data management for this project was achieved through a comprehensive ETL (Extract, Transform, Load) pipeline built in Azure Data Factory (ADF). This pipeline was designed to ensure data quality, facilitate efficient transformations, and prepare datasets for advanced analysis. Below are the detailed steps undertaken:

1. Extraction:

   ○ The sales data, sourced from a structured CSV file, was ingested into Azure Data Factory from the StructuredData container in the Azure Storage Account.

   ○ The XML log data, from the UnstructuredData container, was also ingested for further processing.

2. Data Transformation

   **Sales data:**

   ○ Type Casting: The quantity column was cast to integer and the amount column to double to standardize the data types for accurate calculations.

   ○ Data Filtering: A set of filters was applied to ensure data consistency and reliability: Only rows where Amount > 0 and Qty > 0 were retained.

   ○ Null or empty values in the currency field were excluded (!isNull(currency)).

   ○ Column Selection and Renaming: Only necessary columns were selected and renamed to align with naming conventions. This step reduced redundancy and improved clarity for downstream processes.

   ○ Derived Columns: Additional calculated fields were introduced: totalAmount was computed as the product of price and quantity.

   ○ Date Transformation: The date column was reformatted from dd-mm-yy to the standardized yyyy-mm-dd format to ensure compatibility with analysis tools.

   **log data:**

   ○ Data Flattening: The nested structure of the XML log data was flattened to extract relevant details into a tabular format.

   ○ Column Selection: A subset of critical columns was chosen for further analysis to minimize processing overhead.

   ○ Derived Field: A unique identifier, logId, was created by concatenating userId and timestamp. This step ensured a unique key for tracking individual log events across the dataset.

3. Loading and Storage

   The transformed datasets were stored in the Curated container of the Azure Storage Account in Parquet format. This format was chosen for its efficient storage and querying capabilities.
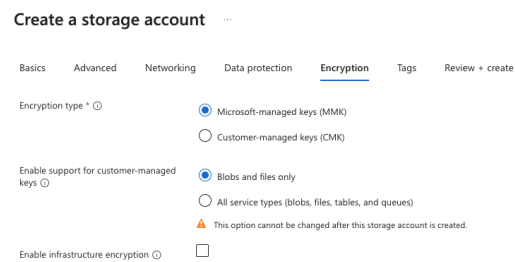
## ETL Process Visual Representation

To illustrate the data management workflow, the following diagram summarizes the key steps involved in the ETL pipeline:
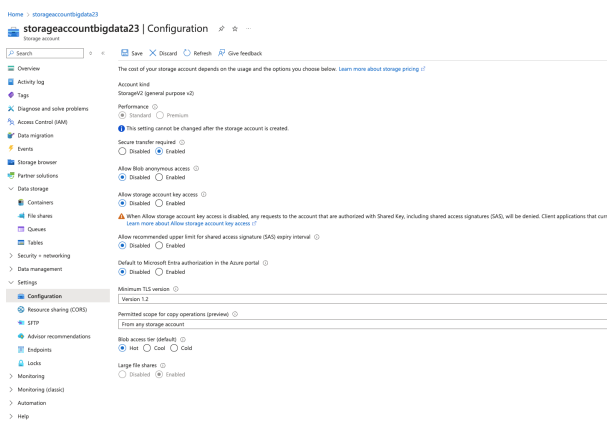


# Security

- Data in the storage account is encrypted at rest.



  Morevoer, Azure uses Secure Socket Layer (SSL)/Transport Layer Security (TLS) protocols to encrypt data as it travels between applications and Azure services. This safeguards data from interception during transmission.

