



Leveraging Microsoft Fabric for Advanced Data Engineering Solutions





Philip Goldman

Field Sr. Data Solution Architect
EPAM

Introducing Microsoft Fabric



```
graph TD; A[Introducing Microsoft Fabric] --> B[Deep Dive into Data Engineering]; B --> C[DEMO];
```

Deep Dive into Data Engineering

DEMO

Introducing Microsoft Fabric for Data Engineering

Microsoft Fabric

End-to-end analytics data fabric
From the data lake to the business user

Complete Analytics Platform

Everything, unified

Unified SaaS Solution

Low Code Plus Pro Dev

Lake-centric and Open

OneLake

One Copy

Always Synced

Empower Every Office User

Familiar and Intuitive

Built Into Office 365

Insight to Action

Persistent Security and Governance

End-to-End Visibility

Always Governed

Secure by Default

Introducing Microsoft Fabric for Data Engineering

- **Complete analytics platform**

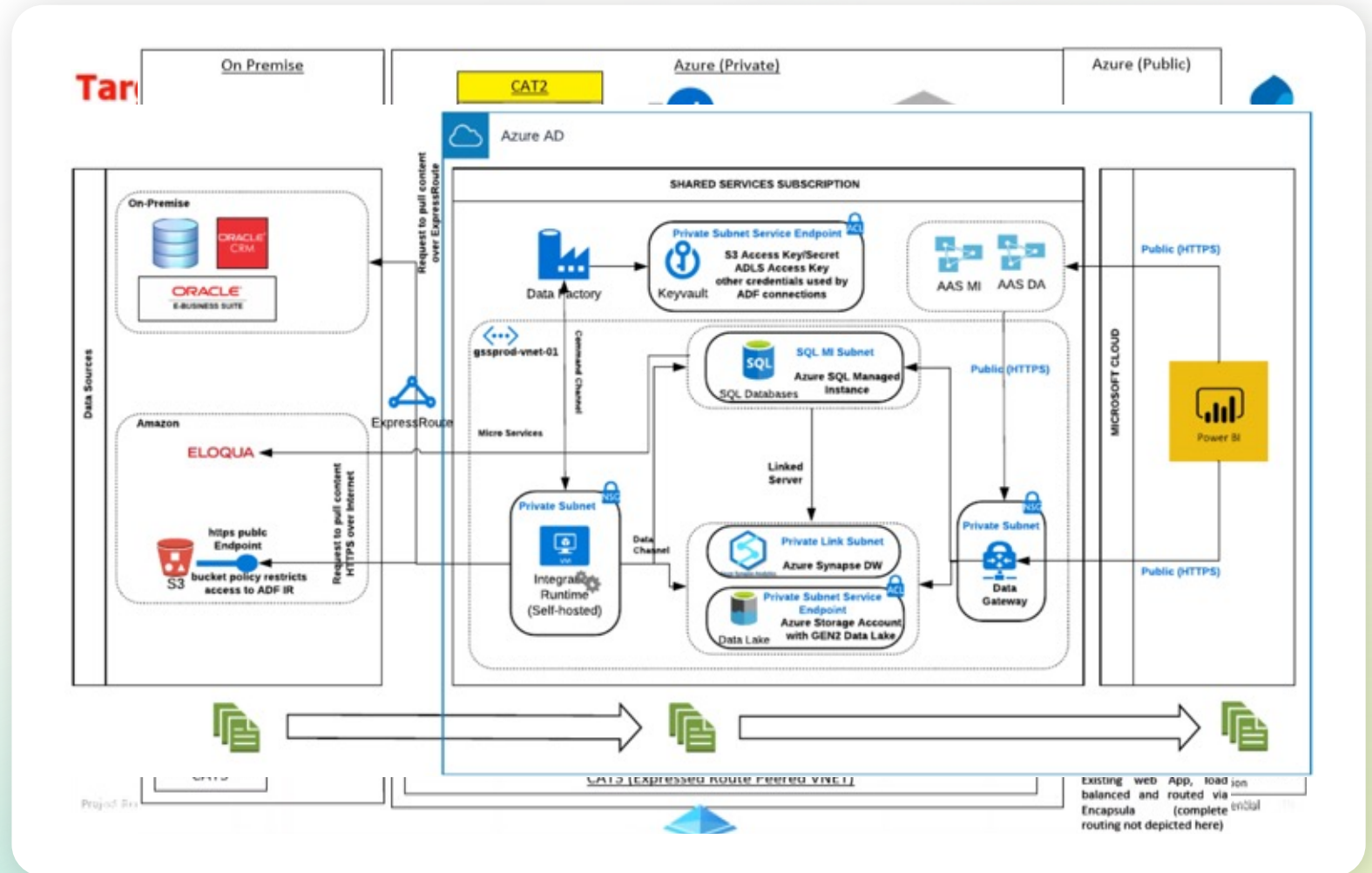
Scalable analytics are complex and fragmented

Every analytics project has many subsystems

Every subsystem need a different class of product

Products often come from multiple vendors

Integration at scale across products is complex, fragile, and expensive



Scalable analytics are complex and fragmented

Every analytics project
has many subsystems

Every subsystem need a
different class of product

Products often come
from multiple vendors

**Integration at scale across
products is complex,
fragile and expensive**

//

Simplify,

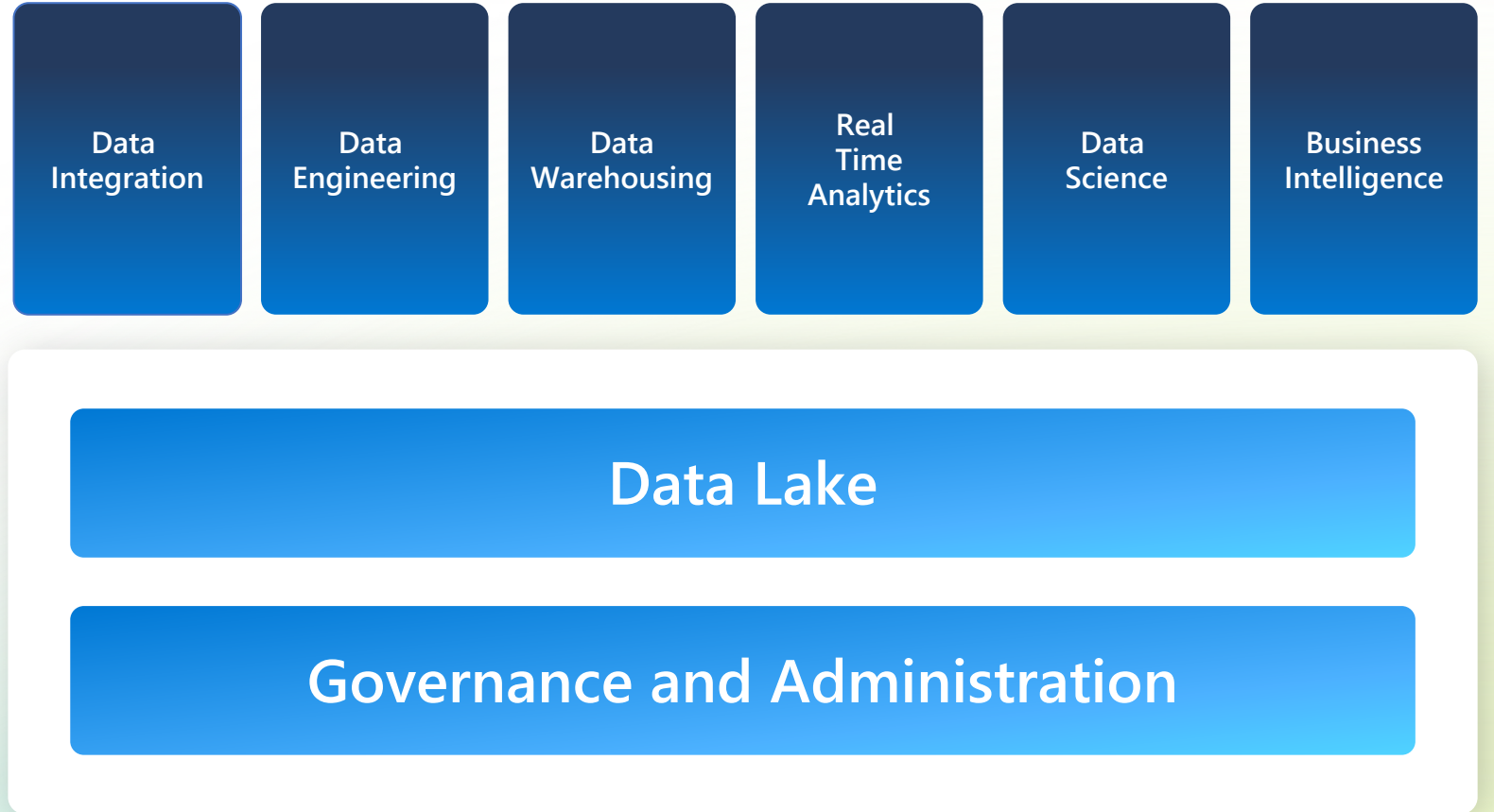
I am the Chief Data Officer
and don't want to be the
Chief Integration Officer.”

Every CDO, Every Enterprise

A silver lining?

Analytics systems have very predictable patterns

Microsoft has all the products with the right scale needed to build a complete analytics system



A silver lining?

Analytics systems have very predictable patterns

Microsoft has all the products with the right scale needed to build a complete analytics system



Still far too complex

Many Products

Different Experiences

Proprietary and Open

Dedicated and Serverless

PaaS and SaaS

Different Business Models

Steep Learning Curves

Deep Expertise Needed

High Integration Effort



Power BI



Synapse



Kusto



Azure AI



Data Factory



Spark

Microsoft Fabric



Data
Integration



Data
Lake



Spark
Engines



Data
Warehouse



Real Time
Analytics



Data
Science



Business
Intelligence

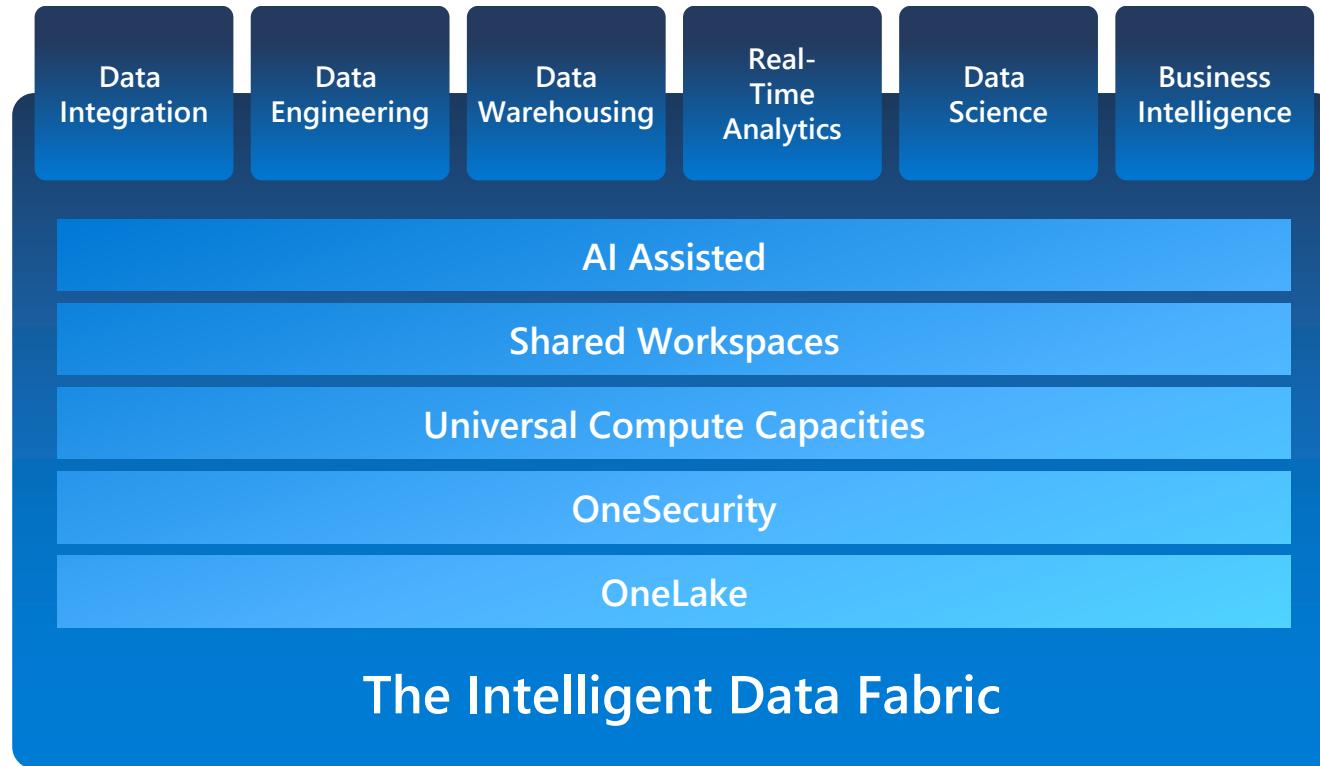


Governance

Unified analytics fabric

End-to-end analytics data fabric
From the data lake to the business user

Microsoft Fabric



Single...

- Onboarding and trials
- Sign-on
- Navigation model
- UX model
- Workspace organization
- Collaboration experience
- Data Lake
- Storage format
- Data copy for all engines
- Security model
- CI/CD
- Monitoring hub
- Data Hub
- Governance & compliance

Persona Centric Experiences

The screenshot displays the Microsoft Fabric user interface. At the top, there is a navigation bar with the text "Microsoft Fabric Home" and a user profile icon. Below this, the main heading "Microsoft Fabric" is centered, followed by the tagline "All your data. In one location. Organize. Collaborate. Create." and the sub-heading "Explore the experience".

