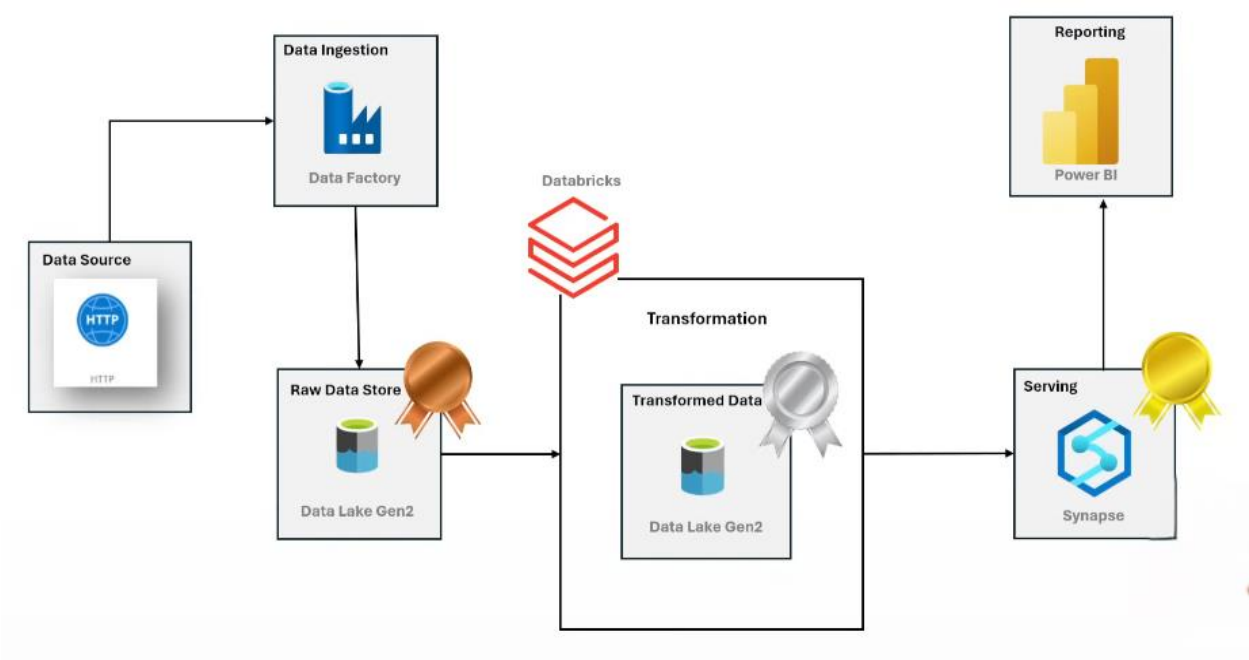


Azure End-To-End Data Engineering Project from Scratch

ARCHITECTURE



Home >

dataeng-e2e-dev-rg What are the best practices for managing this resource group? How do I troubleshoot issues with this resource group? +1

Search

Create Manage view Delete resource group Refresh Export to CSV Open query Assign tags Move Delete

Overview

Activity log

Access control (IAM)

Tags

Resource visualizer

Events

Settings

Cost Management

Monitoring

Automation

Help

Essentials

Resources Recommendations (14)

Filter for any field... Type equals all Location equals all Add filter

Showing 1 to 6 of 6 records. Show hidden types No grouping List view

Name	Type	Location
dataeng-e2e-dev-databricks	Azure Databricks Service	West US 2
dataeng-e2e-dev-kvs	Key vault	West US 2
dataeng-e2e-dev-synapse	Synapse workspace	West US 2
dataaenge2edevadf	Data factory (V2)	West US 2
dataaenge2edevstore	Storage account	West US 2
dataaenge2edevsynstore	Storage account	West US 2

Data Factory | Validate all | Publish all | Auto Save | Preview experience | On

Factory Resources

Filter resources by name

- Pipelines 1
 - DynamicGitToRaw
- Change Data Capture (preview) 0
- Datasets 5
 - ds_git_dynamic
 - ds_git_parameters
 - ds_http
 - ds_raw
 - ds_sink_dynamic
- Data flows 0
- Power Query 0