- To enhance the security posture of the data architecture, key access to the storage account was disabled. The screenshot highlights this configuration, ensuring that access is managed solely through role-based permissions.
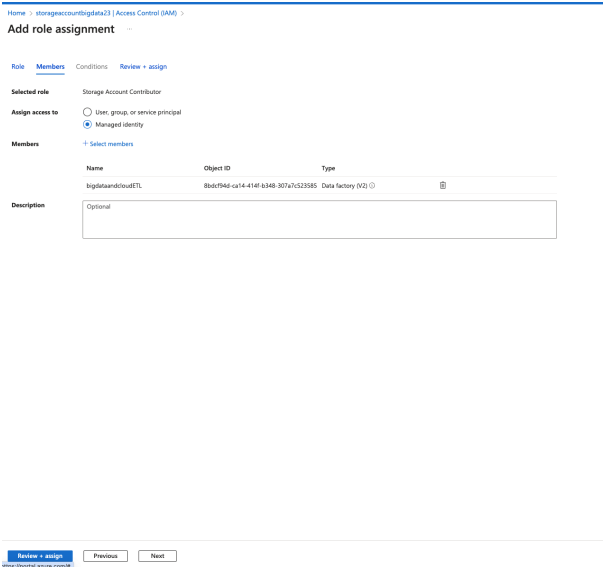


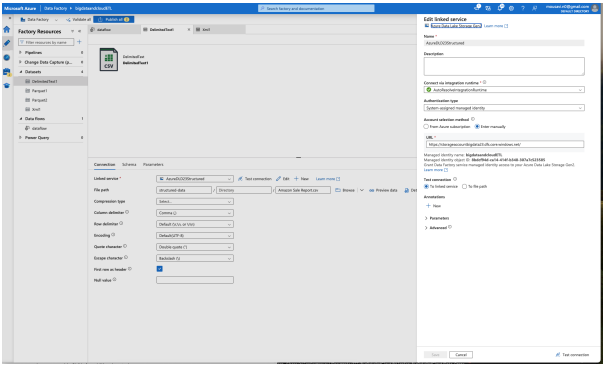- RBAC was used to access data in the blob storage:

  **Data factory's access to data:

  - To ensure proper security measures, the following settings were configured:

    - **Blob Contributor Role**: This role allows ADF to access the storage account without exposing key access. The screenshot shows the assignment of the Blob Contributor role
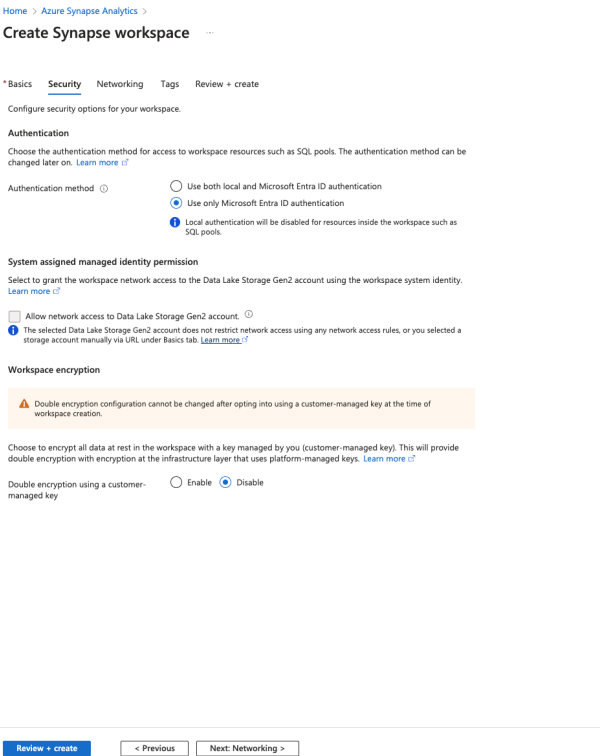
to the managed identity of ADF.



- ■ **Dataset Authentication**: The authentication type is selected to be 'system-assigned-managed-identity'.



**RBAC for Synapse**: The screenshot shows how access is managed, ensuring only authorized users can interact with Synapse resources.

# Scalability

**Serverless SQL Pool**: Utilizing the serverless option in Synapse allows for on-demand querying of data stored in the Parquet files without the need for provisioning dedicated resources. The screenshot displays the setup options available for creating serverless SQL pools, emphasizing scalability, cost-effectiveness, and flexibility in querying large datasets.



# Power BI Report

**Sales Insights**:

- Finally, the Power BI dashboard provides visual insights derived from the sales and logs data.