



Nawatech

Beyond Engineering

nawatech.co



KAI Workshop– Day 1

Prepared for: PT Kereta Api Indonesia



Explore end-to-end analytics
with Microsoft Fabric





Microsoft Fabric



**Data
Integration**
Data Factory



**Data
Engineering**
Synapse



**Data
Warehouse**
Synapse



**Data
Science**
Synapse



**Real Time
Analytics**
Synapse



**Business
Intelligence**
Power BI



Observability
Data Activator

- **Data Factory:** Data integration combining Power Query with the scale of Azure Data Factory to move and transform data.
- **Synapse Data Engineering:** Data engineering with a Spark platform for data transformation at scale.
- **Synapse Data Warehouse:** Data warehousing with industry-leading SQL performance and scale to support data use.
- **Synapse Data Science:** Data science with Azure Machine Learning and Spark for model training and execution tracking in a scalable environment.
- **Synapse Real-Time Analytics:** Real-time analytics to query and analyze large volumes of data in real-time.
- **Power BI:** Business intelligence for translating data to decisions.
- **Data Activator:** Real-time detection and monitoring of data that can trigger notifications and actions when it finds specified patterns in data.



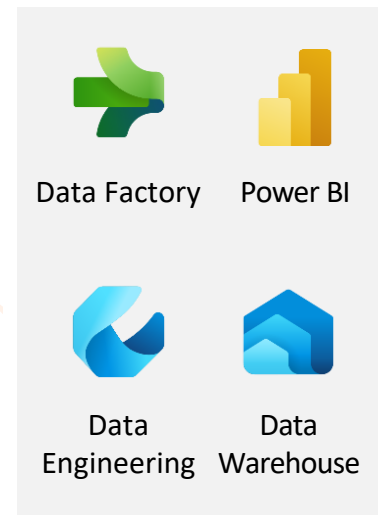
Data teams and Fabric

Fabric's unified management and governance make it easier for data professionals to work together.