Activities

Search activities

- Move and transform
- Synapse
- Azure Data Explorer
- Azure Function
- Batch Service
- Databricks
- Data Lake Analytics
- General
- HDInsight
- Iteration & conditionals
- Machine Learning
- Power Query

✓ Validate | ▶ Debug | ⚡ Add trigger

Parameters | Variables | Settings | Output

Home > dataange2edevstore | Containers >

bronze

Container

Search

+ Add Directory | ↑ Upload | ↻ Refresh | 🗑 Delete | 📄 Copy | 📄 Paste | 🔄 Rename | 🔄 Acquire lease | 🔄 Break lease | 🛠 Edit columns

Overview

🔧 Diagnose and solve problems

🔑 Access Control (IAM)

⌵ Settings

🔑 Shared access tokens

🔑 Manage ACL

🔑 Access policy

📄 Properties

📄 Metadata

bronze

Authentication method: Access key ([Switch to Microsoft Entra user account](#))

🔍 Search blobs by prefix (case-sensitive)

Only show active objects

Showing all 10 items

<input type="checkbox"/>	Name	Last modified	Access tier	Blob type	Size	Lease state
<input type="checkbox"/>	AdventureWorks_Calendar	8/12/2025, 10:04:51 AM				...
<input type="checkbox"/>	AdventureWorks_Customers	8/12/2025, 10:05:09 AM				...
<input type="checkbox"/>	AdventureWorks_Product_Cat...	8/12/2025, 10:04:33 AM				...
<input type="checkbox"/>	AdventureWorks_Products	8/12/2025, 10:05:41 AM				...
<input type="checkbox"/>	AdventureWorks>Returns	8/12/2025, 10:06:04 AM				...
<input type="checkbox"/>	AdventureWorks_Sales_2015	8/12/2025, 10:06:21 AM				...
<input type="checkbox"/>	AdventureWorks_Sales_2016	8/12/2025, 10:06:36 AM				...
<input type="checkbox"/>	AdventureWorks_Sales_2017	8/12/2025, 10:06:53 AM				...

Microsoft Azure databricks

Search data, notebooks, recents, and more... CTRL + P

dataeng-e2e-dev-databricks

New

- Workspace
- Recents
- Catalog
- Jobs & Pipelines
- Compute
- Marketplace

SQL

- SQL Editor
- Queries
- Dashboards
- Genie
- Alerts
- Query History
- SQL Warehouses

Data Engineering

- Job Runs
- Data Ingestion

Workspace

- Home
- Shared with me New
- Workspace
 - Repos
 - Shared
- Users
 - ombabujollireddy@com
 - Favorites
 - Trash

Workspace > Users > ombabujollireddy@com >

dataeng-e2e-dev-notebooks ☆

Search

Type Owner Last modified

Name ↕	Type	Owner	Created at
silver_layer	Notebook	Ombabu Jollireddy	Aug 11, 2025, 04:...

Send feedback Share Create

2 minutes ago (<1s)

```
from pyspark.sql.functions import *
from pyspark.sql.types import *
```

1

SILVER LAYER SCRIPT

DATA ACCESS USING MOUNT

3 minutes ago (1s)

```
spark.conf.set(
    "fs.azure.account.key.dataenge2edevstore.dfs.core.windows.net",
    dbutils.secrets.get(scope="dataeng-e2e-dev-scope", key="dataeng-storage-secret")
)
```

4

```
▶ ✓ 2 minutes ago (<1s) 5

mount_point = "/mnt/dataenge2edevstore/bronze"

if not any(mount.mountPoint == mount_point for mount in dbutils.fs.mounts()):
    dbutils.fs.mount(
        source = "wasbs://bronze@dataenge2edevstore.blob.core.windows.net",
        mount_point = mount_point,
        extra_configs = {
            "fs.azure.account.key.dataenge2edevstore.blob.core.windows.net": dbutils.secrets.get
            (scope="dataeng-e2e-dev-scope", key="dataeng-storage-secret")
        }
    )
else:
    print(f"Already mounted: {mount_point}")
```