The dashboard features eight service tiles arranged in a 2x4 grid:

- Power BI**: Find insights, track progress, and make decisions faster using rich visualizations.
- Data Factory**: Solve the most complex data integration and ETL scenarios with cloud-scale data movement and data transformation services.
- Data Activator**: Monitor data to trigger alerts and automated actions so your organization adapts to changing conditions in real time.
- Industry Solutions**: Get a head start with industry-relevant connectors, transformations, and scenario-specific tools.
- Synapse Data Engineering**: Create a lakehouse, and use Apache Spark to transform and prepare organizational data to share with the business.
- Synapse Data Science**: Explore your data, and build machine learning models to infuse predictive insights into your analytics solutions and applications.
- Synapse Data Warehouse**: Scale up your insights by storing and analyzing data in a secure, open-data-format SQL warehouse with top performance at PB scale.
- Synapse Real-Time Analytics**: Rapidly ingest, transform, and query any data source and format, from 1 GB to 1 PB, and then visualize and share the insights.

At the bottom of the dashboard, there are two additional action tiles:

- Read documentation**: Represented by a document icon.
- Explore community**: Represented by a group of people icon.

The Microsoft Fabric logo is visible in the bottom-left corner of the interface.

Data Integration

The screenshot displays the Microsoft Fabric Data Factory interface. At the top, the breadcrumb navigation shows 'pipeline2' and 'EPAM Confidential\Confidential'. The main navigation bar includes 'Home', 'Activities', 'Run', and 'View'. Below this, a toolbar contains icons for 'Validate', 'Run', 'Schedule', 'View run history', 'Copy data', 'Dataflow', 'Notebook', 'Lookup', and 'Invoke pipeline'. The left sidebar shows navigation options: 'Home', 'Create', 'Browse', 'OneLake data hub', 'Monitoring hub', 'Workspaces', 'Fabric Demo', and 'pipeline2'. The bottom left corner features the 'Data Factory' logo.

The central focus is the 'Copy data' wizard, which is currently on the 'Choose data source' step. The wizard's progress bar on the left shows five steps: 'Choose data source' (selected), 'Connect to data source', 'Choose data destination', 'Connect to data destination', and 'Review + save'. The 'Choose data source' step is titled 'Choose data source' and includes the instruction: 'Select a connector. Then enter the connection information.'

The 'Choose data source' panel is titled 'Choose data source' and contains the following content:

- Build your data ingestion task to move objects from a data source to a data destination. [Learn more](#)**
- Sample data**
- COVID-19 Data Lake**
Varied per format (CSV, JSON, JSON Lines, Parquet)
The COVID-19 Data Lake contains COVID-19 related datasets from various sources. It covers testing and patient outcome tracking data, social...
- NYC Taxi - Green**
2 GB (Parquet)
The green taxi trip records include fields capturing pick-up and drop-off dates/times, pick-up and drop-off locations, trip distances, itemized fares,...
- Diabetes**
14 KB (Parquet)
The Diabetes dataset has 442 samples with 10 features, making it ideal for getting started with machine learning algorithms.
- Public Holidays**
500 KB (Parquet)
Worldwide public holiday data sourced from PyPI holidays package and Wikipedia, covering 38 countries or regions from 1970 to 2099.
- Retail Data Model from Wide World Importers**
352MB (Parquet)
Wide World Importers (WWI) is a fictional novelty goods importer and distributor operating from the San Francisco Bay Area. WWI sells to retail...

Below the sample data, the 'Data sources' section is visible, featuring a search bar and a list of connectors categorized by 'All categories', 'Workspace', 'Azure', 'Database', 'File', 'Generic protocol', 'NoSQL', and 'Services and apps'. The connectors listed include:

- Amazon RDS for SQL Server Database
- Amazon Redshift Database
- Amazon S3 File
- Amazon S3 Compatible File
- Azure Blob Storage Azure
- Azure Cosmos DB for NoSQL Azure
- Azure Data Explorer Azure
- Azure Data Lake Storage Gen2 Azure
- Azure Database for PostgreSQL Azure
- Azure SQL Database
- Azure SQL Database Managed Instance
- Azure Synapse Analytics

At the bottom right of the wizard, there are 'Back' and 'Next' buttons.

Data Engineering

The screenshot displays the Microsoft Fabric Lakehouse interface. The top navigation bar shows the workspace name 'wwilakehouse' and a search bar. Below the navigation bar, there are several tabs and buttons, including 'Home', 'Get data', 'New semantic model', 'Open notebook', and 'Manage OneLake data access'. The main content area is divided into three sections: a left sidebar with an 'Explorer' view showing a tree of tables and files, a central 'Code' editor with a PySpark notebook, and a right 'Visuals' pane.

The 'Code' editor shows a PySpark notebook with the following code:

```
1 sampled_df['hour'] = sampled_df['tpepDropoffDateTime'].dt.hour
2 sampled_df['dayofweek'] = sampled_df['tpepDropoffDateTime'].dt.dayofweek
3 sampled_df['dayname'] = sampled_df['tpepDropoffDateTime'].dt.day_name()
4 sns.histplot(data=sampled_df, x='hour', stat='count', discrete=True, kde=True)
5 plt.title("Distribution by Hour of the day")
6 plt.xlabel('Hours')
7 plt.ylabel('Count of trips')
```

The 'Visuals' pane displays a histogram titled 'Distribution by Hour of the day'. The x-axis is labeled 'Hours' and ranges from 0 to 24. The y-axis is labeled 'Count of trips' and ranges from 0 to 3000. The histogram shows a distribution of trip counts across the 24 hours of the day, with a peak around 18-19 hours. A KDE curve is overlaid on the histogram.

Below the histogram, there is a caption: 'Visual 8: Analyze average taxi trip duration by hour and day of the week using a heatmap'.

The bottom of the interface shows a 'Data Engineering' logo and a status bar with 'Not connected' and 'AutoSave: On'.

Data Warehouse

The screenshot shows the Microsoft Fabric SQL analytics endpoint interface. The top navigation bar includes the user name 'wwilakehouse', a search bar, and a 'Fabric Trial: 54 days left' notification. The main area is divided into three sections: Explorer, SQL query editor, and Results.

Explorer: Shows a tree view of the 'wwilakehouse' data lake. Under 'Schemas', the 'dbo' schema is expanded to show 'Tables'. The following tables are listed: 'aggregate_sal...', 'aggregate_sal...', 'dimension_city', 'dimension_cu...', 'dimension_da...', 'dimension_e...', 'dimension_st...', 'fact_sale', and 'Views'. Other schemas like 'guest', 'INFORMATION_SCHE...', 'queryinsights', 'sys', and 'Security' are also visible.

SQL query editor: Contains the following SQL query:

```
1 SELECT BuyingGroup, Count(*) AS Total
2 FROM dimension_customer
3 GROUP BY BuyingGroup
4
```

Results: The query results are displayed in a table with 5 rows and 2 columns. The first column is 'BuyingGroup' and the second is 'Total'.

BuyingGroup	Total
Toddler Toys	23
N/A	1
Wingtip Toys	201
Kids Toys	15
Tailspin Toys	163

The status bar at the bottom indicates 'Succeeded (1 sec 713 ms)' and 'Columns: 2 Rows: 5'.

Real-Time Analytics

The screenshot displays the Microsoft Real-Time Analytics interface for a database named 'scaletestdxt'. The interface is organized into several sections:

- Header:** Includes the user profile 'Confidential\Microsoft...', a search bar, and navigation options like 'Home', 'Manage', 'Editing', and 'Share'.
- Object tree (Left Panel):** Shows a hierarchical view of the database structure, including 'Database', 'Tables', 'Logs', 'IoTDevices', 'FHV_Trips', 'Functions', 'Materialized views', 'Shortcuts', and 'DataStreams'.
- Database: scaletestdxt (Main Panel):**
 - Database details:** A summary card showing 'Created by: Amir Leshman', 'Region: EastUS2', 'Created on: 1/30/23, 08:49', 'Last ingestion Query URI', and 'Ingestion URI: OneLake folder'.
 - Size:** A card displaying '99.99 GB' (Compressed size), '402.46 GB' (Original size), and '4.03' (Compression ratio).
 - Top tables:** A table listing the largest tables in the database:

Name	Size
TaxiRides	401.82 GB
IoTDevices	657.95 MB
Logs	147.43 KB
FHV_Trips	0 Bytes
 - Most active users:** A table showing the top users by queries run last month:

Name	Queries run last month
tzgitin@microsoft.com	79
vladikb@microsoft.com	22
 - Recently updated functions:** A list of functions updated recently:

Function Name	Updated
MyFunction3	2/8/23, 16:31
MyFunction2	2/8/23, 16:31
MyFunction1	2/8/23, 16:31
 - Recently used Querysets:** A list of querysets used recently:

Queryset Name	Used
TaxiRides	Yesterday, 17h ago
ScaleQueryDemo	2/8/23, 18:03
 - Recently created data connections:** A list of data connections created recently:

Connection Name	Created
scaletestdxt-StocksData	Yesterday, 22h ago
scaletestdxt-liram-eh	1/30/23, 09:06

Data Science

Train Evaluate | EPAM Confidential\Confidential · Saved

Search

Fabric Trial: 54 days left

Home Edit Run View

Editing Comments Share

Run all Connect PySpark (Python) Workspace default Data Wrangler Copilot

Explorer

```
5 X_res, y_res = sm.fit_resample(X_train, y_train)
6 new_train = pd.concat([X_res, y_res], axis=1)
```

PySpark (Python)

Model Training

- Train the model using Random Forest with maximum depth of 4 and 4 features

```
1 mlflow.sklearn.autolog(registered_model_name='rfc1_sm') # Register the trained model with autologging
2 rfc1_sm = RandomForestClassifier(max_depth=4, max_features=4, min_samples_split=3, random_state=1) # Pass hyperparameters
3 with mlflow.start_run(run_name="rfc1_sm") as run:
4     rfc1_sm_run_id = run.info.run_id # Capture run_id for model prediction later
5     print("run_id: {}; status: {}".format(rfc1_sm_run_id, run.info.status))
6     # rfc1.fit(X_train,y_train) # Imbalanced training data
7     rfc1_sm.fit(X_res, y_res.ravel()) # Balanced training data
8     rfc1_sm.score(X_val, y_val)
9     y_pred = rfc1_sm.predict(X_val)
10    cr_rfc1_sm = classification_report(y_val, y_pred)
11    cm_rfc1_sm = confusion_matrix(y_val, y_pred)
12    roc_auc_rfc1_sm = roc_auc_score(y_res, rfc1_sm.predict_proba(X_res)[: , 1])
```

PySpark (Python)

- Train the model using Random Forest with maximum depth of 8 and 6 features

```
1 mlflow.sklearn.autolog(registered_model_name='rfc2_sm') # Register the trained model with autologging
2 rfc2_sm = RandomForestClassifier(max_depth=8, max_features=6, min_samples_split=3, random_state=1) # Pass hyperparameters
3 with mlflow.start_run(run_name="rfc2_sm") as run:
4     rfc2_sm_run_id = run.info.run_id # Capture run_id for model prediction later
5     print("run_id: {}; status: {}".format(rfc2_sm_run_id, run.info.status))
6     # rfc2.fit(X_train,y_train) # Imbalanced training data
7     rfc2_sm.fit(X_res, y_res.ravel()) # Balanced training data
8     rfc2_sm.score(X_val, y_val)
9     y_pred = rfc2_sm.predict(X_val)
10    cr_rfc2_sm = classification_report(y_val, y_pred)
11    cm_rfc2_sm = confusion_matrix(y_val, y_pred)
12    roc_auc_rfc2_sm = roc_auc_score(y_res, rfc2_sm.predict_proba(X_res)[: , 1])
```

PySpark (Python)

Selected Cell 1 of 42 cells

Not connected AutoSave: On

Power BI

The screenshot displays the Microsoft Power BI interface. At the top, the navigation bar includes 'Home', 'Manage', and 'Editing' options. The main content area is titled 'Contoso Daily Sales' and features a 'Sales Overview' dashboard. This dashboard includes four key metrics: Revenue Won (\$7,811,851), Close % (37.7%), AVG Days to Close (121), and Opportunities Won (526). Below these metrics are three charts: 'Revenue Won by Month' (line chart), 'Close % by Month' (bar chart), and 'Close % by Region' (map). A 'Copilot' chat window is open on the right, with a prompt to 'Create a report with Copilot' and a progress bar indicating 'Analyzing your data...'. The left sidebar shows the 'Object tree' with a 'Database' section containing tables like 'TaxiRides', 'Logs', 'IoTDevices', and 'FHV_Trips'. The bottom status bar shows 'Power BI' and the current report name 'Sales Overview'.

Sales Overview

Date: 06/01/2022 - 01/12/2023

Metric	Value
Revenue Won	\$7,811,851
Close %	37.7%
AVG Days to Close	121
Opportunities Won	526

Revenue Won by Month

Month	Revenue Won (\$M)
June 2022	~0.2
July 2022	~0.2
Aug 2022	~0.8
Sep 2022	~0.5
Oct 2022	~1.0
Nov 2022	~1.5
Dec 2022	~2.0
Jan 2023	~1.8

Close % by Month

Month	Close %
May	23%
Jun	32%
Jul	26%
Aug	37%
Sep	37%
Oct	38%
Nov	52%
Dec	49%
Jan	39%

Close % by Region

Map showing Close % by region across the United States.

AVG Days to Close by Month

Month	AVG Days to Close
May	116
Jun	118
Jul	119
Aug	125
Sep	121
Oct	122
Nov	123
Dec	120
Jan	121

Copilot Preview

Create a report with Copilot

Describe the report you want, in your own words, and Copilot will create it quickly.