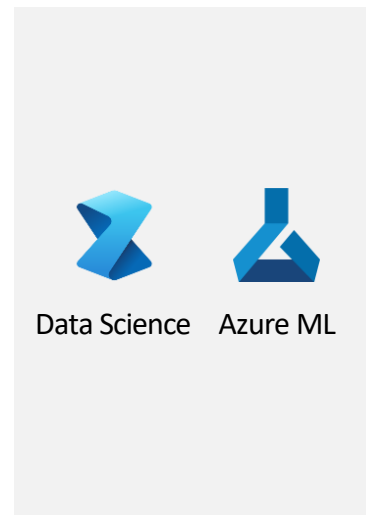
Data Engineers



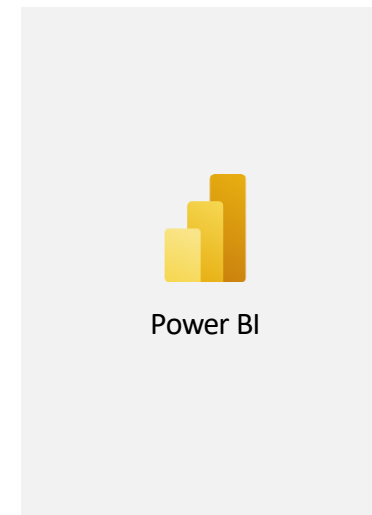
Analytics Engineers



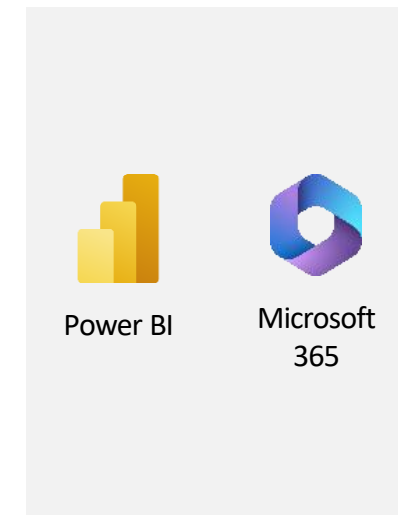
Data Scientists



Data Analysts



Decision Makers



Data Stewards

Get started with lakehouses in Microsoft Fabric



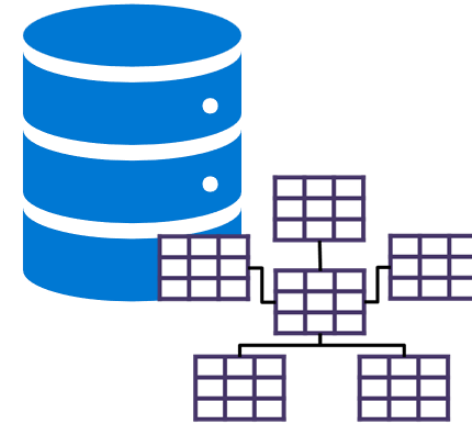


What is a lakehouse?



Data Lake

- Scalable, distributed file storage
- Flexible schema-on-read semantics
- Big data technology compatibility



Data Warehouse

- Relational schema modeling
- SQL-based querying
- Proven basis for reporting and analytics



Work with a Fabric lakehouse

Home

Create

Browse

OneLake data hub

Monitoring hub

Workspaces

DP-601

DP601_Bronze

Data Engineering

DP-601

+ New

Upload

Create app

Manage access

Details for DP601_Bronze

Name	Type	Owner
DP601_Bronze	Dataset (default)	
DP601_Bronze	SQL endpoint	
DP601_Bronze	Lakehouse	

Home

Create

Browse

OneLake data hub

Monitoring hub

Workspaces

DP-601

DP601_Bronze

DP601_Bronze

Data Engineering

File

Refresh

Create a report

Analyze in Excel

Lineage

Open data model

Details for DP601_Bronze

Location

DP-601

Refreshed

7/5/23, 4:45:00 PM

Sensitivity

Confidential\Microsoft Extended

Visualize this data

Create an interactive report, or a table, to discover and share business insights. [Learn more](#)

+ Create a report

Share this data

Give people access to the dataset and set their permissions to work with it. [Learn more](#)

Share dataset

See what already exists

These items use the same data source as DP601_Bronze.

Filter by keyword

Filter

Name	Type	Relation	Location	Refreshed	Endorsement	Sensitivity
DP601_Bronze	SQL endpoint	Upstream	DP-601	—	—	Confidential\Microsoft Extend...
DP601_Bronze	Lakehouse	Upstream	DP-601	—	—	Confidential\Microsoft Extend...



Explore, transform, and visualize data in the lakehouse

Tools and techniques to explore and transform data:

- ☆ Apache Spark
 - 📄 Notebooks
 - 📅 Spark Job Definitions
- 🏠 SQL analytics endpoint
- 🔗 Dataflows Gen2
- 📺 Data Pipelines

Visualize lakehouse data using Power BI

Use Apache Spark in Microsoft Fabric






Run Spark in Fabric


- To edit and run Spark code in Fabric, use *notebooks* or create a *Spark Job Definition*

Notebook (Preview)

