# Retail Sales Analysis
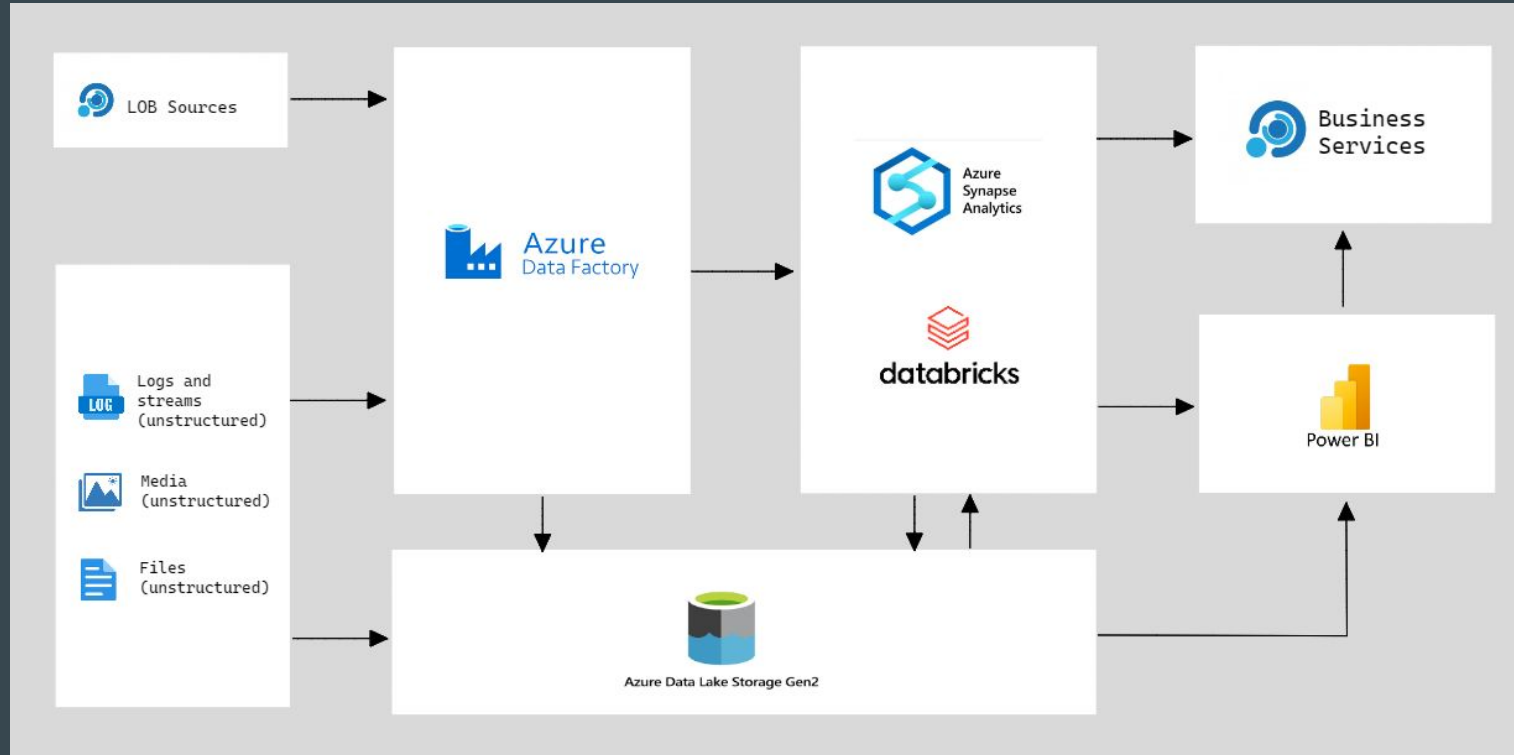
...

Group 06
Roshan Rawat, Gangothri Gadige, Chandan Kumar,
Narendra Reddy, Faryar Memon