Help me build a sales report summarizing our key metrics and trends

Sales overview page added

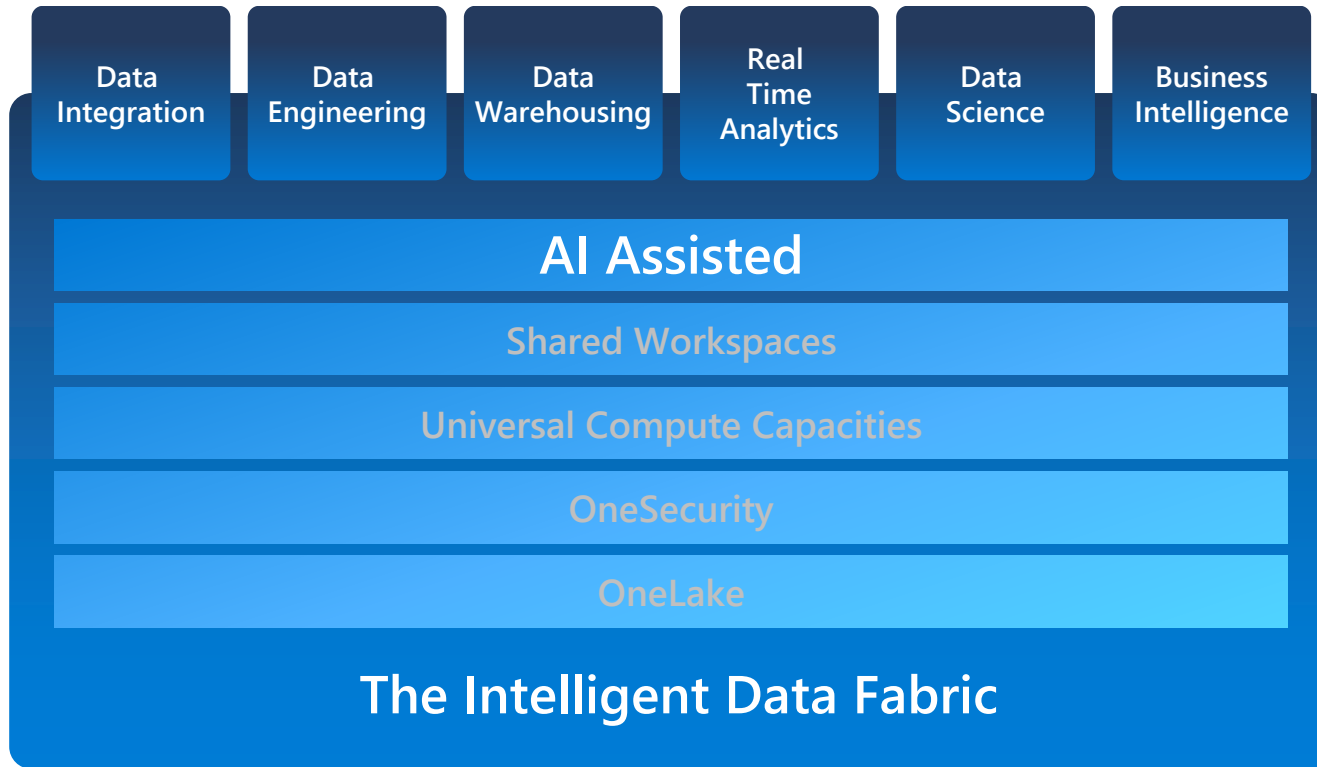
What are the biggest drivers for close %?

Analyzing your data...

Cancel

AI-generated content can have mistakes. Make sure it's accurate and appropriate before using it. [Read preview terms](#)

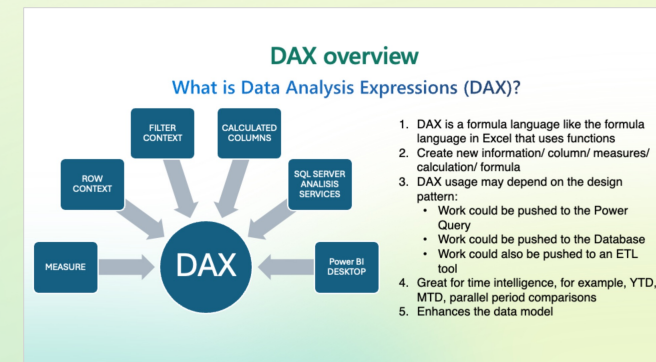
AI Assisted Creation in Microsoft Fabric



The Fabric platform will include built in Azure Open AI based assistant that will serve all the workloads

First GPT-based feature is already shipping in Power BI - NL2DAX – DAX calculation creation based on natural language prompts

Ongoing major ramp-up for pervasive AOAI based product-wide AI assistance

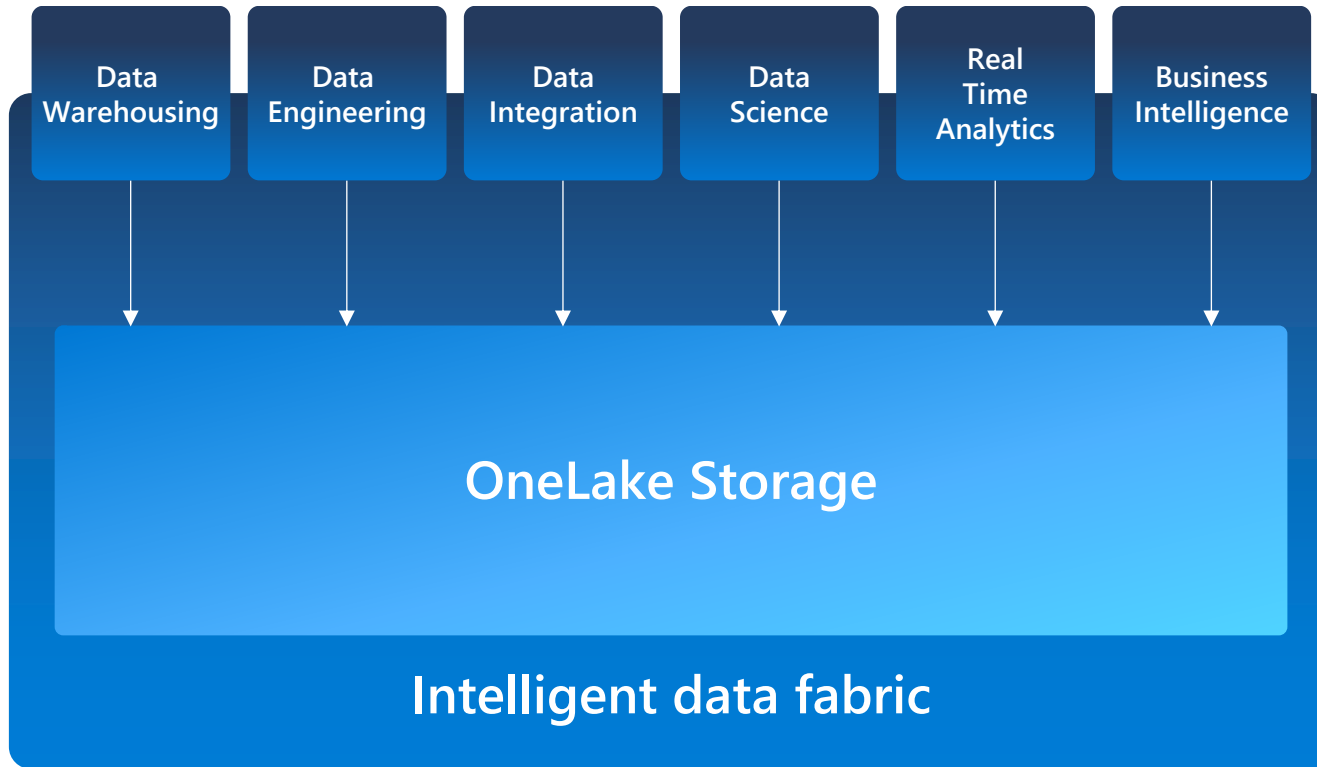


Introducing Microsoft Fabric for Data Engineering

- Complete analytics platform
- **Lake-centric and open architecture**

OneLake for all Data

“The OneDrive for Data”



A single SaaS lake for the whole organization

Provisioned automatically with the tenant

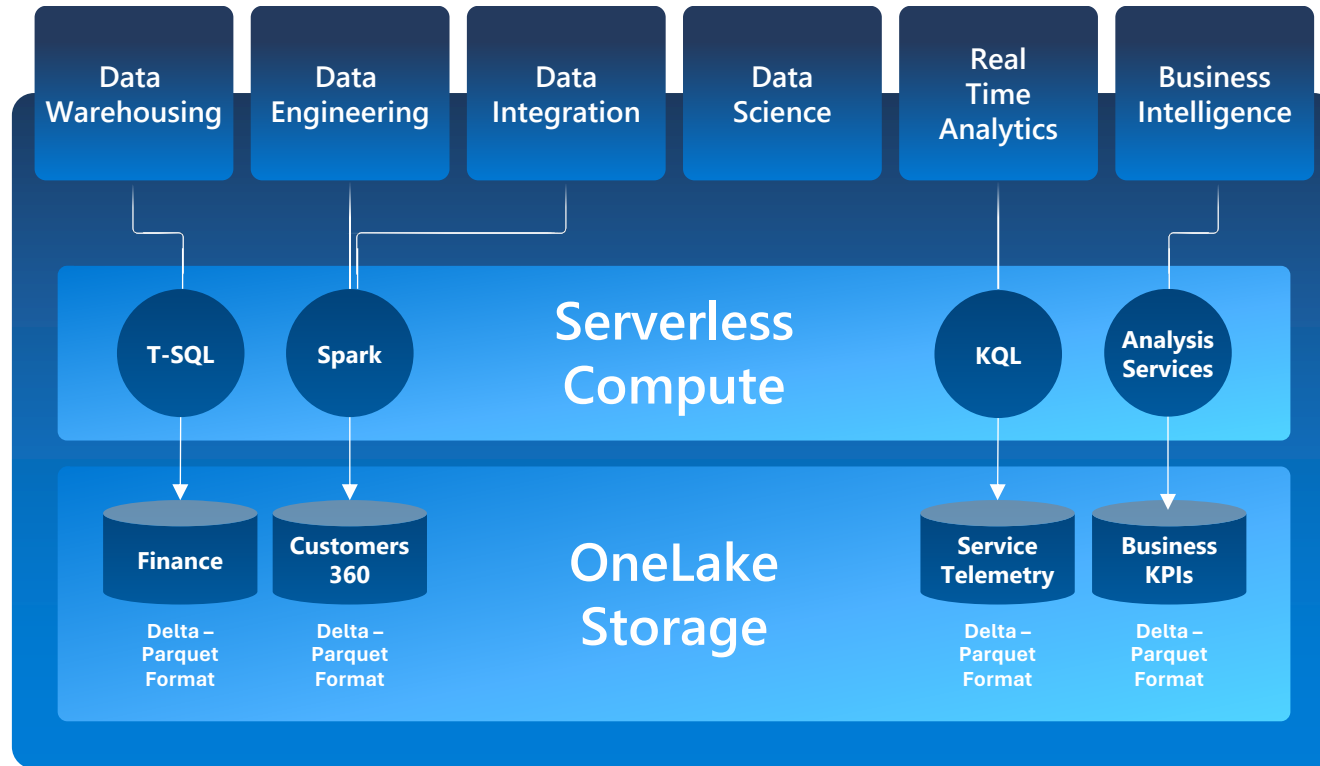
All workloads automatically store their data in the OneLake workspace folders

All the data is organized in an intuitive hierarchical namespace

The data in OneLake is automatically indexed for discovery, MIP labels, lineage, PII scans, sharing, governance and compliance

One Copy for all computes

Real separation of compute and storage



All the compute engines store their data automatically in OneLake

The data is stored in a single common format

Delta – Parquet, an open standards format, is the storage format for all tabular data in Fabric

Once data is stored in the lake, it is directly accessible by all the engines without needing any import/export

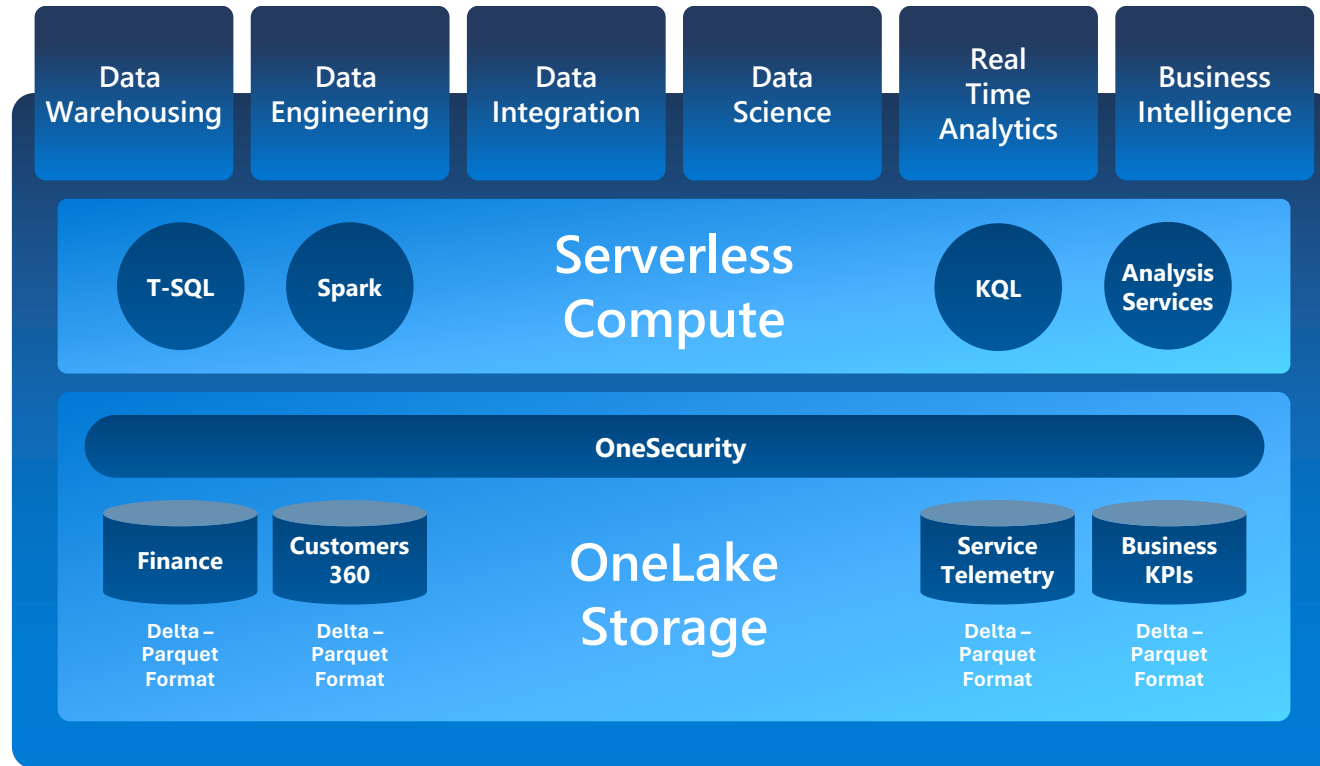
All the compute engines have been fully optimized to work with Delta Parquet as their native format