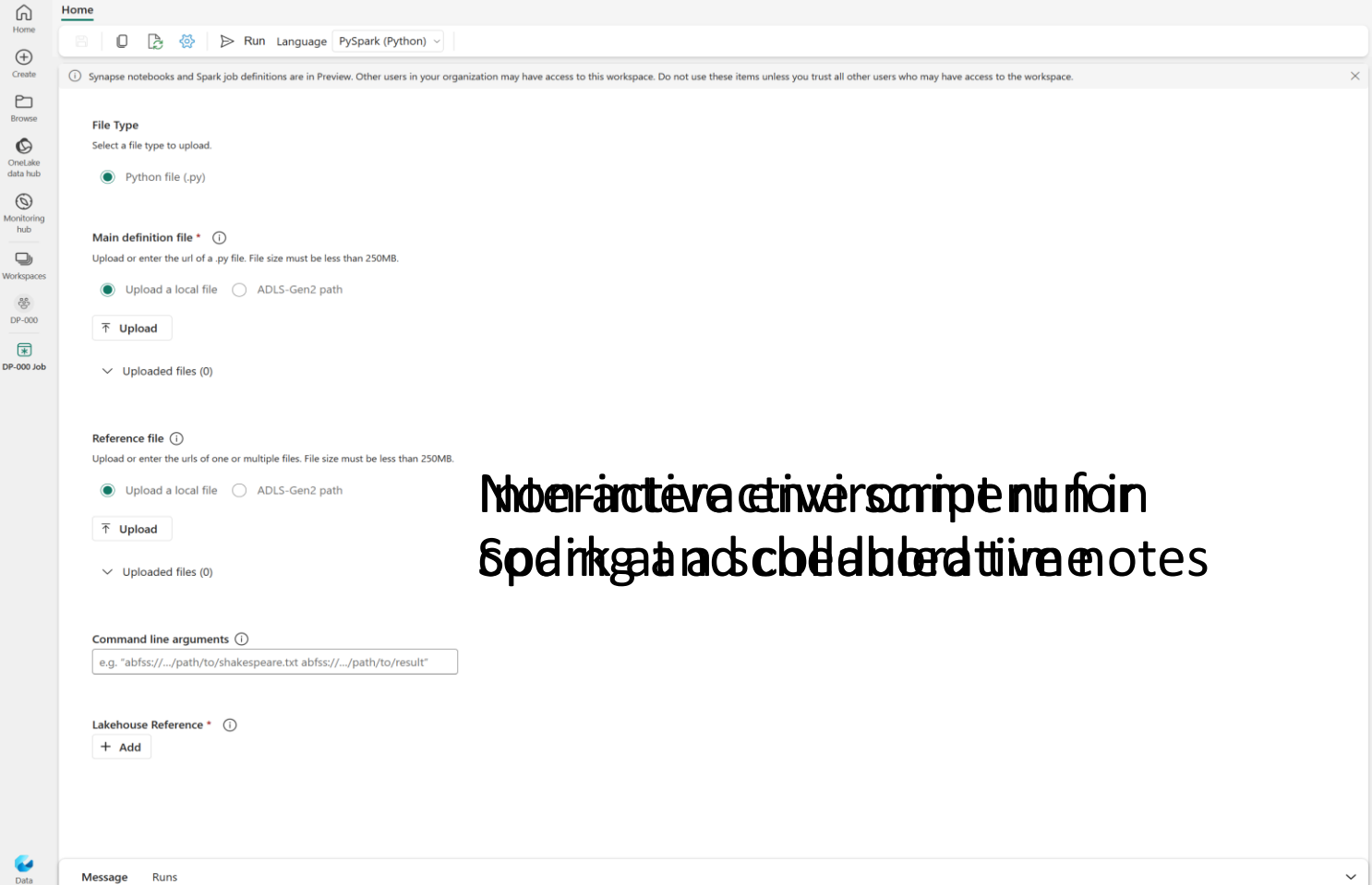


Explore data and build machine learning solutions with Apache Spark applications.

Spark Job Definition (Preview)



Define, schedule, and manage your Apache Spark jobs for big data processing.



The screenshot shows the 'Home' page of the Databricks Spark Job Definition interface. The left sidebar contains navigation links: Home, Create, Browse, OneLake data hub, Monitoring hub, Workspaces, DP-000, and DP-000 Job. The main content area is titled 'File Type' and 'Select a file type to upload.' It features a 'Python file (.py)' option. Below this, the 'Main definition file' section allows uploading a local file or an ADLS-Gen2 path. The 'Reference file' section also allows uploading a local file or an ADLS-Gen2 path. The 'Command line arguments' section has a text input field with a placeholder example. The 'Lakehouse Reference' section has an 'Add' button. The bottom of the interface shows a 'Message' and 'Runs' tab.