Already mounted: /mnt/dataenge2edevstore/bronze

DATA LOADING

Read the Data

```
▶ ✓ 3 minutes ago (1s) 8

df_cal = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/AdventureWorks_Calendar")
```

▶ (2) Spark Jobs

▶ df_cal: pyspark.sql.dataframe.DataFrame = [Date: date]

```
▶ ✓ 3 minutes ago (1s) 9

df_cus = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/AdventureWorks_Customers")
```

▶ (2) Spark Jobs

▶ df_cus: pyspark.sql.dataframe.DataFrame = [CustomerKey: integer, Prefix: string ... 11 more fields]

```
▶ ✓ 3 minutes ago (1s) 10

df_procat = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/AdventureWorks_Product_Categories")
```

▶ (2) Spark Jobs

▶ df_procat: pyspark.sql.dataframe.DataFrame = [ProductCategoryKey: integer, CategoryName: string]

▶ ✓ 4 minutes ago (1s)

11

```
df_pro = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/AdventureWorks_Products")
```

▶ (2) Spark Jobs

▶ df_pro: pyspark.sql.dataframe.DataFrame = [ProductKey: integer, ProductSubcategoryKey: integer ... 9 more fields]

▶ ✓ 4 minutes ago (1s)

12

```
df_ret = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/AdventureWorks>Returns")
```

▶ (2) Spark Jobs

▶ df_ret: pyspark.sql.dataframe.DataFrame = [ReturnDate: date, TerritoryKey: integer ... 2 more fields]

▶ ✓ 4 minutes ago (1s)

13

```
df_sales = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/AdventureWorks_Sales")
```

▶ (2) Spark Jobs

▶ df_sales: pyspark.sql.dataframe.DataFrame = [OrderDate: date, StockDate: date ... 6 more fields]

▶ ✓ 4 minutes ago (1s)

14

```
df_ter = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/AdventureWorks_Territories")
```

▶ (2) Spark Jobs

▶ df_ter: pyspark.sql.dataframe.DataFrame = [SalesTerritoryKey: integer, Region: string ... 2 more fields]

▶ ✓ 10:18 AM (1s)

15

```
df_subcat = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("/mnt/dataenge2edevstore/bronze/Product_Subcategories")
```

▶ (2) Spark Jobs

▶ df_subcat: pyspark.sql.dataframe.DataFrame = [ProductSubcategoryKey: integer, SubcategoryName: string ... 1 more field]

TRANSFORMATIONS

Calendar

▶ ✓ 10:18 AM (<1s) 18

```
df_cal.display()
```

▶ (1) Spark Jobs

Table ▾ + 🔍 ⚙️ 📄

	📅 Date
1	2015-01-01
2	2015-01-02
3	2015-01-03
4	2015-01-04
5	2015-01-05

▶ ▾ ✓ 10:18 AM (<1s) 19 Python ⚙️ 📄

```
df_cal = df_cal.withColumn('Month', month(col('Date')))\
| | | .withColumn('Year', year(col('Date')))\
df_cal.display()
```

▶ (1) Spark Jobs

▶ 📄 df_cal: pyspark.sql.dataframe.DataFrame = [Date: date, Month: integer ... 1 more field]

Table ▾ + 🔍 ⚙️ 📄

	📅 Date	1 ² ₃ Month	1 ² ₃ Year
1	2015-01-01	1	2015
2	2015-01-02	1	2015
3	2015-01-03	1	2015
4	2015-01-04	1	2015
5	2015-01-05	1	2015

```
▶ ✓ 10:18 AM (<1s) 20

mount_point = "/mnt/dataenge2devstore/silver"

if not any(mount.mountPoint == mount_point for mount in dbutils.fs.mounts()):
    dbutils.fs.mount(
        source = "wasbs://silver@dataenge2devstore.blob.core.windows.net",
        mount_point = mount_point,
        extra_configs = {
            "fs.azure.account.key.dataenge2devstore.blob.core.windows.net": dbutils.secrets.get
            (scope="dataeng-e2e-dev-scope", key="dataeng-storage-secret")
        }
    )
else:
    print(f"Already mounted: {mount_point}")

Already mounted: /mnt/dataenge2devstore/silver
```

```
▶ ✓ 10:18 AM (2s) 21

df_cal.write.format('parquet')\
    .mode('append')\
    .save('/mnt/dataenge2devstore/silver/AdventureWorks_Calendar')

▶ (1) Spark Jobs
```