Shared universal security model is enforced across all the engines



One Copy for all computes

Universal security makes it real



All the compute engines store their data automatically in OneLake

The data is stored in a single common format

Delta - Parquet, an open standards format, is the storage format for all tabular data in Fabric

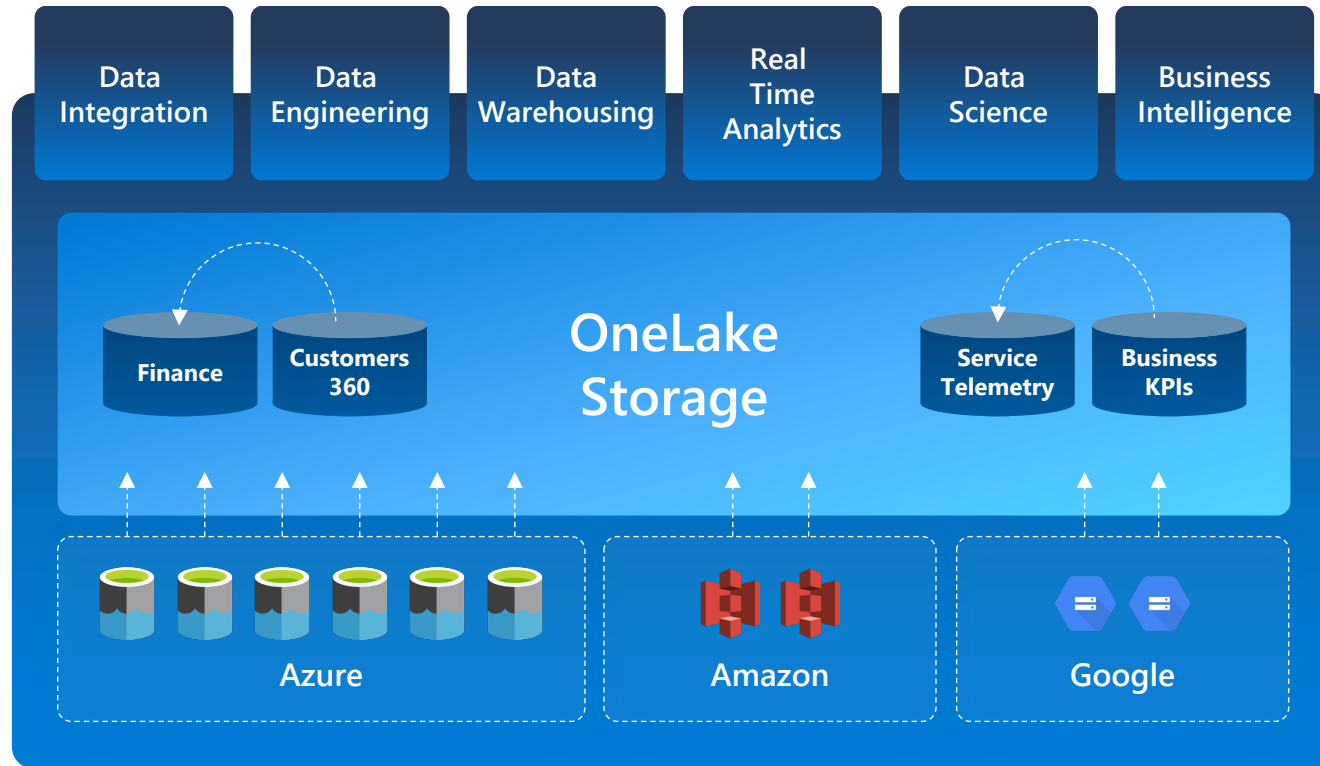
Once data is stored in the lake, it is directly accessible by all the engines without needing any import/export

All the compute engines have been fully optimized to work with Delta Parquet as their native format

Shared universal security model is enforced across all the engines

Taking One Copy to the next level

Shortcuts



Sharing data in OneLake is as easy as sharing files in OneDrive, removing the needs for data duplication

With **shortcuts**, data throughout OneLake can be composed together without any data movement

Shortcuts also allow instant linking of data already existing in Azure and in other clouds, without any data duplication and movement, making **OneLake the first multi-cloud data lake**

With support for industry standard APIs, OneLake data can be directly accessed by any application or service

Introducing Microsoft Fabric for Data Engineering

- Complete analytics platform
- Lake-centric and open architecture
- **Empower every office user**

Office Integration

The screenshot displays the Microsoft Power BI web interface. At the top, there is a search bar and a user profile for 'EPAM'. The left sidebar contains navigation options: Activity, Chat, Teams, Calendar, Calls, Power BI (selected), Data Activ..., and Apps. The main content area is titled 'Power BI Home' and includes a 'New report' button. Below this, a 'Recommended' section features three cards: 'People Science & Analytics' (opened by Volodymyr Tomurka), 'Teams activity analytics' (popular in the org), and another 'People Science & Analytics' card (frequently opened). A 'Recent' tab is active, showing a table of recent items with columns for Name, Type, Opened, Location, Endorsement, and Sensitivity.

Name	Type	Opened	Location	Endorsement	Sensitivity
Learning_DWH	Semantic m...	5 minutes ago	My workspace	—	EPAM Confi... ⓘ
My workspace	Workspace	a day ago	Workspaces	—	—
People Science & Analytics	App	a month ago	Apps	—	—
Item Sales Report	Report	2 months ago	My workspace	—	EPAM Confi... ⓘ
edd_sentry_metadata	Workbook	2 years ago	My workspace	—	—
SENTRY_ROLE	Report	2 years ago	My workspace	—	—

Introducing Microsoft Fabric for Data Engineering

- Complete analytics platform
- Lake-centric and open architecture
- Empower every office user
- **Persistent security and governance**

Regulatory compliance

Data residency

- Fabric will be available in every Azure region
- Data at rest: compliant with EUDB and other single-geo data residency regulations
- Multi-geo capacities allow control over content storage location in most Azure data centers world-wide

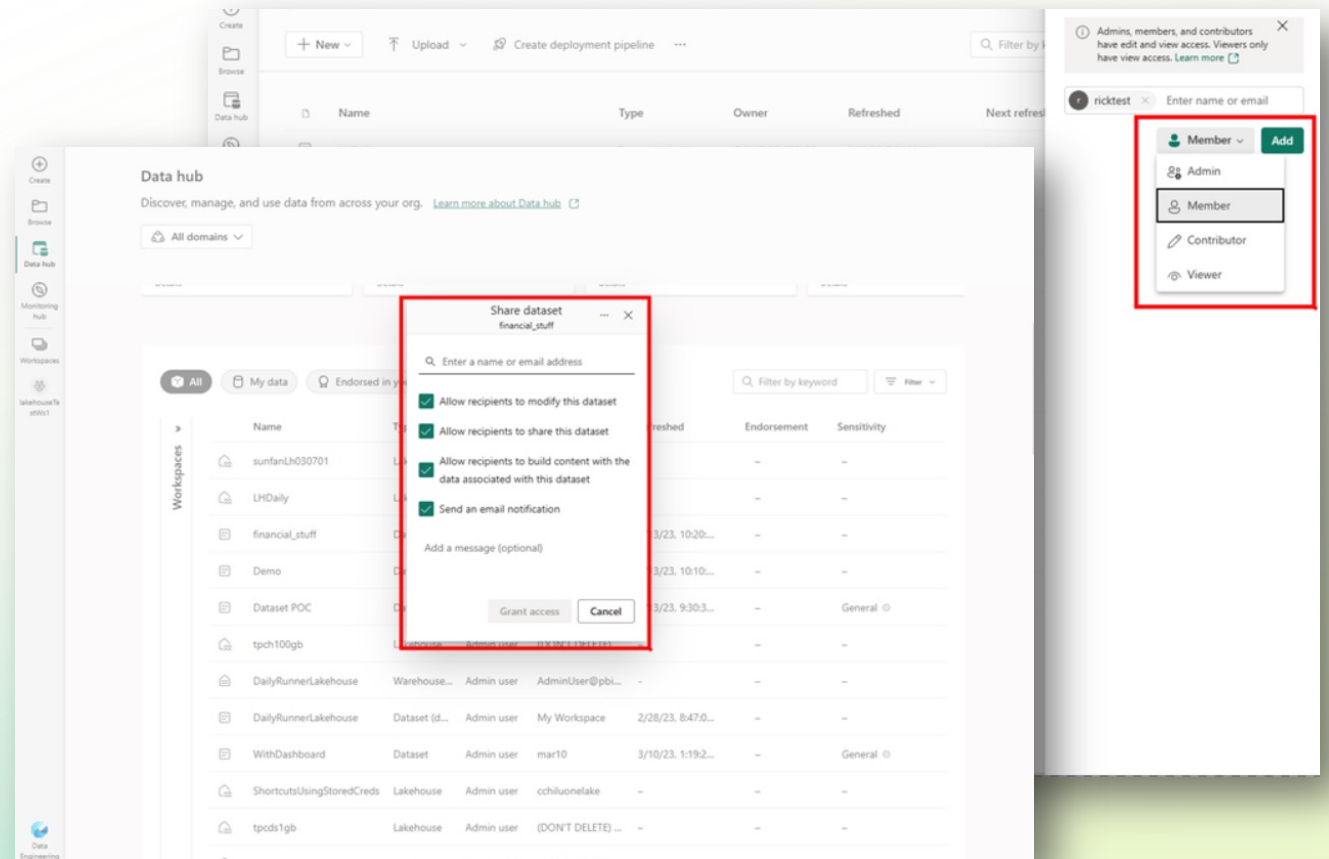


[Microsoft pledges support for EU Data Boundary](#)

Access control

Workspace roles and artifact permissions

- Fabric workspace roles define default permissions on workload items on the Control Plane
- Workload item permissions can be modified and managed via sharing
- On the Data Plane, Universal Security defines access policies on the delta tables directly and all workload compute engines will respect such policies



Defining user access via workspace roles and sharing

Microsoft Fabric

End-to-end analytics data fabric
From the data lake to the business user

Complete Analytics Platform

Everything, unified

SaaS Solution

Low Code Plus Pro Dev

Lake-centric and Open

OneLake

One Copy

Always Synced

Empower Every Office User

Familiar and Intuitive

Built Into Office 365

Insight to Action

Persistent Security and Governance

End-to-End Visibility

Always Governed

Secure by Default

Introducing Microsoft Fabric for Data Engineering

- Complete analytics platform
- Lake-centric and open architecture
- Empower every office user
- Persistent security and governance

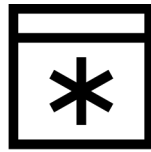


Deep Dive into Data Engineering

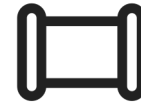
- **Data pipelines and movement**

The choice of tools for your data transformations

Notebooks



Spark Job
Definitions



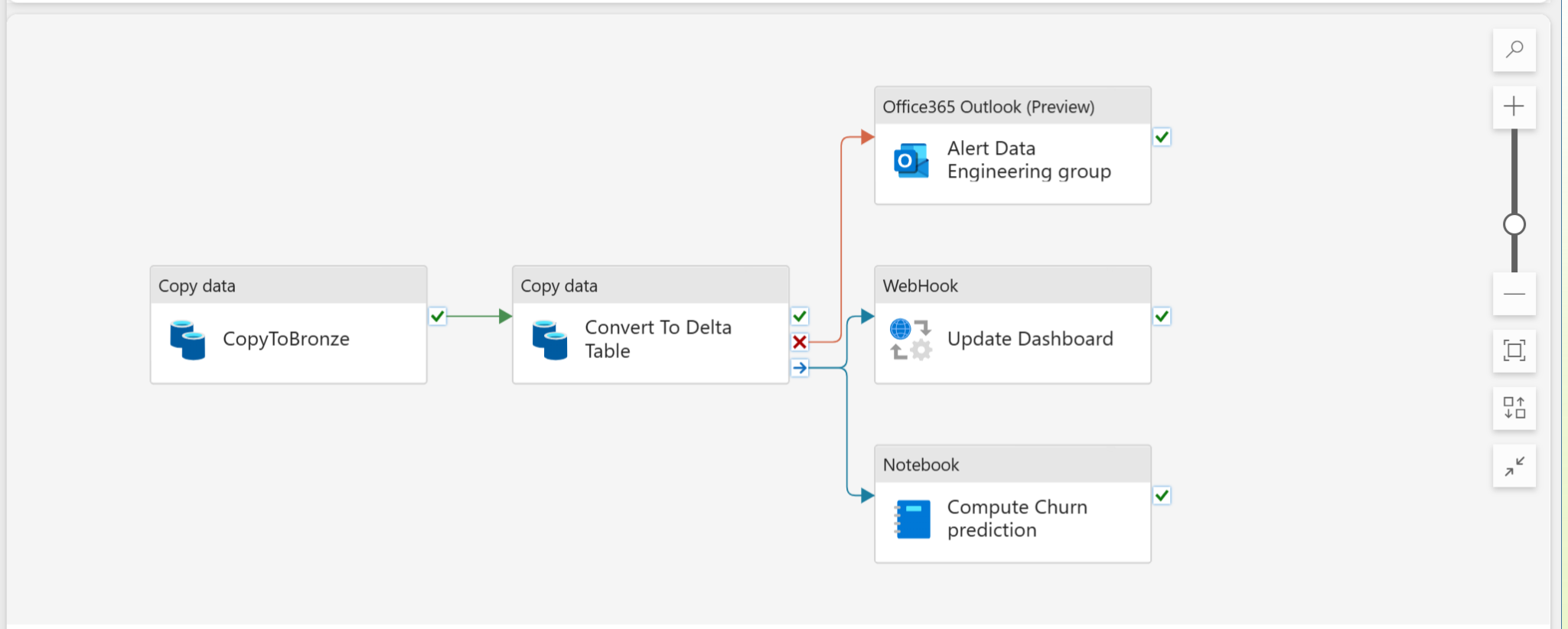
Data pipelines



Data Flows

Code-based pipelines

UI-based pipelines



- All sources
- Lakehouses
 - Lakehouse
- wwilakehouse
 - Tables
 - aggregate_sale_by_date_city
 - aggregate_sale_by_date_emp...
 - dimension_city
 - dimension_customer
 - dimension_date
 - dimension_employee
 - dimension_stock_item
 - fact_sale
 - Files

```

1 df = spark.read.parquet("Files/wwi-raw-data/full/fact_sale_1y_full/part-00000-ced648ca-e8c8-46e5-8526-5ca85d56e67e-c000.snappy.parquet")
2 df.printSchema()