Interactive environment for
Spark and collaborative notes



Load data in a Spark Dataframe

- Schema can be either inferred or specified

Specify a schema

```
1 from pyspark.sql.types import *
2 from pyspark.sql.functions import *
3
4 productSchema = StructType([
5     StructField("ProductID", IntegerType()),
6     StructField("ProductName", StringType()),
7     StructField("Category", StringType()),
8     StructField("ListPrice", FloatType())
9 ])
10
11 df = spark.read.load('Files/data/product-data.csv',
12     format='csv',
13     schema=productSchema,
14     header=False)
15 display(df.limit(3))
```

1 sec - Command executed in 950 ms by Shannon Lindsay (SHE/HER) on 1:22:24 PM, 7/06/23

Spark jobs (1 of 1 succeeded) Log

Table	Chart	Export results		
Index	ProductID	ProductName	Category	ListPrice
1	771	Mountain-100 Silver, 38	Mountain Bikes	3399.99
2	772	Mountain-100 Silver, 42	Mountain Bikes	3399.99
3	773	Mountain-100 Silver, 44	Mountain Bikes	3399.99

Session ready Save option: Automatic 1 Selected 1 of 4 cells



Transform data in a Spark Dataframe

The screenshot shows the Databricks workspace interface. On the left is the 'Lakehouse explorer' sidebar with 'Files' selected. The main area displays a notebook cell with the following code and execution results:

```
1 bikes_df.write.mode("overwrite").parquet('Files/product_data/bikes.parquet')
```

The execution log shows a successful Spark job:

ID	Description	Status	Stages	Tasks	Duration
Job 30	parquet at NativeMethodAccessorImpl.java:0	✓ Succeeded	1/1	1/1 succeeded	1 sec

A large grey arrow points from the text 'Save dataframe for further analysis' to the code cell.