Customers

```
▶ ✓ 10:18 AM (1s) 23

df_cus.display()

▶ (1) Spark Jobs
```

	1 ² ₃ CustomerKey	A ^B _C Prefix	A ^B _C FirstName	A ^B _C LastName	📅 BirthDate	A ^B _C MaritalStatus	A ^B _C Gender
1	11000	MR.	JON	YANG	1966-04-08	M	M
2	11001	MR.	EUGENE	HUANG	1965-05-14	S	M
3	11002	MR.	RUBEN	TORRES	1965-08-12	M	M
4	11003	MS.	CHRISTY	ZHU	1968-02-15	S	F
5	11004	MRS.	ELIZABETH	JOHNSON	1968-08-08	S	F

10:18 AM (1s) 24

```
df_cus = df_cus.withColumn('fullname', concat_ws(' ', col('Prefix'), col('FirstName'), col('LastName')))
df_cus.display()
```

(1) Spark Jobs

df_cus: pyspark.sql.dataframe.DataFrame = [CustomerKey: integer, Prefix: string ... 12 more fields]

	CustomerKey	Prefix	FirstName	LastName	BirthDate	MaritalStatus	Gender
1	11000	MR.	JON	YANG	1966-04-08	M	M
2	11001	MR.	EUGENE	HUANG	1965-05-14	S	M
3	11002	MR.	RUBEN	TORRES	1965-08-12	M	M
4	11003	MS.	CHRISTY	ZHU	1968-02-15	S	F
5	11004	MRS.	ELIZABETH	JOHNSON	1968-08-08	S	F

10:18 AM (2s) 25

```
df_cus.write.format('parquet')\
    .mode('append')\
    .save('/mnt/dataenge2edevstore/silver/AdventureWorks_Customers')
```

(1) Spark Jobs

Sub Categories

10:18 AM (<1s) 27

```
df_subcat.display()
```

(1) Spark Jobs

	ProductSubcategoryKey	SubcategoryName	ProductCategoryKey
1	1	Mountain Bikes	1
2	2	Road Bikes	1
3	3	Touring Bikes	1
4	4	Handlebars	2
5	5	Bottom Brackets	2

10:18 AM (1s) 28

```
df_subcat.write.format('parquet')\
    .mode('append')\
    .save('/mnt/dataenge2edevstore/silver/AdventureWorks_Subcategories')
```

(1) Spark Jobs

Products

10:18 AM (<1s) 30

```
df_pro.display()
```

(1) Spark Jobs

	ProductKey	ProductSubcategoryKey	ProductSKU	ProductName	ModelName	ProductDe
1	214	31	HL-U509-R	Sport-100 Helmet, Red	Sport-100	Universa
2	215	31	HL-U509	Sport-100 Helmet, Black	Sport-100	Universa
3	218	23	SO-B909-M	Mountain Bike Socks, M	Mountain Bike Socks	Combina
4	219	23	SO-B909-L	Mountain Bike Socks, L	Mountain Bike Socks	Combina
5	220	31	HL-U509-B	Sport-100 Helmet, Blue	Sport-100	Universa

10:18 AM (<1s) 31

```
df_pro = df_pro.withColumn('ProductSKU', split(col('ProductSKU'), '-')[0])\
                .withColumn('ProductName', split(col('ProductName'), ' ')[0])\
df_pro.display()
```

(1) Spark Jobs

df_pro: pyspark.sql.dataframe.DataFrame = [ProductKey: integer, ProductSubcategoryKey: integer ... 9 more fields]