```

PySpark (Python)

TRIDENT: SQLTHREADS DEV

- Notebook
 - Notebook 1
 - 00 - Initial Setup
 - Notebook 2
 - Compute Visitor fingerprint
 - Notebook 3
- Spark job definition
 - TestSparkJob
 - TestSparkJob2
- Lakehouse
 - SQLThreadsLakehouse
 - Tables
 - Files
 - DemoTable.csv File
 - wwilakehouse

C: > src > trident-root > 138a0420-3e20-471c-a436-5df0a515d57e > 00 - Initial Setup > 00 - Initial Setup.ipynb > Sample data copy from public storage to your lakehouse

Code | Markdown | Run All | Clear All Outputs | Outline | .NET Interactive

```
# Copyright (c) Microsoft Corporation.
# Licensed under the MIT License.

import os

bronze_zone = "//lakehouse/default/Files/bronze/"
# silver_zone = "//lakehouse/default/Tables/silver/"
# gold_zone = "//lakehouse/default/Tables/gold/"

zones = [bronze_zone] # [bronze_zone, silver_zone, gold_zone]
for path in zones:
    if not os.path.exists(path):
        print(path, "zone doesn't exist, creating it now!")
        os.mkdir(path)
    else:
        print(path, "zone exists already, skipping it now!")
```


[] Python

```
... StatementMeta(, 78fe1cca-f90e-484a-83d4-34beb734248d, 4, Finished, Available)

//lakehouse/default/Files/bronze/ zone doesn't exist, creating it now!
```

Sample data copy from public storage to your lakehouse

```
source_wasb = "wasbs://sampledata@azuresynapsestorage.blob.core.windows.net/WideWorldImportersDW"
```

 **ImportBlobToLakehouse**
Data pipeline

Search

About

Sensitivity label

Endorsement

Schedule

Last success is in

April 13, 2023 at 4:34:14 PM
(UTC+01:00) Brussels, Copenhagen, Madrid, Paris

Next refresh in

15 hour(s) 47 minute(s)

 Run

 Schedule

Scheduled run

On Off

Repeat


Daily

Time

08:00 AM  

[+ Add a time](#)

Start

mm/dd/yyyy --:-- -- 

End

mm/c

Time zone

(UTC+01:00) Brussels, Copenhagen, Madrid, ...

 **Compute Churn prediction**
Notebook



Search

About

Sensitivity label

Endorsement

Schedule

Last success is in

April 13, 2023 at 2:38:08 PM
(UTC) Coordinated Universal Time

The scheduled refresh is turned off

 Schedule

Scheduled run

On Off


Repeat

Hourly


Every

5 hour(s)

Start

mm/dd/yyyy --:-- -- 

End

mm/dd/yyyy --:-- -- 

Time zone

(UTC+01:00) Brussels, Copenhagen, Madrid, ...

Apply

Discard

Introducing Microsoft Fabric for Data Engineering

- Complete analytics platform
- Lake-centric and open architecture
- Empower every office user
- Persistent security and governance



Deep Dive into Data Engineering

- Data pipelines and movement
- **Data storage and architecture**



One Lake

Unified lake house

Based on open standards

Accessible from any workload

Lake view Table view

Search

Table

- > Deltatable_1
- > Deltatable_2
- > Account.parquet
- > Customer.csv

Files

- > TestFolder1
- > RandomFiles
- > CustomReview.csv
- > Username.parquet
- > Deltafolder

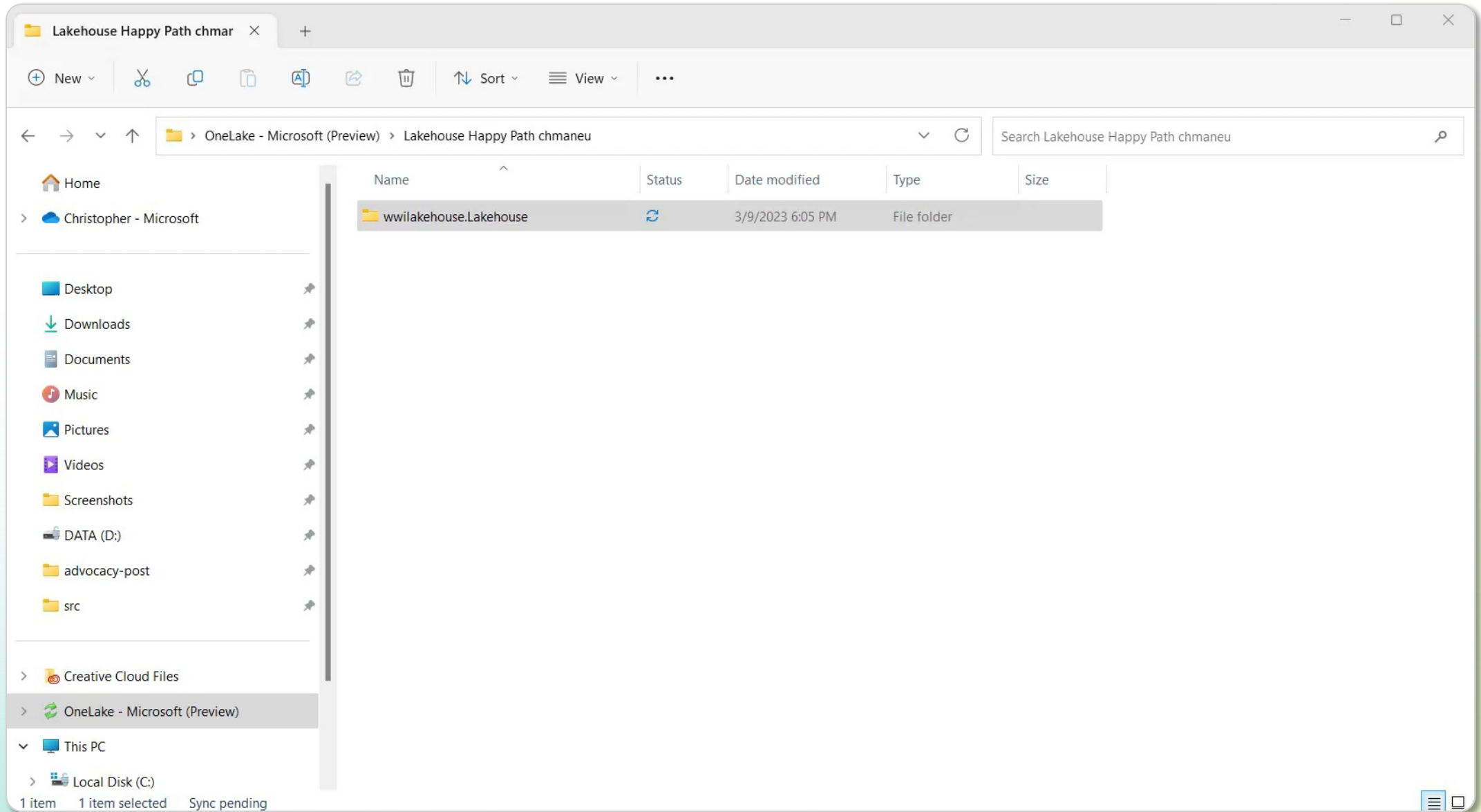
Lake view Table view

Search

Tables

- Deltatable_1
- Deltatable_2
- Account
- Customer





Data
warehousing

Data
engineering

Data
integration

Data
science

Real-time
analytics


Business
intelligence


OneLake

SaaS
foundation


All Azure Database File Generic protocol Services and apps


Search


 Amazon RDS for SQL Server
Database


 Amazon Redshift
Database


 Amazon S3
File


 Amazon S3 Compatible
File


 Apache Impala
Database


 Azure Blob Storage
Azure


 Azure Cosmos DB for NoSQL
Azure


 Azure Data Explorer (Kusto)
Azure


 Azure Data Lake Storage Gen1
Azure

 Azure Data Lake Storage Gen2
Azure


 Azure Database for PostgreSQL
Azure


 Azure SQL Database
Azure


 Azure SQL Database Managed Instance
Azure


 Azure Synapse Analytics
Azure

 Azure Table Storage
Azure

 Dataverse
Services and apps


 Dynamics CRM
Services and apps

 Google Cloud Storage
File


 HTTP
Generic protocol


 Hive
Database


 Microsoft 365
Services and apps


 OData
Generic protocol

 PostgreSQL
Database

 REST
Generic protocol, Services and apps

 SQL server
Database

 SharePoint Online List
Services and apps

 Snowflake
Services and apps

 Spark
Database

Choose data source

Select a connector or directly drag a file from your computer.

All categories

File

























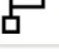















Database

Power Platform

Azure

Online services

Other

 Excel workbook File	 Text/CSV File	 XML File	 JSON File
 Folder File	 PDF File	 Parquet File	 SharePoint folder File
 SQL Server database Database	 Access Database	 SQL Server Analysis Services Database	 Oracle database Database
 IBM Db2 database Database	 MySQL database Database	 PostgreSQL database Database	 Teradata database Database
 SAP HANA database Database	 SAP BW Application Server Database	 SAP BW Message Server Database	 Snowflake Database
 Google BigQuery Database	 Amazon Redshift Database	 Impala Database	 Vertica Database
 Dataflows Power Platform	 Power BI dataflows (Legacy) Power Platform	 Dataverse Power Platform	 Common Data Service (Legacy) Power Platform
 Azure SQL database Azure	 Azure Synapse Analytics (SQL DW) Azure	 Azure Analysis Services Azure	 Azure Blobs Azure
 Azure Tables Azure	 Azure Data Explorer (Kusto) Azure	 Azure Data Lake Storage Gen2 Azure	 Azure HDInsight Spark Azure
 SharePoint Online list	 Microsoft Exchange Online	 Salesforce objects	 Salesforce reports

Introducing Microsoft Fabric for Data Engineering

- Complete analytics platform
- Lake-centric and open architecture
- Empower every office user
- Persistent security and governance



Deep Dive into Data Engineering

- Data pipelines and movement
- Data storage and architecture
- **Delivering data to data analysts and data scientists**

Notebook 1 | EPAM Confidential\Confidential · Saved ▾

Search

Fabric Trial: 52 days left

Home Edit Run View

Editing ▾ Comments Share

Standard session ▾ PySpark (Python) ▾ Workspace default ▾ Data Wrangler ▾ Copilot

← All sources <<

+ Code + Markdown

Lakehouses

+ Lakehouse

wwilakehouse ↗ ↻

Tables

- > aggregate_sale_by_date_city
- > aggregate_sale_by_date_emp...
- > dimension_city
- > dimension_customer
- > dimension_date
- > dimension_employee
- > dimension_stock_item
- > fact_sale

Files

- > wwi-raw-data

Session ready AutoSave: On

Selected Cell 0 of 0 cells

NubesGen - M... | Search

File | Refresh | Share | Create a report | Analyze in Excel | Lineage | Open data model

Home | Create | Browse | Data hub | Workspaces | My workspace | NubesGen - Monthly... | More... | Power BI

Details for NubesGen - Monthly Report

+ Add description

Location	Refreshed	Sensitivity
My Workspace	4/20/23, 6:35:59 AM	General

Visualize this data

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

+ Create from scratch

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 NubesGen - Monthly Report

Filter by keyword | Filter



Home



Create



Browse



OneLake data hub



Monitoring hub



Workspaces



Fabric DEMO



wwilakehouse



02 - Data Transform...



01 - Create Delta Tables



IngestDataFromSource...

Home

SQL analytics endpoint



New SQL query

New visual query

Explorer

Warehouses

wwilakehouse

Schemas

dbo

Tables

aggregate_sale_by_dat...

aggregate_sale_by_dat...

dimension_city

dimension_customer

dimension_date

dimension_employee

EmployeeKey

WWIEmployeeID

Employee

PreferredName

IsSalesperson

ValidFrom

ValidTo

LineageKey

SQL query 1

SQL query 2



Save as view

```

1 SELECT TOP (100) [EmployeeKey]
2     , [WWIEmployeeID]
3     , [Employee]
4     , [PreferredName]
5     , [IsSalesperson]
6     , [ValidFrom]
7     , [ValidTo]
8     , [LineageKey]
9 FROM [wwilakehouse].[dbo].[dimension_employee]

```

Get the full current results in an Excel worksheet.

