

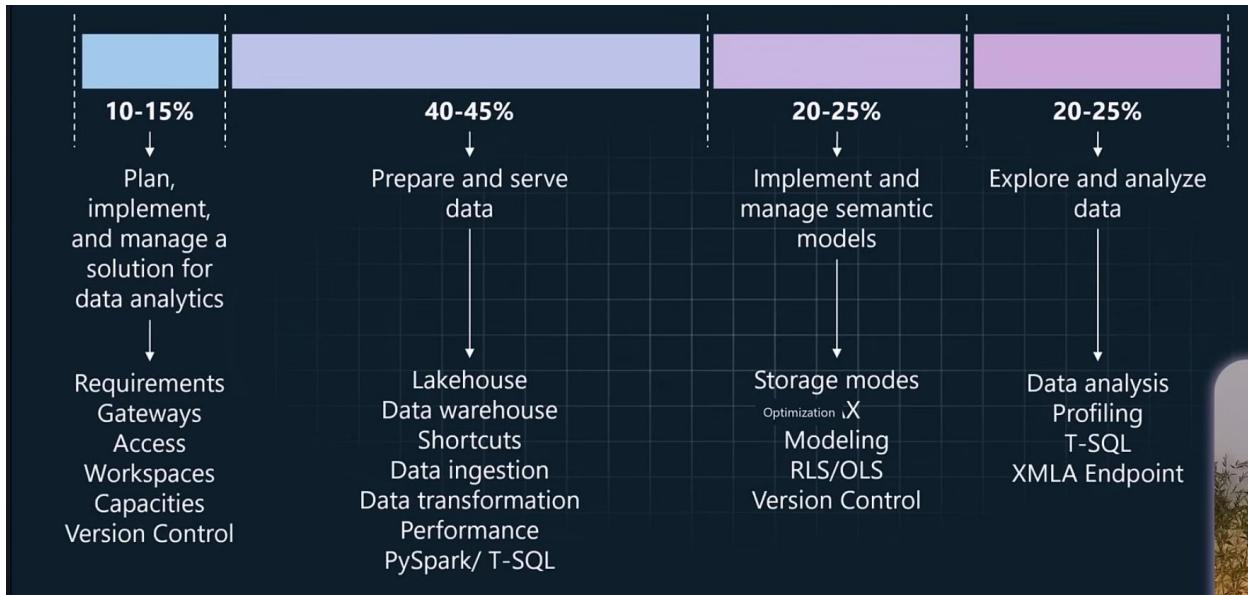
Fabric Analytics Concepts Notes

Fabric Analytics Concepts

Ch1. Intro	2
Ch2. Plan Data Analytics Env	2
Ch3. Implement & Manage a Data Analytics Env	11
Capacity settings in Azure Portal	12
Capacity settings in Fabric Portal.....	13
Ch4. Manage analytics dev lifecycle.....	22
Setup deployment pipeline in MS Fabric	24
Impact Analysis Tool of Fabric (Lineage View).....	28
Ch5. Getting Data to Fabric	31
File Sharing Protocols	36
Ch6. SQL, DWH, Scheduling	38
Azure DataFactory	43
Data Pipeline	43
Ch7. Transform Data	47
Transformations in PowerQuery	48
Transformations in SQL Endpoint	49
Transformations in PySpark.....	50
Data Modeling.....	54
SCD.....	55
JOIN Types	61
Ch8. Optimizing Performance.....	64
Ch9. Design & build semantic models	80
Normalized vs de-normalized data.....	80
KQL	82
Storage modes	83
DAX Variables, functions & parameters	86
Ch10. Secure & Optimize Semantic Models	97
DAX Studio	100
Tabular Editor	102
Ch11. Perform exploratory analytics.....	108
PBI Visuals	109
PBI Features.....	112
Data Profiling with Power Query.....	114
Ch12. Query data using SQL.....	118

Fabric Analytics Concepts Notes

Ch1. Intro



S

Ch2. Plan Data Analytics Env

Data ingestion requirements

The requirements we need:

- the Fabric items/ features you will need to get data into Fabric
- how these items/ features will need to be configured

Deciding factors:

Where is the external data stored?

- ADLS Gen2, Amazon S3 (or S3 backed), Google Cloud Storage or Dataverse?
- Azure SQL, Azure Cosmos DB, Snowflake?
- On-premise SQL?
- Real-time events?
- Other?

Some options:

- Shortcut
- Database mirroring
- ETL - Dataflow
- ETL – Data pipeline
- ETL – Notebook
- Eventstream

A video call interface is visible in the top right corner, showing a man speaking. A yellow arrow points from the 'Real-time events?' list item to the 'Eventstream' option in the list of options.

⊕ Data ingestion requirements

The requirements we need:

- the Fabric items/ features you will need to get data into Fabric
- how these items/ features will need to be configured



Deciding factors:

How is the data secured?

- On-premise SQL _____
- Azure virtual network / private endpoint _____

Some features:

On-premise data gateway

VNet data gateway

Fast copy

Staging

What is the volume of the data?