	ProductKey	ProductSubcategoryKey	ProductSKU	ProductName	ModelName	ProductDe
1	214	31	HL	Sport-100	Sport-100	Universal fit, w
2	215	31	HL	Sport-100	Sport-100	Universal fit, w
3	218	23	SO	Mountain	Mountain Bike Socks	Combination o
4	219	23	SO	Mountain	Mountain Bike Socks	Combination o
5	220	31	HL	Sport-100	Sport-100	Universal fit, w

10:18 AM (1s) 32

```
df_pro.write.format('parquet')\
        .mode('append')\
        .save('/mnt/dataenge2edevstore/silver/AdventureWorks_Products')
```

(1) Spark Jobs

10:18 AM (<1s) 33

```
df_ret.display()
```

(1) Spark Jobs

	ReturnDate	TerritoryKey	ProductKey	ReturnQuantity
1	2015-01-18	9	312	1
2	2015-01-18	10	310	1
3	2015-01-21	8	346	1
4	2015-01-22	4	311	1
5	2015-02-02	6	312	1

10:18 AM (1s) 34

```
df_ret.write.format('parquet')\
  .mode('append')\
  .save('/mnt/dataenge2edevstore/silver/AdventureWorks_Returns')
```

(1) Spark Jobs

Territories

10:18 AM (<1s) 36

```
df_ter.display()
```

(1) Spark Jobs

	SalesTerritoryKey	Region	Country	Continent
1	1	Northwest	United States	North America
2	2	Northeast	United States	North America
3	3	Central	United States	North America
4	4	Southwest	United States	North America
5	5	Southeast	United States	North America

10:18 AM (1s) 37

```
df_ter.write.format('parquet')\
  .mode('append')\
  .save('/mnt/dataenge2edevstore/silver/AdventureWorks_Territories')
```

(1) Spark Jobs

Sales

▶ ✓ 10:18 AM (<1s) 39

```
df_sales.display()
```

▶ (1) Spark Jobs

	📅 OrderDate	📅 StockDate	Ⓐ OrderNumber	1 ² ₃ ProductKey	1 ² ₃ CustomerKey	1 ² ₃ TerritoryKey	1 ² ₃ O
1	2017-01-01	2003-12-13	SO61285	529	23791	1	
2	2017-01-01	2003-09-24	SO61285	214	23791	1	
3	2017-01-01	2003-09-04	SO61285	540	23791	1	
4	2017-01-01	2003-09-28	SO61301	529	16747	1	
5	2017-01-01	2003-10-21	SO61301	377	16747	1	

▶ ✓ 10:18 AM (<1s) 40

```
df_sales = df_sales.withColumn('StockDate', to_timestamp('StockDate'))
```

▶ 📖 df_sales: pyspark.sql.dataframe.DataFrame = [OrderDate: date, StockDate: timestamp ... 6 more fields]

▶ ✓ 10:18 AM (<1s) 41

```
df_sales = df_sales.withColumn('OrderNumber', regexp_replace(col('OrderNumber'), 'S', 'T'))
```

▶ 📖 df_sales: pyspark.sql.dataframe.DataFrame = [OrderDate: date, StockDate: timestamp ... 6 more fields]

▶ ✓ 10:18 AM (<1s) 42

```
df_sales = df_sales.withColumn('multiply', col('OrderLineItem') * col('OrderQuantity'))
```

▶ 📖 df_sales: pyspark.sql.dataframe.DataFrame = [OrderDate: date, StockDate: timestamp ... 7 more fields]