Messages

Results

Open in Excel

Explore this data (preview)

Search

	EmployeeKey	WWIEmployeeID	Employee	PreferredName	IsSalesperson	ValidFrom
1	0	0	Unknown	N/A	0	2013-01-01 00:00:00.000000
2	175	20	Jack Potter	Jack	1	2016-05-31 23:13:00.000000
3	176	19	Jai Shand	Jai	0	2016-05-31 23:13:00.000000
4	177	18	Katie Darwin	Katie	0	2016-05-31 23:13:00.000000
5	178	17	Piper Koch	Piper	0	2016-05-31 23:13:00.000000
6	179	16	Archer Lamble	Archer	1	2016-05-31 23:13:00.000000
7	180	15	Taj Shand	Taj	1	2016-05-31 23:13:00.000000
8	181	14	Lily Code	Lily	1	2016-05-31 23:13:00.000000
9	182	13	Hudson Hollinworth	Hudson	1	2016-05-31 23:13:00.000000
10	183	12	Henry Forlonge	Henry	0	2016-05-31 23:13:00.000000
11	184	11	Ethan Onslow	Ethan	0	2016-05-31 23:13:00.000000

DBeaver 24.0.3 - dimension_customer

Database Navigator | Projects

Enter a part of object name here

- > AclaraOne - jdbc:sqlserver://mdmdev.database.windows.net:1433;datab...
- > clickhouse.apps.dataplatform.westus2.aroapp.io - clickhouse.apps.dat...
- > clickhouse.apps.dataplatformqa2.westus2.aroapp.io - clickhouse.apps...
- > dev-cis-hub-db - jdbc:sqlserver://dev-cis-hub-db.database.windows.ne...
- > InsightsMetaData - mdmdev.database.windows.net:1433
- > master - 2bzbxne7jytezctj7fe7gz6jdu-tou4qy6fjnruvm6docsqk3x3pi.dat...
- > Databases
 - > DataflowsStagingLakehouse
 - > DataflowsStagingWarehouse
 - > wwilakehouse
 - > Schemas
 - > _rsc
 - > db_accessadmin
 - > db_backupoperator
 - > db_datareader
 - > db_datawriter
 - > db_ddladmin
 - > db_denydatareader
 - > db_denydatawriter
 - > db_owner
 - > db_securityadmin
 - > dbo
 - > Tables
 - dimension_customer 72K
 - External Tables
 - Views
 - Indexes
 - Procedures
 - Synonyms
 - Data Types
 - guest
 - queryinsights

dimension_customer X

Properties | Data | ER Diagram

dimension_customer | Enter a SQL expression to filter results (use Ctrl+Space)

	123 CustomerKey	123 WWICustomerID	ABC Customer	ABC BillToCustomer	ABC Category	ABC BuyingGroup	ABC
3	103	103	Tailspin Toys (Kalvesta, KS)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Nas
4	104	104	Tailspin Toys (Wallagrass, ME)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Lab
5	105	105	Tailspin Toys (Tomnolen, MS)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Sur
6	106	106	Tailspin Toys (Tumacacori, AZ)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Shi
7	107	107	Tailspin Toys (Glen Avon, CA)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Kar
8	108	108	Tailspin Toys (Bernie, MO)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Bha
9	109	109	Tailspin Toys (South Laguna, CA)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Ae-
10	110	110	Tailspin Toys (North Crows Nest, IN)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Din
11	111	111	Tailspin Toys (Oriole Beach, FL)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Ada
12	112	112	Tailspin Toys (Sallyards, KS)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Ing
13	113	113	Tailspin Toys (Dahlia, NM)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Jae
14	114	114	Tailspin Toys (Cherry Grove Beach, SC)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Duc
15	115	115	Tailspin Toys (Bethania, NC)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Jag
16	116	116	Tailspin Toys (Rafael Capó, PR)	Tailspin Toys (Head Office)	Novelty Shop	Kids Toys	Tom
17	79	79	Tailspin Toys (Page City, KS)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Mar
18	80	80	Tailspin Toys (Valdese, NC)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Bha
19	81	81	Tailspin Toys (Big Moose, NY)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Ver
20	82	82	Tailspin Toys (La Cueva, NM)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Rak
21	83	83	Tailspin Toys (Absecon, NJ)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Sar
22	84	84	Tailspin Toys (Aceitunas, PR)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Eek
23	85	85	Tailspin Toys (Andrix, CO)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Tom
24	86	86	Tailspin Toys (New Lexington, OH)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	And
25	87	87	Tailspin Toys (Sauquoit, NY)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Isak
26	88	88	Tailspin Toys (Dracut, MA)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Am
27	89	89	Tailspin Toys (Victory Gardens, NJ)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Aak
28	90	90	Tailspin Toys (Tolna, ND)	Tailspin Toys (Head Office)	Novelty Shop	Toddler Toys	Sta

app.fabric.microsoft.com

wwilakehouse | EPAM Confidential\Confidential | Search

Fabric Trial: 52 days left

Home | Lakehouse | Share

Get data | New semantic model | Open notebook | Manage OneLake data access (preview)



Explorer

- wwilakehouse
 - Tables
 - aggregate_sale_by_date...
 - aggregate_sale_by_date...
 - dimension_city
 - dimension_customer
 - dimension_date
 - dimension_employee
 - dimension_stock_item
 - fact_sale
 - Files
 - wwi-raw-data
 - full
 - incremental
 - tables

Get data in your lakehouse

- Upload files**
Upload data from your local machine.
- New Dataflow Gen2**
Prep, clean, transform, and ingest data.
- New data pipeline**
Ingest data at scale and schedule data workflows.
- Open notebook**
Transform and ingest data using code in Apache Spark.
- New shortcut**
Access data through an external lake.

Data Engineering

File ▾ + Create a report ▾  Analyze in Excel  Lineage ▾  Open data model ...

Details for NubesgenLH

+ Add description

 Location

NubesGen 

 Endorsement

Promoted

 Refreshed

4/12/23, 2:41:42 PM

 Sensitivity

Confidential\Microsoft Exter

Introducing Microsoft Fabric for Data Engineering

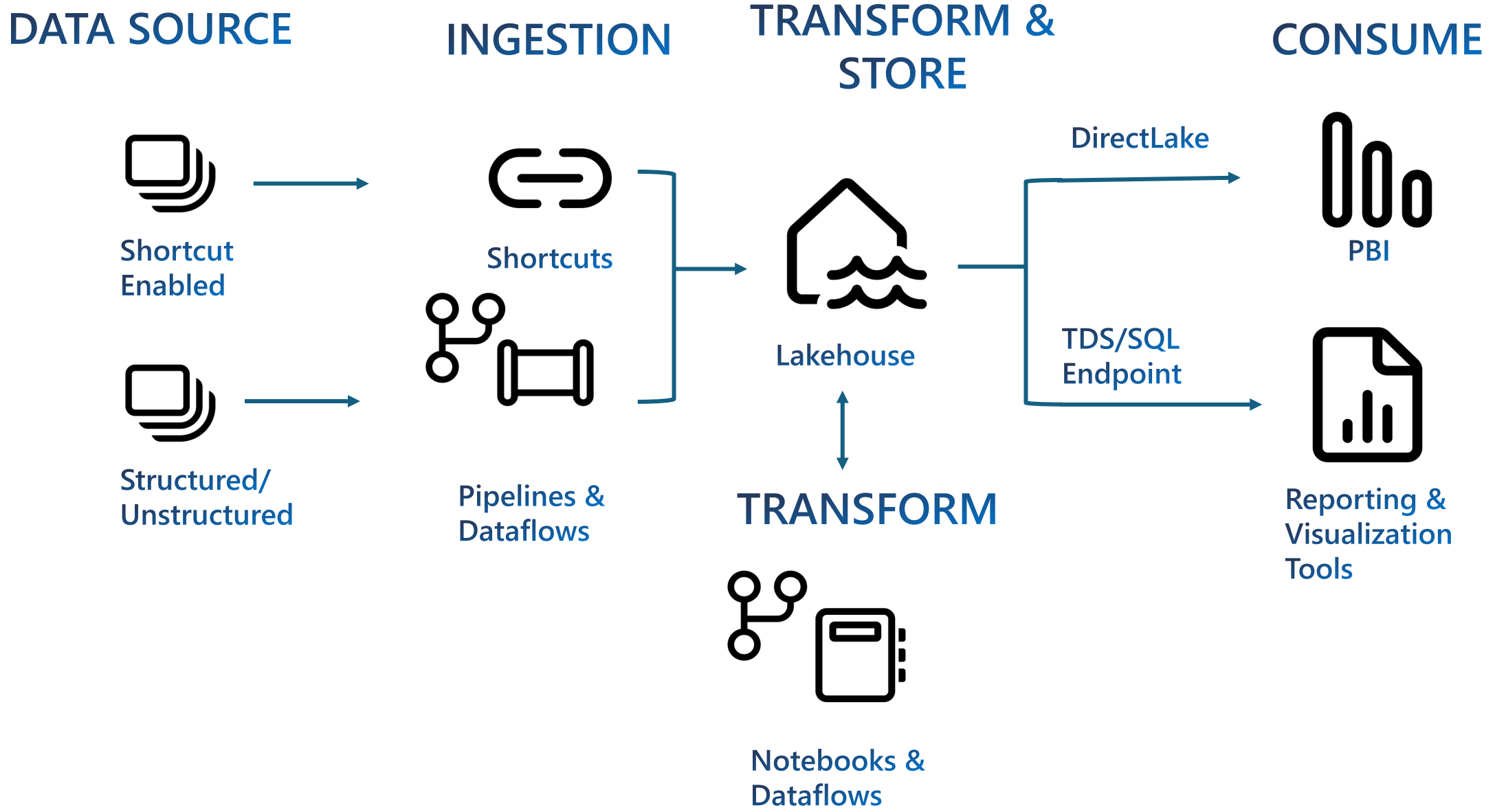
- Complete analytics platform
- Lake-centric and open architecture
- Empower every office user
- Persistent security and governance

Deep Dive into Data Engineering

- Data pipelines and movement
- Data storage and architecture
- Delivering data to data analysts and data scientists