# Modern Data Warehouse Architecture

# Tasks

| | | |
|---|---|---|
| **Source** | Azure SQL | ● Collecting the data from the client source |
| **ETL** | Azure Data Factory | ● Using **ADF** to ingest data from client source **Azure SQL Database** to **Azure Data Lake Storage** |
| **Analysis** | Azure Synapse Analytics / databricks / Spark | ● Integrating ADLS with **Synapse** using **Linked Services** & mounting ADLS to **Azure Databricks** |
| **Visualization** | Power BI | ● Integrating the analysis results with **Power BI** for impactful visualizations and reporting |

# Managed by DE Expert Venkat Sir✨

## Our Sub-Teams:

| IAM Team | Data Integration Team | Data Analytics & BI Team |
|---|---|---|
| Managing Azure AD resources and access controls, configuring authentication methods and monitoring access activities | Extracting data from various sources, transforming and loading data into data warehouse; maintaining ETL workflows and pipelines | Perform big data analysis and create dashboards to communicate findings |
| - Faryar<br>- Roshan | - Gangothri<br>- Roshan | - Chandan<br>- Narendra<br>- Roshan |

# Why add multiple users to a single Resource Group?

- Centralized Access Management
- Seamless Collaboration and Resource Sharing
- Improved Security
- Cost Optimization

This simplifies administration, enhances collaboration, improves security, and optimizes costs

# Creating Users Under One Azure Active Directory

# Providing Privileges Resource Group and Subscriptions

# Azure Data Factory

## Code-Free ETL as a service

### Ingest

- Multi-cloud and on-premise hybrid copy data
- 100+ native connectors
- Serverless and auto-scale
- Use wizard for quick copy jobs

### Control Flow

- Design code-free data pipelines
- Generate pipelines via SDK
- Utilize workflow constructs: loops, branches, conditional execution, variables, parameters, …

### Data Flow

- Code-free data transformations that execute in Spark
- Scale-out with Azure Integration Runtimes
- Generate data flows via SDK
- Designers for data engineers and data analysts

### Schedule

- Build and maintain operational schedules for your data pipelines
- Wall clock, event-based, tumbling windows, chained

### Monitor

- View active executions and pipeline history
- Detail activity and data flow executions
- Establish alerts and notifications

←— ADF Pipeline "SQL to ADLS" in action

ADLS Storage -—>

# Azure Synapse



**Synapse Workspace** is the overall container that holds all the resources and settings related to Azure Synapse, while **Synapse Studio** is the web-based interface within the workspace that provides tools and features for development, management, and collaboration on Synapse components.

# Analysis

```
df=spark.read.load('abfss://source@retailsales123.dfs.core.windows.net/retail.csv
', format='csv'
, header=True,inferSchema=True
)
```

# Pre-processing:

```
retail_df = df.withColumn('order date', to_date(col('Order Date'), 'MM/dd/yyyy'))
            .withColumn('ship date', to_date(col('Ship Date'), 'MM/dd/yyyy'))
```

# User Stories

## User Story: Monthly Sales Report

```
monthly_sales = retail_df.groupBy([year("Order
Date").alias("Years"), month("Order
Date").alias("Month")]).agg(sum("Total
Cost").alias("Sales")).orderBy("Years","Month")

## Output:
+-----+-----+------------------+
|Years|Month|   Sales          |
+-----+-----+------------------+
| 2010|    1| 5.6117500219999984E7|
| 2010|    2| 6.149703725999998E7|
| 2010|    3| 3.2563351479999997E7|
| 2010|    4| 5.008999801000001E7|
| 2010|    5| 6.251231557999999E7|
| 2010|    6| 5.518854506000001E7|
| 2010|    7| 2.8828290519999996E7|
| 2010|    8| 5.163789438999999E7|
| 2010|    9| 4.10263253E7|
| 2010|   10| 4.209123058E7|
| 2010|   11| 3.881841093000001E7|
| 2010|   12| 5.583074315E7|
| 2011|    1| 3.531044195000001E7|
| 2011|    2| 4.812870418000015E7|
| 2011|    3| 4.932200064999999E7|
| 2011|    4| 3.0250918640000004E7|
| 2011|    5| 5.099740488000001E7|
| 2011|    6| 4.028547157E7|
| 2011|    7| 6.0774376599999994E7|
| 2011|    8| 5.4067997589999996E7|
+-----+-----+------------------+
```