Save dataframe for further analysis

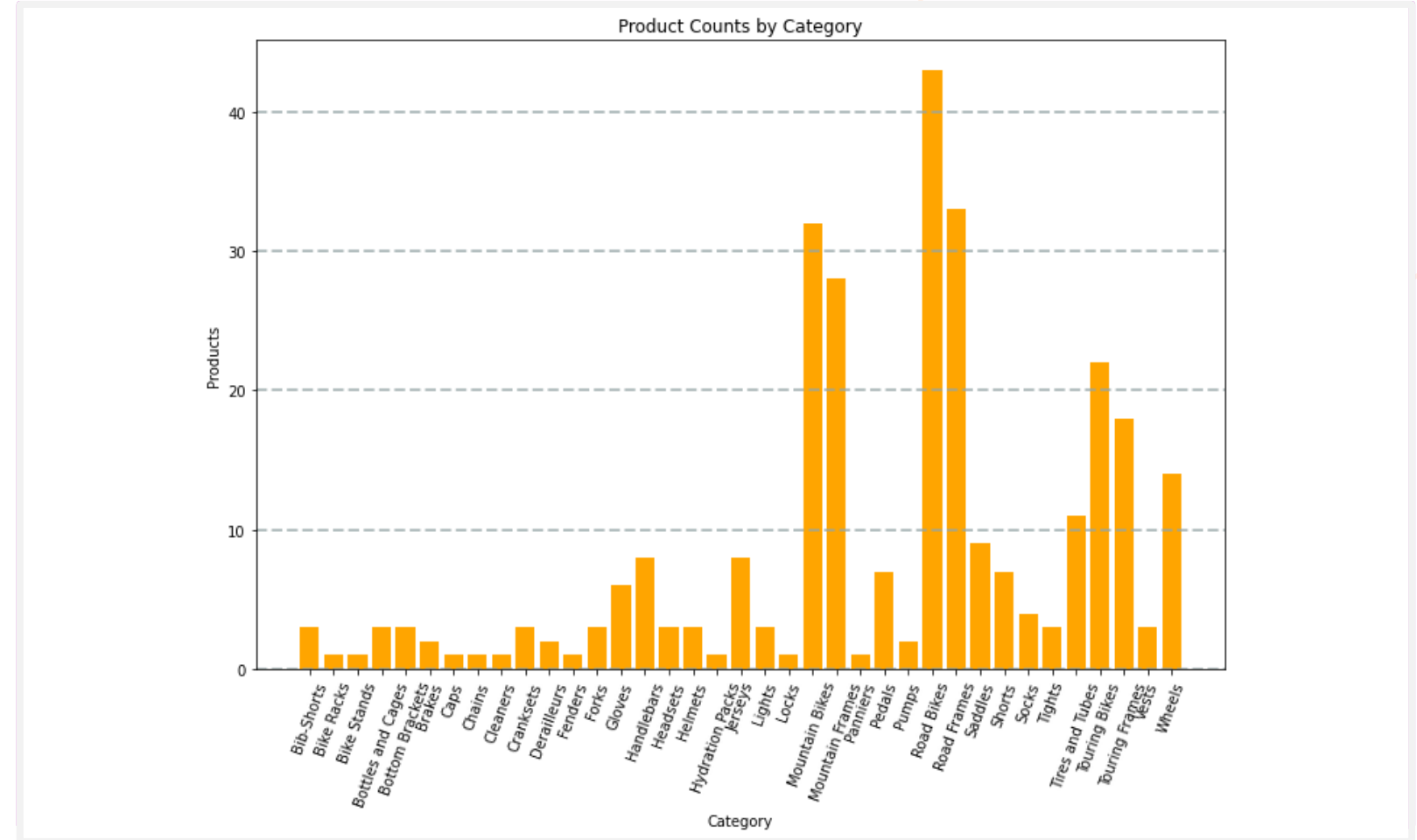
Copyright Nawatech 2024



Visualize data

- Two ways to visualize data in notebooks:

1. Use built-in notebook charts
2. Use graphics packages in code
 - This example uses Matplotlib





Hands-on: Apache Spark

Access GitHub : <https://s.id/qogLM>

Run Spark

Load Data

Transform Data

Visualize Data

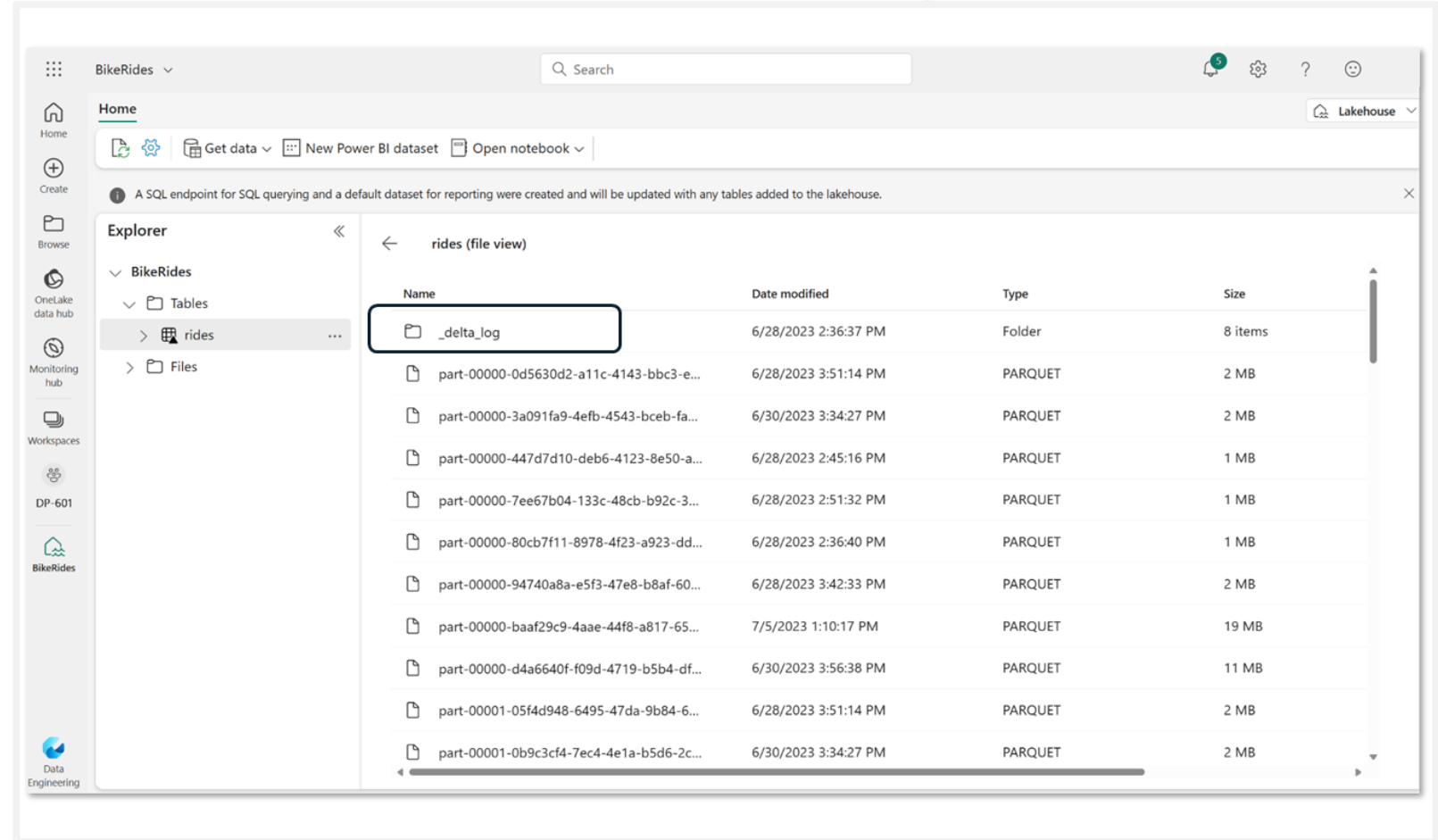
Work with Delta Lake tables in Microsoft Fabric





Understand Delta Lake

- Relational tables that support querying and data modification
- Support for ACID transactions
- Data versioning and time travel
- Standard formats and interoperability



Create Delta tables using code in Spark



1. Save a dataframe as a managed table

```
# Load a file into a dataframe
df = spark.read.load('Files/mydata.csv', format='csv', header=True)

# Save the dataframe as a delta table
df.write.format("delta").saveAsTable("mytable")
```

2. Use Spark SQL

```
%%sql

CREATE TABLE salesorders
(
  Orderid INT NOT NULL,
  OrderDate TIMESTAMP NOT NULL,
  CustomerName STRING,
  SalesTotal FLOAT NOT NULL
)
USING DELTA
```

3. Save a dataframe in delta format in an explicit path

```
delta_path = "Files/mydatatable"
df.write.format("delta").save(delta_path)
```

Managed vs external tables



Managed tables

- Defined without a specific location – Files are created in the default metastore folder (*Tables/...*)
- Dropping the table deletes the files

```
# Save a dataframe as a delta table  
df.write.format("delta").saveAsTable("mytable")
```

External tables

- Defined with an explicit file location outside of the default metastore folder
- Dropping the table does not delete the files

```
df.write.format("delta").saveAsTable("myexternaltable", path="Files/myexternaltable")
```

Work with Delta tables in Spark



1. Use Spark SQL to embed a SQL statement in PySpark

- Embed SQL statements in other languages using the **spark.sql** library.

```
spark.sql("INSERT INTO products VALUES (1, 'Widget', 'Accessories', 2.99)")
```