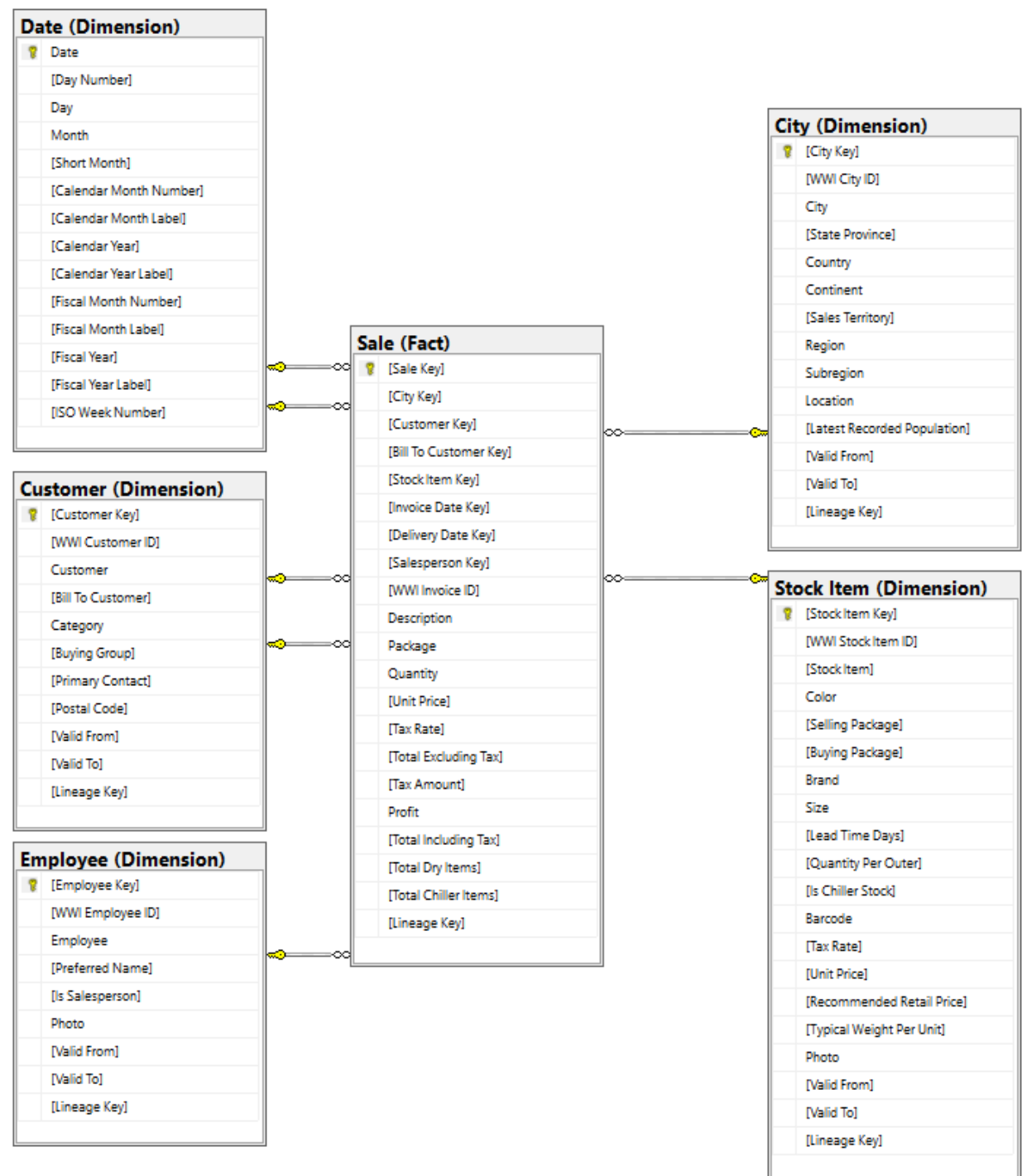
DEMO

Architecture



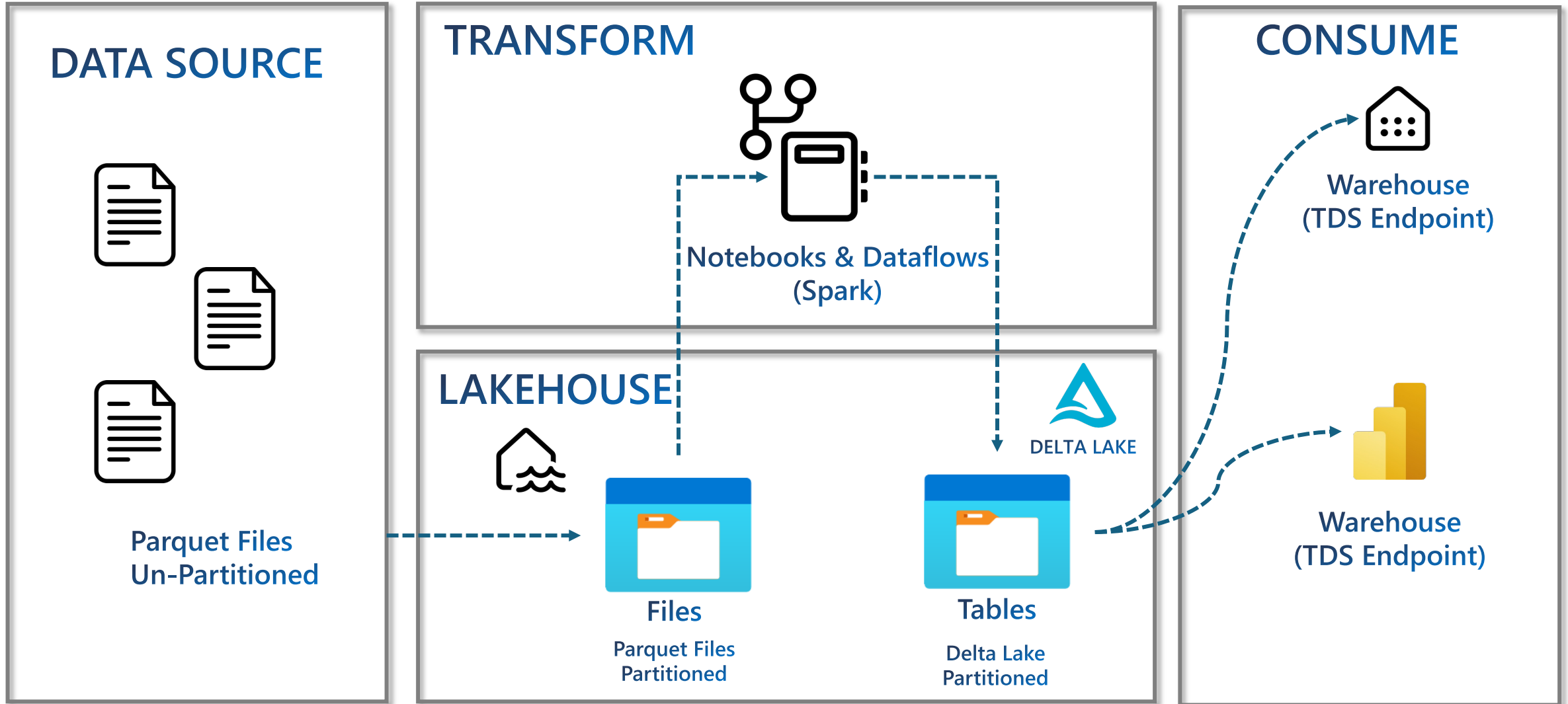
Data model

Wide World Importers (WWI) data model



See [Wide World Importers sample databases for Microsoft SQL](#).

Data and Transformation flow



References

Learning Fabric

- [Introduction](#) End-to-End Analytics in Microsoft Fabric
- [Lakehouse](#) Get Started Lakehouses
- [Spark on Lakehouse](#) : Use Apache Spark in Lakehouse
- [Work with Delta](#) Delta Lake Tables in Microsoft Fabric
- [Data Factory Pipelines](#) Pipelines, Activities, Templates
- [Data Warehouses](#) Get started with Data Warehouses
- [Real-Time Analytics](#) Analyze real-time data
- [Data Science](#) Get started with data science in Microsoft Fabric
- [Administer](#) Administration, Security, and Govern data in Microsoft Fabric
- [Medallion Architecture](#) Design Fabric Medallion Architecture with Bronze, Silver and Gold layers of Lakehouse
- [DataFlow Gen2](#) Ingest with Dataflows in Microsoft Fabric
- [Data Analysis with Kusto Query Language](#) Explore the fundamentals of data analysis
- [Azure Data Engineer - free online training from Azure](#)

Q&A AND SOCIAL ADS

Slides



X @philg0ld

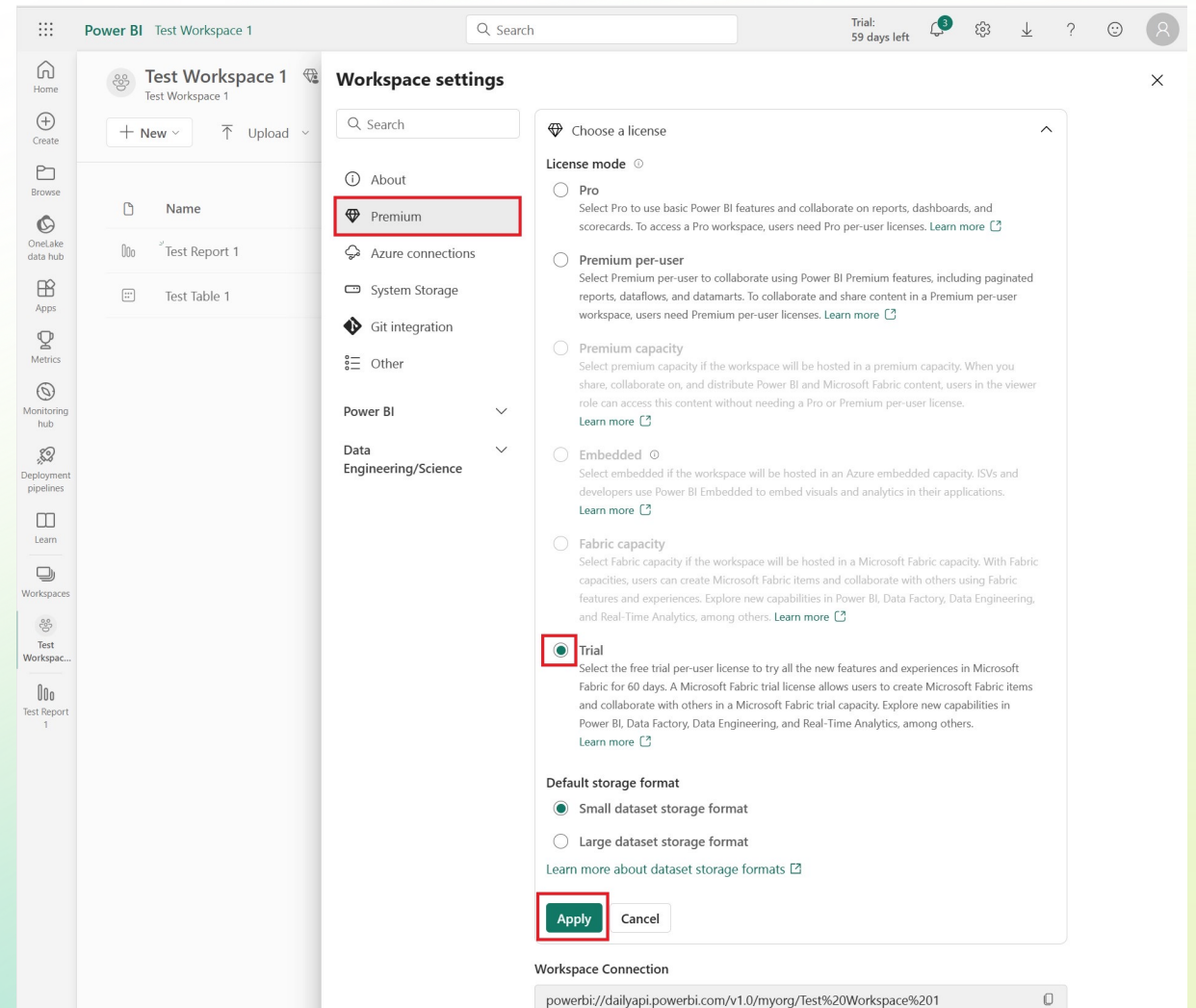
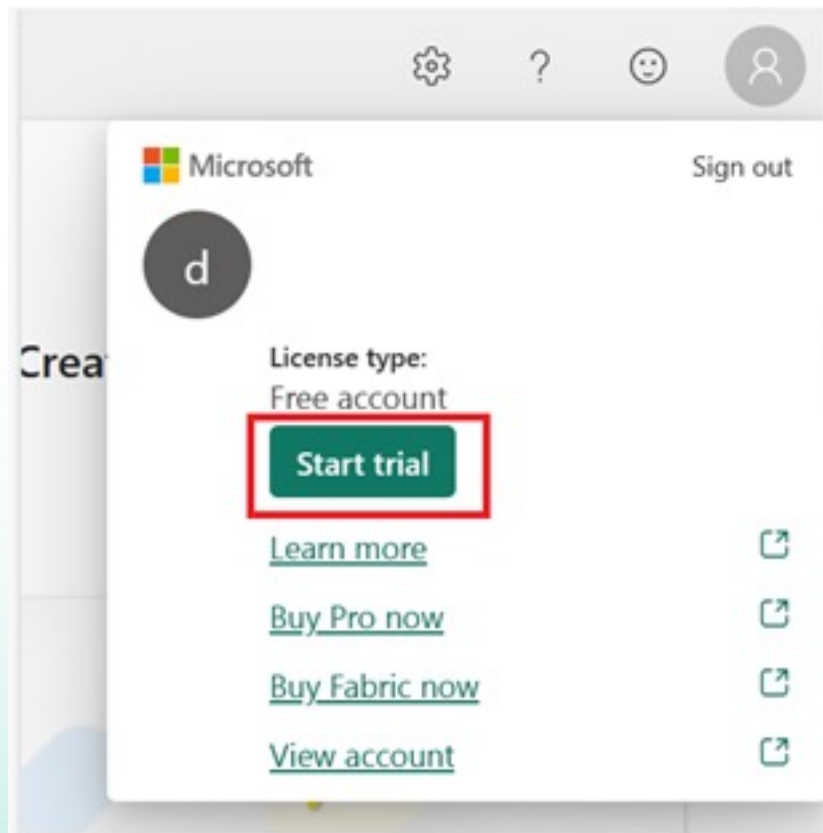
[#NashvilleDataEngineering](#)



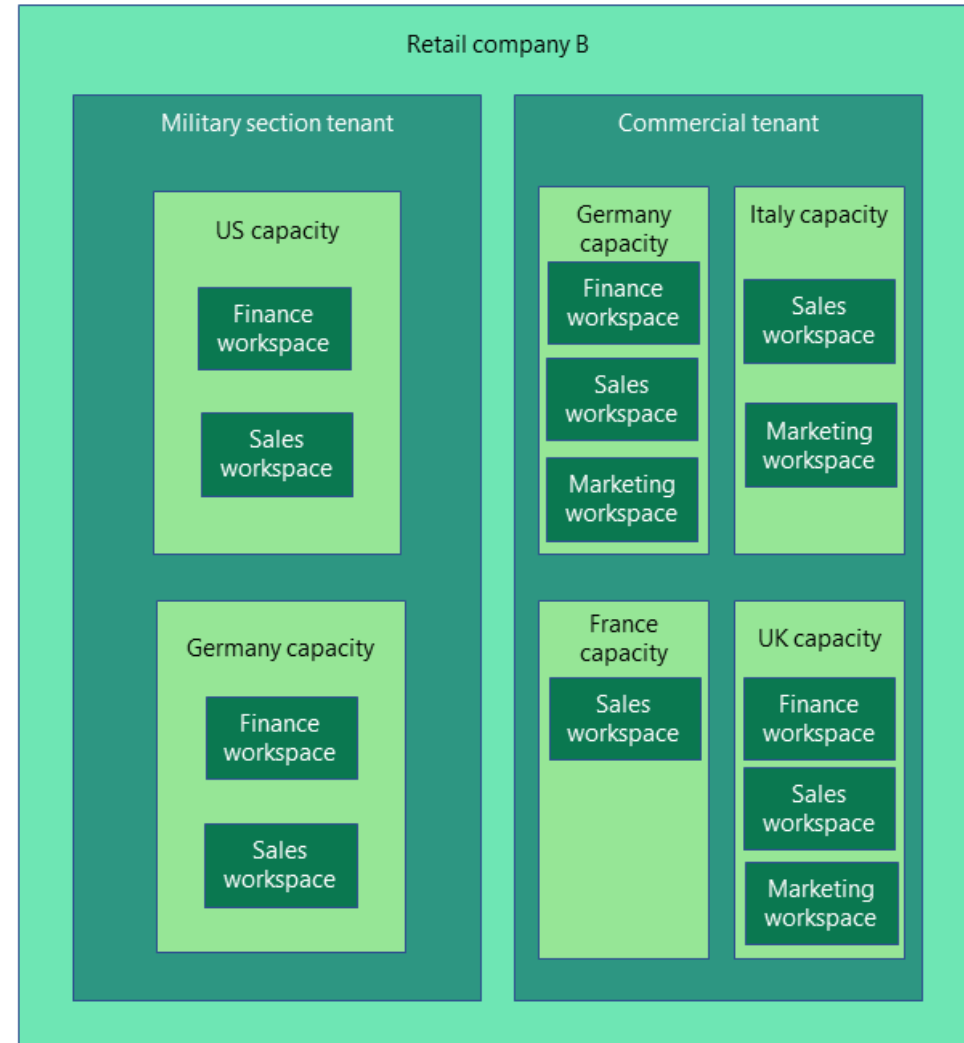
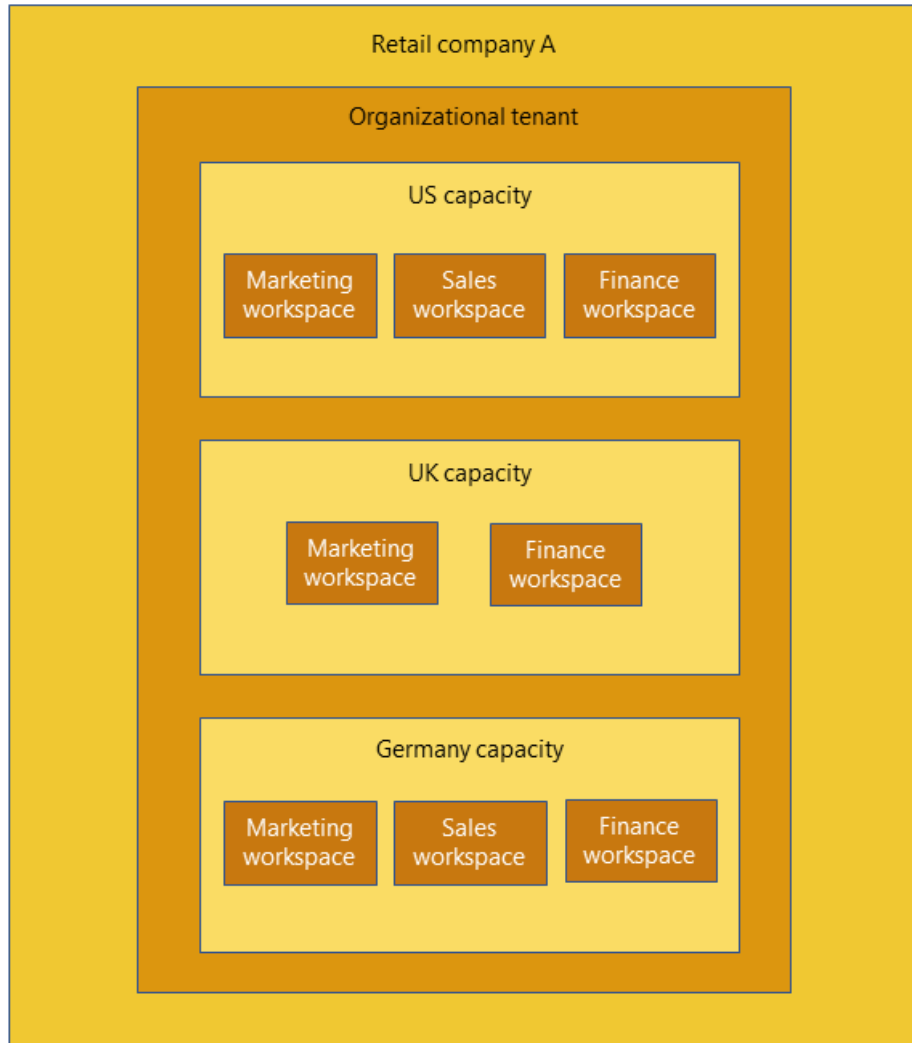
<http://linkedin.com/in/philg0ld>