## User Story: Quarterly Sales Report

```
quarterly_sales_total_revenue =
retail_df.groupBy([year("Order
Date").alias("Year"), quarter("Order
Date").alias("Quarter")]).agg(sum("Total
Revenue")).orderBy("Year","Quarter")

## Output:
+----+-------+------------------+
|Year|Quarter| sum(Total Revenue) |
+----+-------+------------------+
|2010|      1|2.1568648214000002E8|
|2010|      2| 2.328106550899999E8|
|2010|      3|1.7718539354999992E8|
|2010|      4| 1.920828373899999E8|
|2011|      1| 1.958348384799999E8|
|2011|      2|     1.7647134479E8|
|2011|      3|       2.3487859321E8|
|2011|      4|2.3043744375999996E8|
|2012|      1|2.5948719668999994E8|
|2012|      2|2.2292645934999993E8|
|2012|      3|2.1255167141999993E8|
|2012|      4|2.0868209533999994E8|
|2013|      1|2.3419277444999996E8|
|2013|      2|1.8661791348999992E8|
|2013|      3|2.4164184929999998E8|
|2013|      4|2.4042816672000012E8|
|2014|      1|2.3005481114000008E8|
|2014|      2|2.1748007145999995E8|
|2014|      3|2.1132196698999998E8|
|2014|      4|1.9754089597000015E8|
+----+-------+------------------+
```

# User Stories

**User Story: Display the number of orders for each item**

```python
each_item = rdd_df.map(lambda l: (l[2], 1))
item_ocount=each_item.reduceByKey(add).sortBy
(    lambda x: x[1],ascending=False)


 ## Output:
[('Beverages', 447),
('Fruits', 447),
('Baby Food', 445),
('Cosmetics', 424),
('Household', 424),
('Office Supplies', 420),
('Personal Care', 415),
('Vegetables', 410),
('Meat', 399),
('Snacks', 398),
('Clothes', 386),
('Cereal', 385)]


item_type_orders_count =
spark.createDataFrame(item_ocount,
['item_type', 'orders_count'])
```

**User Story: Display the country with highest sale**

```python
country_order = rdd_df.map(lambda l: (l[1], l[11]))
country_revenue=country_order.reduceByKey(add).sortBy(
lambda x: x[1], ascending=False)

countrywise_sales =
spark.createDataFrame(country_revenue, ['country',
'sales'])


 ## Output:
 +--------------------+--------------------+
 | country| sales|
 +--------------------+--------------------+
 | Rwanda| 6.039873958999999E7|
 | Myanmar| 5.883846785E7|
 | South Korea| 5.743435542999999E7|
 | Ghana| 5.627138265000001E7|
 | Niger| 5.529821127999998E7|
 | Grenada| 5.498818455000001E7|
 |Republic of the C...| 5.437980838999998E7|
 | Kosovo| 5.383314279000001E7|
 | Czech Republic| 5.354393214E7|
 | Ukraine| 5.325231754E7|
 | Vanuatu| 5.129172321E7|
```

# Databricks

```
dbutils.fs.mount(
        source='wasbs://mycontainer151@mystorageaccount151160.blob.core.windows.net/',
    mount_point='/mnt/capstun',
    extra_configs ={'fs.azure.account.key.mystorageaccount151160.blob.core.windows.net':"y33+zi2yumQufcCd+G5DfbU5iHt/73uCPQ91wZpCpViqDwCLst4quOE/
1onBQNMC4WJPkYEwIaa1+AStJ9vKHg=="}
        )
```

```
1   sales_report=df.groupBy(year(df['Order Date'])).agg({'Total Cost':'sum','Total Revenue':'sum','Total Profit':'sum'})
```

▸ ▦ sales_report: pyspark.sql.dataframe.DataFrame = [year(Order Date): integer, sum(Total Cost): double ... 2 more fields]

Command took 0.15 seconds -- by narendrareddi2000@gmail.com at 6/10/2023, 11:05:55 PM on Narendra Reddy's Cluster