▶ ✓ 10:18 AM (<1s) 43

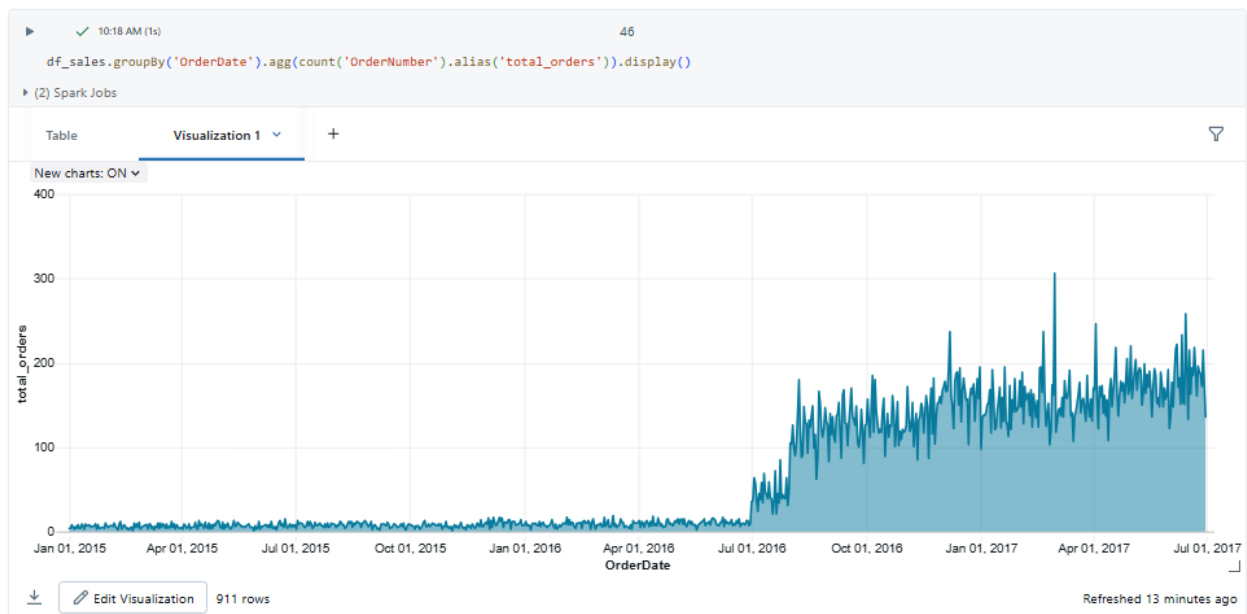
```
df_sales.display()
```

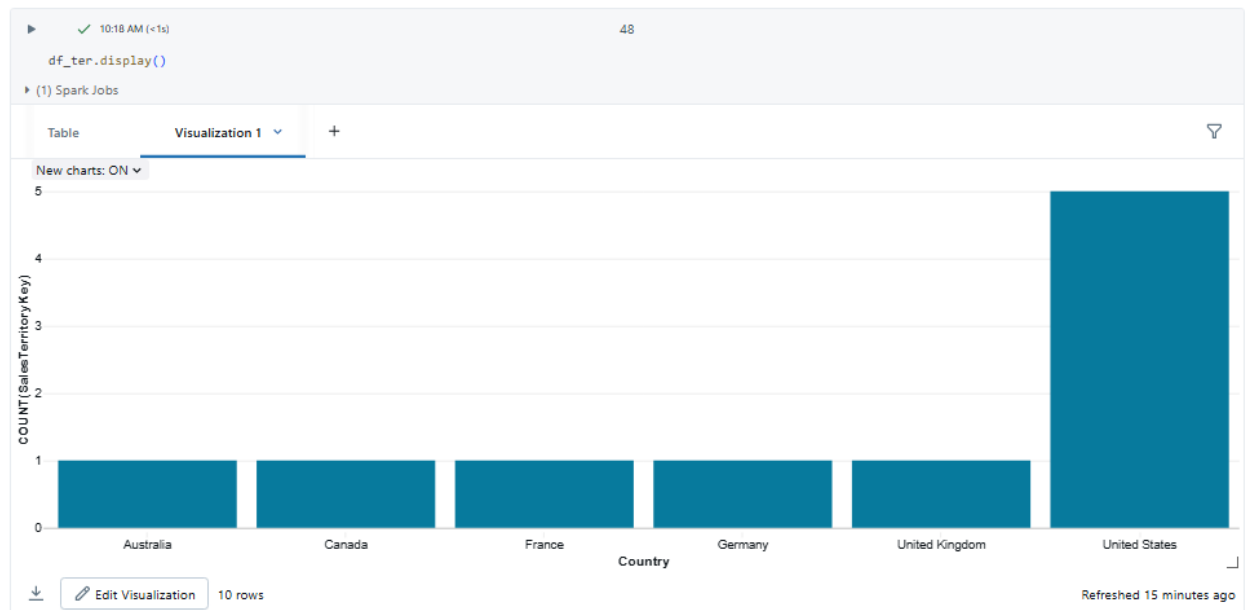
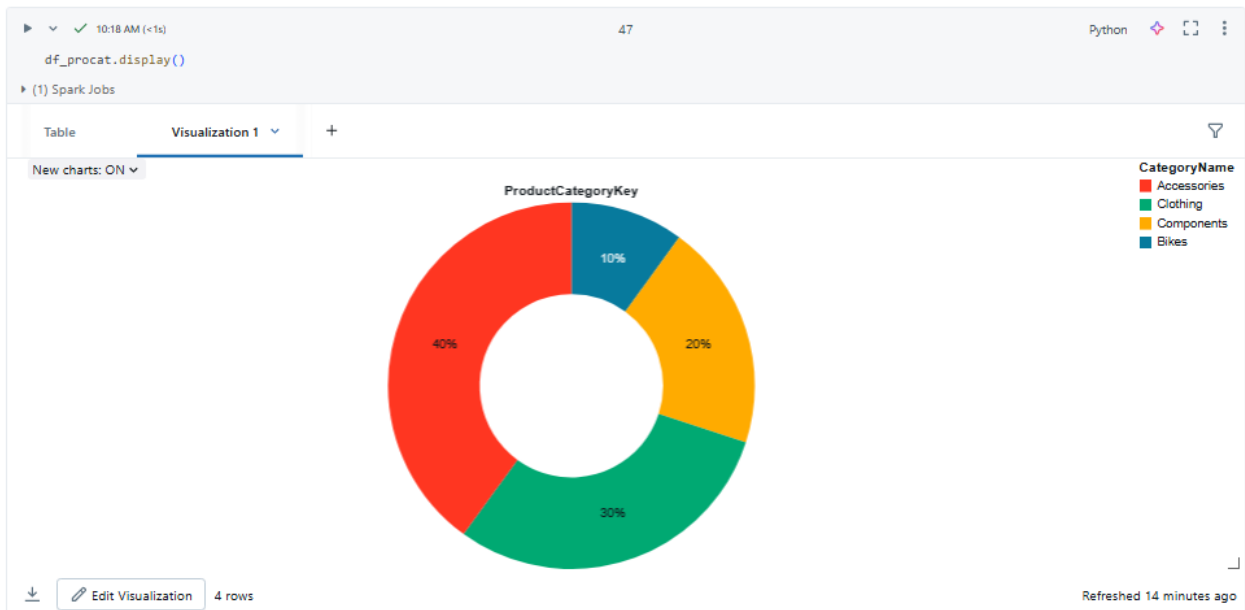
▶ (1) Spark Jobs

	📅 OrderDate	📅 StockDate	Ⓐ OrderN...	1 ² ₃ ProductKey	1 ² ₃ CustomerKey	1 ² ₃ TerritoryKey	1 ² ₃ C
1	2017-01-01	2003-12-13T00:00:00.000+00:...	TO61285	529	23791	1	
2	2017-01-01	2003-09-24T00:00:00.000+00:...	TO61285	214	23791	1	
3	2017-01-01	2003-09-04T00:00:00.000+00:...	TO61285	540	23791	1	
4	2017-01-01	2003-09-28T00:00:00.000+00:...	TO61301	529	16747	1	
5	2017-01-01	2003-10-21T00:00:00.000+00:...	TO61301	377	16747	1	

```
▶ 10:18 AM (2s) 44
df_sales.write.format('parquet')\
  .mode('append')\
  .save('/mnt/dataenge2edevstore/silver/AdventureWorks_Sales')
▶ (1) Spark Jobs
```