APPENDIX

Sign up for a free Fabric trial

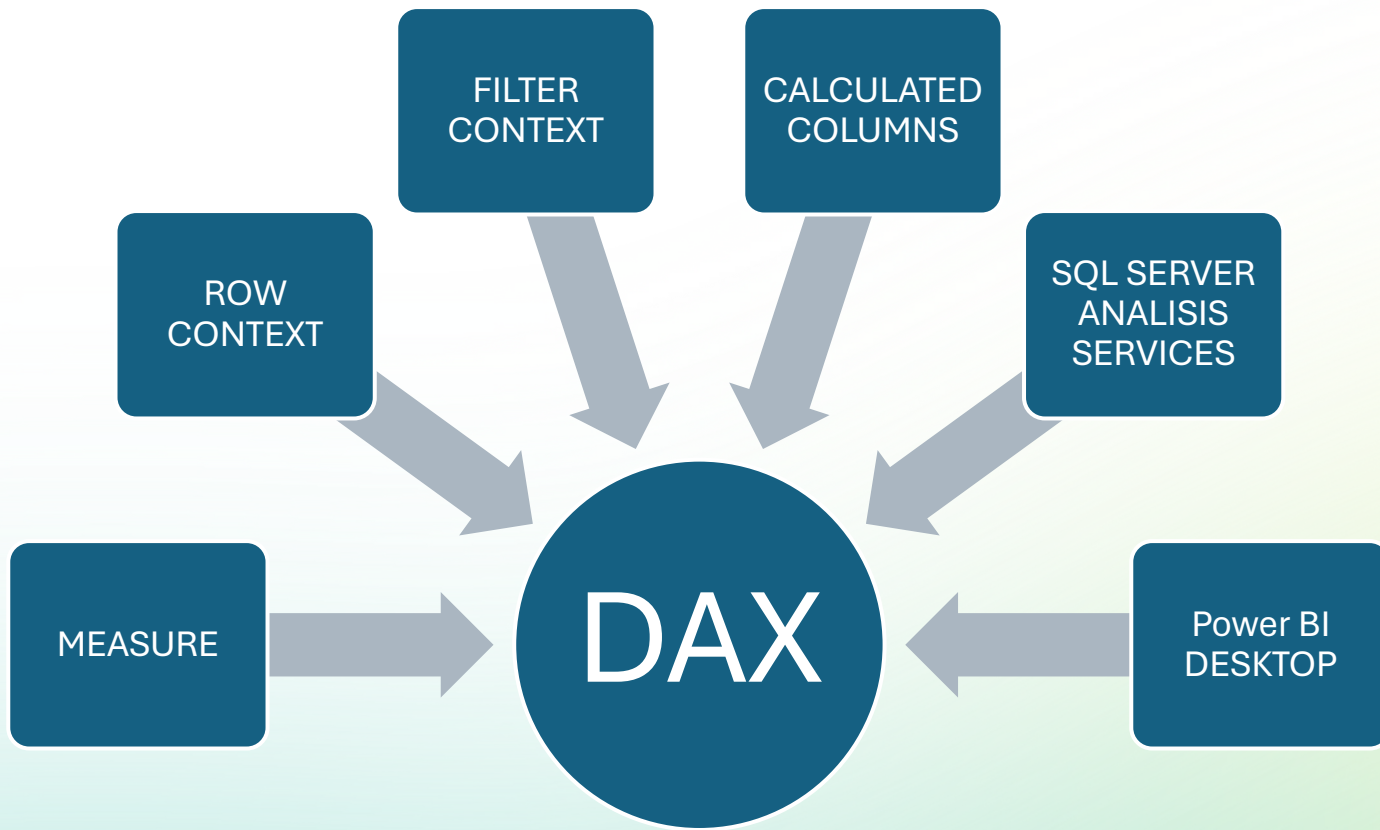


Microsoft Fabric concepts and licenses



DAX overview

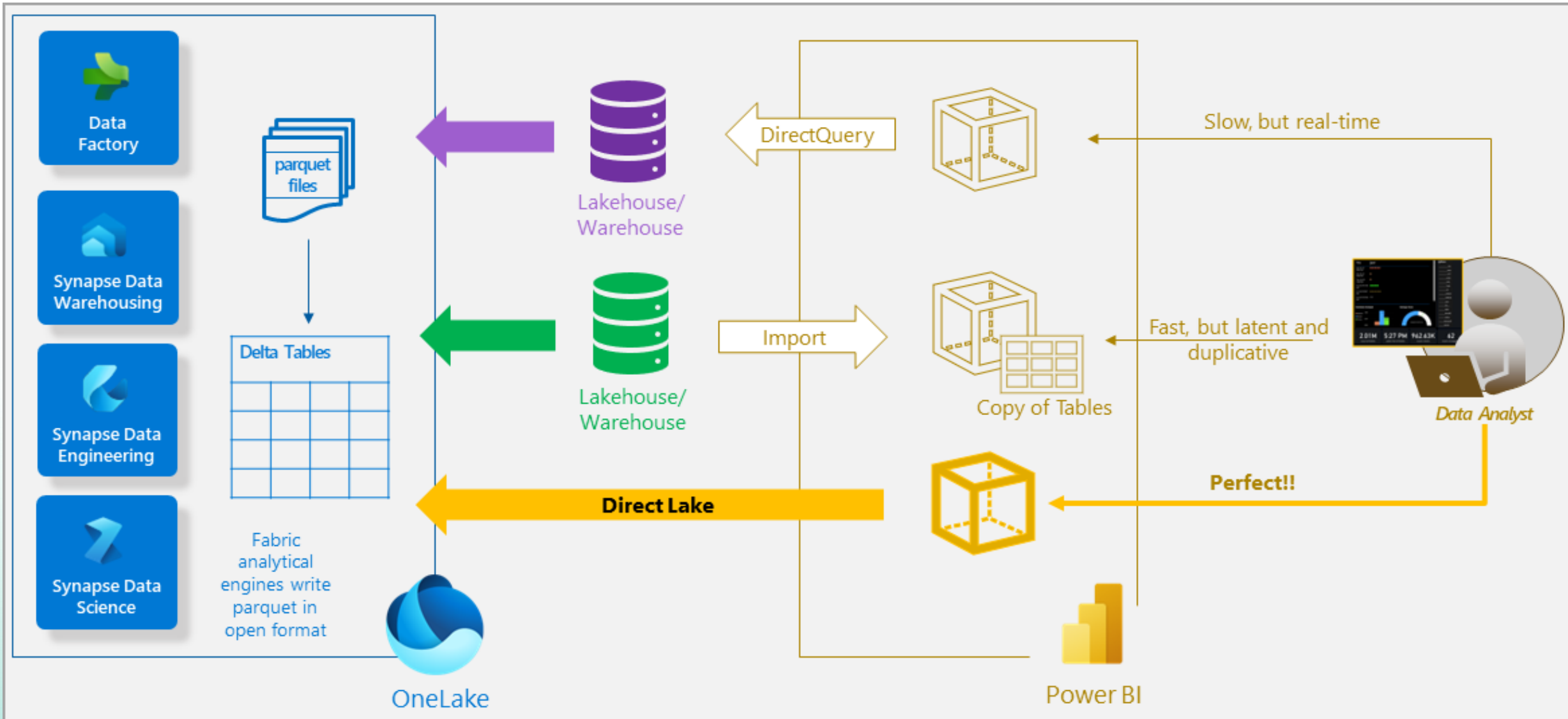
What is Data Analysis Expressions (DAX)?



1. DAX is a formula language like the formula language in Excel that uses functions
2. Create new information/ column/ measures/ calculation/ formula
3. DAX usage may depend on the design pattern:
 - Work could be pushed to the Power Query
 - Work could be pushed to the Database
 - Work could also be pushed to an ETL tool
4. Great for time intelligence, for example, YTD, MTD, parallel period comparisons
5. Enhances the data model

Power BI Architecture - Fabric

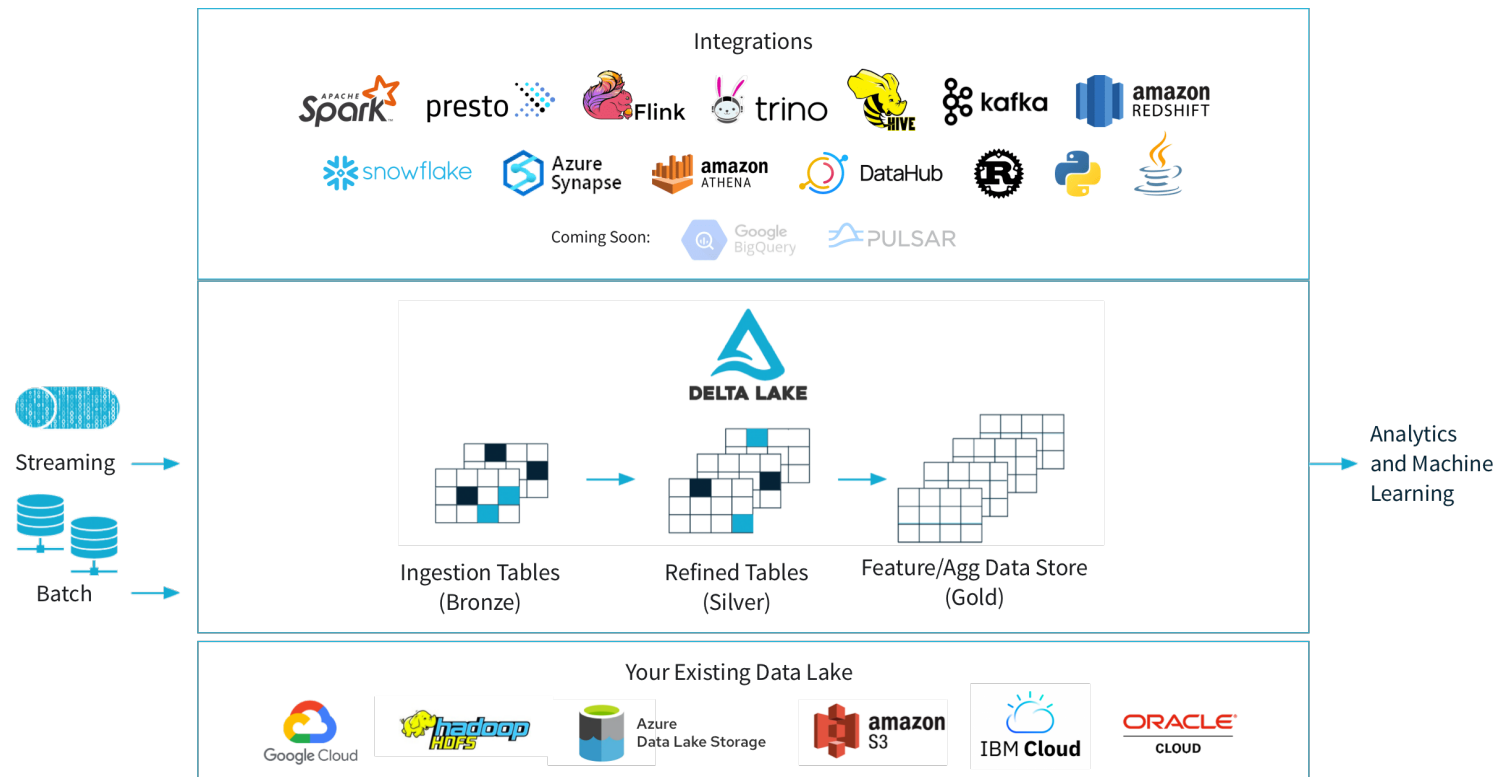
Understanding Direct Lake in Microsoft Fabric



What is Delta Lake?



- Open-source storage framework that enables **building a Lakehouse architecture** with compute engines, including Spark, PrestoDB, Flink, Trino, and Hive, and APIs for Scala, Java, Rust, Ruby, and Python.
- **ACID-compliant storage layer** that runs on top of cloud object stores such as MinIO, Hadoop HDFS, Amazon S3, Azure Data Lake Storage, and Google Cloud Storage.
- Provides features such as **scalable metadata handling for** petabyte-scale tables with billions of partitions and files with ease.
- **Provides time travel access/reverts** to earlier versions of data for audits, rollbacks, or reproduce.
- **Production-ready** and has been battle-tested in over 10,000+ production environments.



Delta Lake ACID Implementation