Cmd 6

```
1   sales_report.show()
```

▸ (2) Spark Jobs

```
+----------------+-------------------+--------------------+--------------------+
|year(Order Date)|    sum(Total Cost)|   sum(Total Revenue)|    sum(Total Profit)|
+----------------+-------------------+--------------------+--------------------+
|            2015|    7.1175791073E8|9.845037861999993E8| 2.727458754699999E8|
|            2013|6.3788361501999995E8|9.028807039600003E8|2.649970889400001E8|
|            2014|5.999882564799997E8|8.563977455599985E8|2.5640948907999977E8|
|            2012|6.279867956999996E8|9.036474228000004E8|2.756606271000002E8|
|            2016|    6.0137735788E8|8.553197613500007E8|2.5394240346999985E8|
|            2010|5.762016424799999E8|8.177653681700003E8|2.4156372568999994E8|
|            2011|5.793084282600002E8|    8.3762222024E8|2.5831379197999984E8|
|            2017|3.3096197290000004E8|4.7055220025999993E8|1.3959022736000013E8|
+----------------+-------------------+--------------------+--------------------+
```

```
1   country_report=df.groupBy([year(df['Order Date']).alias('year'),df['Country']]).agg({'Total Cost':'sum','Total Revenue':'sum','Total Profit':'sum'})
```

▸ ▦ country_report: pyspark.sql.dataframe.DataFrame = [year: integer, Country: string ... 3 more fields]

Command took 0.09 seconds -- by narendrareddi2000@gmail.com at 6/10/2023, 11:20:48 PM on Narendra Reddy's Cluster

Cmd 9

Python

```
1   country_report.orderBy(country_report.year.asc(),country_report.Country.asc()).show()
```

▸ (2) Spark Jobs

```
+----+-------------------+----------------+--------------------+----------------+
|year|            Country|  sum(Total Cost)|   sum(Total Revenue)|  sum(Total Profit)|
+----+-------------------+----------------+--------------------+----------------+
|2010|            Albania| 5798207.049999999|         8531301.81|      2733094.76|
|2010|            Algeria|         631933.79|        1131473.33|499539.5399999999|
|2010|            Andorra|7220305.570000001|       1.046352798E7|      3243222.41|
|2010|Antigua and Barbuda|        9053361.49|      1.35583839E7|      4505022.41|
|2010|            Armenia|        5741950.42|        7672185.09|      1930234.67|
|2010|          Australia|        2517154.85|4133990.7899999996|1616835.9400000002|
|2010|            Austria|      8142984.28|1.1632719089999998E7|      3489734.81|
|2010|         Azerbaijan|          74388.6|         111033.0|         36644.4|
|2010|            Bahrain|4016249.5999999996|        5940840.98|      1924591.42|
|2010|         Bangladesh|        1163331.22|        1557794.74|       394463.52|
|2010|           Barbados|        1033108.59|        1639914.55|       606805.96|
|2010|            Belarus|        1306786.99|  2080164.1800000002|       773377.19|
|2010|            Belgium|         830578.56|         1300591.92|       470013.36|
+----+-------------------+----------------+--------------------+----------------+
```