2. Native Spark SQL using %%sql magic

- Use %%sql magic in a notebook to run SQL statements.

```
%%sql
```

```
UPDATE products  
SET Price = 2.49 WHERE ProductId = 1;
```

3. Use the Delta API

- Create an instance of a DeltaTable from a folder location containing files in delta format, and then use the API to modify the data in the table.

```
from delta.tables import *  
from pyspark.sql.functions import *  
  
# Create a DeltaTable object  
delta_path = "Files/mytable"  
deltaTable = DeltaTable.forPath(spark, delta_path)  
  
# Update the table (reduce price of accessories by 10%)  
deltaTable.update(  
    condition = "Category == 'Accessories'",  
    set = { "Price": "Price * 0.9" })
```



Hands-on: Delta Lake Tables

Access GitHub : <https://s.id/qogLM>

Create Delta Table in Spark

Save as Managed and External Tables

Work with Delta Table in Spark

Organize a lakehouse with
medallion architecture





Medallion architecture overview

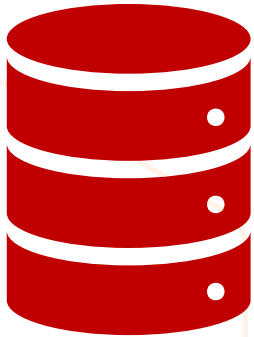




Implement a medallion architecture

- Bronze

Ingest raw data



Pipelines, dataflows, notebooks

- Silver

Cleanse and validate data



Dataflows or notebooks

- Gold

Additional transformations and modeling



SQL analytics endpoint or semantic model



Query and report on data in your Fabric lakehouse

Sales (1) Search Fabric Trial: 59 days left SQL endpoint

Home New SQL query New visual query

A default dataset for faster reporting was created and will be automatically updated with any tables and views added to the lakehouse. [Learn more](#) Manage default Power BI dataset

Explorer

- Warehouses
- Sales
 - Schemas
 - dbo
 - Tables
 - dimcustomer...
 - dimdate_gold
 - dimproduct...
 - factsales_gold
 - sales_silver
 - Views
 - Functions
 - Stored Procedur...
 - guest
 - INFORMATION_SCHE...
 - sys
 - Security
 - Queries
 - My queries
 - SQL query 1
 - Visual query 1
 - Shared queries

Visual query 1 SQL query 1

Run Save as view

```
1 select CustomerName, SUM(Quantity) AS TotalQuantity
2 FROM [sales_silver]
3 GROUP BY CustomerName
4 ORDER BY TotalQuantity DESC
```

Messages Results Download Excel file Visualize results Limited to 10,000 rows Search

	ABC CustomerName	123 TotalQuantity
1	Samantha Jenkins	41
2	Henry Garcia	39
3	Charles Jackson	38
4	April Shan	35
5	Hailey Patterson	34
6	Ryan Thompson	34
7	Mason Roberts	33
8	Dalton Perez	32
9	Jason Griffin	31
10	Fernando Barnes	28
11	Jasmine Powell	27
12	Nicholas Brown	25
13	Ashley Henderson	25
14	Jennifer Simmons	24
15	Gina Martin	24
16	Ana Perry	24
17	Alexandra Jenkins	22

Succeeded (2 sec 671 ms) Columns: 2 Rows: 10,000

Data Query Model



Secure the medallion layers

- **Gold Layer Access Control:** Read-only.
- **Silver Layer Utilization:** Build optional.
- **Bronze Layer Access Control:** Read only.



Thank you
Let's discuss our collaborations

PT Nawa Darsana Teknologi
Gedung Office 8, 18th Floor, Jl. Jend Sudirman Kav 52-53
SCBD Lot 8, Senopati Jakarta Selatan 12190 DKI Jakarta, Indonesia
Website : www.nawatech.co
Phone : +62 21 29552754

Copyright Nawatech 2024

Contact

Innes Sabrina Husa
Senior Business Development
Manager
Email : windy@nawatech.co
Phone : +62 812 8033 2457