Sales Analysis





Home > dataeng2edevstore | Containers >

silver
Container

Search

+ Add Directory Upload Refresh Delete Copy Paste Rename Acquire lease Break lease Edit columns

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

silver

Authentication method: Access key ([Switch to Microsoft Entra user account](#))

Search blobs by prefix (case-sensitive)

Only show active objects

Showing all 8 items

<input type="checkbox"/>	Name	Last modified	Access tier	Blob type	Size	Lease state
<input type="checkbox"/>	AdventureWorks_Calendar	8/12/2025, 10:18:30 AM				...
<input type="checkbox"/>	AdventureWorks_Customers	8/12/2025, 10:18:34 AM				...
<input type="checkbox"/>	AdventureWorks_Products	8/12/2025, 10:18:37 AM				...
<input type="checkbox"/>	AdventureWorks>Returns	8/12/2025, 10:18:39 AM				...
<input type="checkbox"/>	AdventureWorks_Sales	8/12/2025, 10:18:44 AM				...
<input type="checkbox"/>	AdventureWorks_Subcategories	8/12/2025, 10:18:35 AM				...
<input type="checkbox"/>	AdventureWorks_Territories	8/12/2025, 10:18:41 AM				...
<input type="checkbox"/>	azuretempfolder\$	8/12/2025, 10:17:34 AM				...

Synapse Analytics workspace

dataeng-e2e-dev-synapse

New

Ingest
Perform a one-time or scheduled data load.

Explore and analyze
Learn how to get insights from your data.

Visualize
Build interactive reports with Power BI capabilities.

Discover more

Knowledge center

Browse partners

Synapse live Validate all Publish all

Develop

Filter resources by name

- SQL scripts 4
 - Create External Table
 - Create Schema
 - Create Views Gold
 - SQL script 1

Create External Table x Create Schema Create Views Gold SQL script 1

Run Undo Publish Query plan Connect to Built-in Use database dataeng-e2e-dev-sqlldb

```
1 CREATE DATABASE SCOPED CREDENTIAL cred_dev
2 WITH
3     IDENTITY = 'Managed Identity'
4
5 CREATE EXTERNAL DATA SOURCE source_silver
6 WITH
7     (
8         LOCATION = 'https://dataeng2edevstore.blob.core.windows.net/silver',
9         CREDENTIAL = cred_dev
10     )
11
```

Synapse live Validate all Publish all

Develop

Filter resources by name

SQL scripts 4

- Create External Table
- Create Schema
- Create Views Gold
- SQL script 1

Create External Table x Create Schema Create Views Gold SQL script 1

Run Undo Publish Query plan Connect to Built-in Use database dataeng-e2e-dev-sqldb

```
12 CREATE EXTERNAL DATA SOURCE source_gold
13 WITH
14 (
15     LOCATION = 'https://dataenge2edevstore.blob.core.windows.net/gold',
16     CREDENTIAL = cred_dev
17 )
18
19 CREATE EXTERNAL FILE FORMAT format_parquet
20 WITH
21 (
22     FORMAT_TYPE = PARQUET,
23     DATA_COMPRESSION = 'org.apache.hadoop.io.compress.SnappyCodec'
24 )
25
26 CREATE EXTERNAL TABLE gold.extsales
27 WITH
28 (
29     LOCATION = 'extsales',
30     DATA_SOURCE = source_gold,
31     FILE_FORMAT = format_parquet
32 )
33 AS
34 SELECT * FROM gold.sales
35
36 SELECT * FROM gold.extsales
```

Home > dataenge2edevstore | Containers >

gold
Container

Search x <

+ Add Directory Upload Refresh Delete Copy Paste Rename Acquire lease Break lease Edit columns

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

gold

Authentication method: Access key (Switch to Microsoft Entra user account)

Search blobs by prefix (case-sensitive)

Only show active objects

Showing all 1 items

<input type="checkbox"/>	Name	Last modified	Access tier	Blob type	Size	Lease state
<input type="checkbox"/>	extsales	8/12/2025, 9:11:07 AM				