# Store data

```
adls_account_name = "retailsales123"
adls_container = "source"
adls_folder = "output"
adls_output_path =
r'abfs://source@retailsales123.dfs.core.w
indows.net/output/countrywise_sales'

# Save the DataFrame to ADLS
countrywise_sales .write \
    .mode("overwrite") \
    .csv(adls_output_path)
```
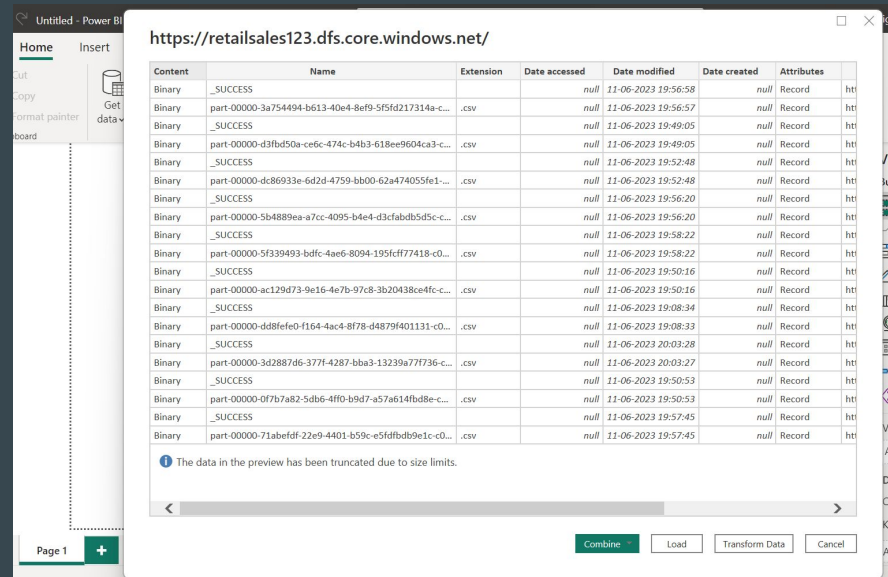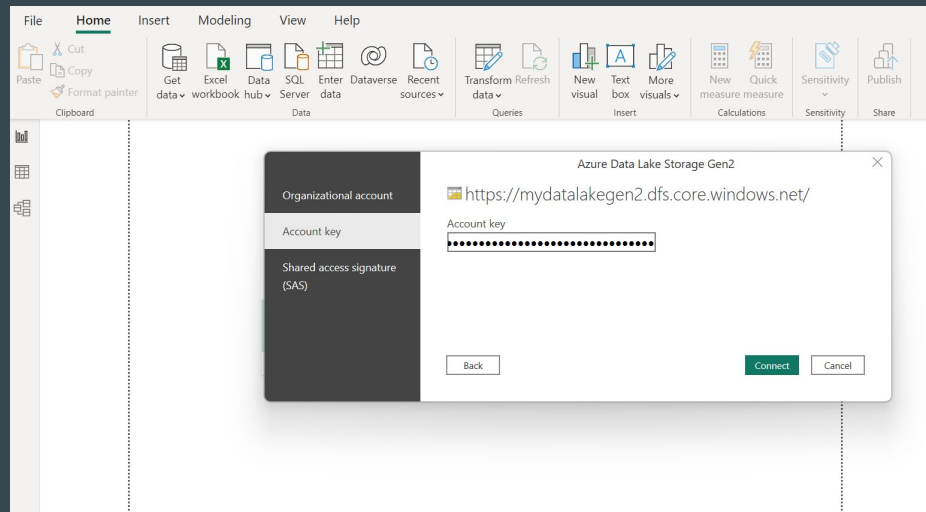
**source** ...
Container

Search

Overview

Diagnose and solve problems

Access Control (IAM)

**Settings**

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

↑ Upload    + Add Directory    ↻ Refresh

Name

☐ 📁 [..]
☐ 📁 countrywise_sales
☐ 📁 countrywise_yearly_sales
☐ 📁 half_yearly_sales
☐ 📁 item_type_orders_count
☐ 📁 monthly_report_2014
☐ 📁 monthly_sales
☐ 📁 products_with_a_occurances
☐ 📁 quarterly_report_2014
☐ 📁 quarterly_sales_total_revenue
☐ 📁 specific_country_yearly_sales
☐ 📁 yearly_sales

# RETAIL SALES DASHBOARD

**COUNTRY**

All

## COST VS REVENUE

● Total Cost ● Total Revenue

1,000M

500M

0M

2010 · 2012 · 2014 · 2016

## CATEGORYWISE PROFIT

Beverages 1.76%
Snacks 5.4%
Meat 5.54%
Vegetables 6.72%
Clothes 7.42%
Cereal 8.8%
Baby Food 11.11%
Office Supplies 13.55%
Household 17.72%
Cosmetics 18.97%

## MONTHWISE PROFIT

Sales Channel ● Offline ● Online

100M

90M

80M

70M

January, February, March, April, May, June, July, August, September, October, November, December

## PROFIT AND REVENUE W.R.T REGION

● Profit ● Revenue

| Region | Profit | Revenue |
|---|---|---|
| Sub-Saharan Africa | 531M | 1815M |
| Europe | 502M | 1704M |
| Asia | 278M | 920M |
| Middle East and North Africa | 231M | 767M |
| Central America and the Caribbean | 205M | 685M |
| Australia and Oceania | 175M | 587M |
| North America | 41M | 151M |

# THANK YOU VENKAT SIR