- Low (megabytes per day) _____
- Medium (gigabytes per day) _____
- High (many GB or terabytes per day) _____

⊕ Data gateways

Before we move on, let's explore data gateways in more detail

On-premise data gateway

1. Install the data gateway on the on-premise server (or update to Fabric version if already exists)
2. In Fabric, create a new On-Premise Data Gateway connection
3. Use the gateway in a Dataflow or Data Pipeline to get data into Fabric

Virtual network (VNet) data gateway

1. Set network configuration in Azure.
 1. Register Power Platform resource provider
 2. Create a private endpoint on your Azure object
 3. Create a subnet
2. In Fabric, create a new virtual network Data Gateway connection
3. Use the gateway in a Dataflow or Data Pipeline to get data into Fabric

Fabric Analytics Concepts Notes

Data storage requirements

The requirements we need:

- the Fabric data store(s)
- overall architectural pattern (medallion, lambda etc)

Deciding factors:

What is the data type?

- Structured, semi-structured and/or unstructured?
- Relational/ structured
- Real-time/ streaming

What skills exist in the team?

- T-SQL
- Spark (Python/ Spark SQL, Scala)
- KQL

Some options:

- Lakehouse
- Data warehouse
- KQL database

```
graph LR; A[Structured, semi-structured and/or unstructured?]; B[Relational/ structured]; C[Real-time/ streaming]; D[T-SQL]; E[Spark (Python/ Spark SQL, Scala)]; F[KQL]; A --> G[Lakehouse]; A --> H[Data warehouse]; A --> I[KQL database]; B --> G; B --> H; C --> I; D --> G; D --> H; E --> I; F --> G; F --> H; F --> I;
```

PBI JSON Theme

Fabric Analytics Concepts Notes

Restricted Mode is intended for safe code browsing. Trust this window to enable all features.

```
{ } my-report-theme.json X Welcome
C: > Users > learn > Downloads > { } my-report-theme.json > [ ] dataColors
1  [
2    "name": "Accessible Orchid",
3    "dataColors": [ ...
9      ],
36      ],
37      "foreground": "#192229",
38      "background": "#FFFFFF",
39      "foregroundNeutralSecondary": "#716E76",
40      "backgroundLight": "#EBE8FA",
41      "foregroundNeutralTertiary": "#96939E",
42      "backgroundNeutral": "#DAD8E8",
43      "tableAccent": "#BD3978",
44      "maximum": "#3F213F",
45      "center": "#663466",
46      "minimum": "#6E98B5",
47      "bad": "#BD3978",
48      "neutral": "#C480A7",
49      "good": "#9C6584"
50  ]
```

31:42 / 6:03:40 • Plan a data analytics environment >

Fabric Analytics Concepts Notes

Question

You are running an F2 capacity and you regularly experiencing throttling.

There are many long-running Spark jobs take on average 3 hours to complete. You need these to complete in under one hour, so you plan to increase the SKU of the capacity.

Where would go to go to make this change?

A) Go to Workspace Settings > Spark settings

B) Go to Admin Portal > Capacity Settings and click through to Azure to update your capacity

C) Go to the monitoring hub, and look at the run history.

D) Use the Capacity Metrics app

Question

Your data governance team would like to certify a Semantic Model to make it discoverable in the organisation.

Only the data governance team should be able to do this.

In what order should you complete the following tasks to certify a Semantic Model?

- 1) Create a security group for the data governance team
- 2) Enable the 'Make certified content discoverable' in Admin Portal > Tenant Settings > Discovery Settings
- 3) Make sure the 'Make certified content discoverable' setting applies only to the Data Governance security group.
- 4) Ask the data governance team to go to the Semantic Model Settings > Endorsement and Discovery, and click Certify
- 5) Ask a business user to visit the OneLake Data Hub to validate they can see the semantic model.

Question

You join a new company and are given a Power BI report theme (JSON file) to use for all new projects.

How do you apply this JSON file theme to the report you are developing?

A) In Power BI Desktop, go to View > Themes > Customize current theme

B) Go to Fabric Admin Portal > Custom branding and set the Default Report Theme

C) Use Tabular Editor 2 to update the theme.

D) In Power BI Desktop, go to View > Themes > Browse for themes

Fabric Analytics Concepts Notes

- **Azure Data Lake Storage (ADLS)**
 - A scalable, secure **cloud storage service** optimized for big data analytics.
 - Stores structured/unstructured data (e.g., files, logs, images) in a hierarchical namespace.
 - Acts as a *data lake* for raw data ingestion and batch/stream processing.
- **Microsoft Fabric**
 - An **end-to-end analytics platform** unifying data engineering, warehousing, science, and BI.
 - Integrates compute (Spark, SQL), storage (OneLake), and services (Data Factory, Power BI) into a single SaaS solution.
 - Built on **OneLake** (a unified data lake for all Fabric workloads).

Feature	ADLS	Microsoft Fabric
Storage	Hierarchical storage (Gen1/Gen2)	OneLake (ADLS Gen2-based, unified for all Fabric workloads)
Compute	Requires external services (e.g., Databricks, Synapse)	Built-in Spark, SQL, and real-time analytics
Data Integration	Needs Azure Data Factory/Synapse	Built-in Data Factory (Data Pipeline)
Governance	Manual setup (RBAC, ACLs)	Centralized governance (Purview integration)
Data Warehousing	Not natively supported	Built-in Warehouse (SQL-based)
BI & Visualization	Requires Power BI	Native Power BI integration
Pricing	Pay-as-you-go (storage + egress)	Unified capacity-based pricing

Fabric Analytics Concepts Notes

3. When to Use Which?

- **Choose ADLS if you need:**
 - A pure storage layer for raw data.
 - Flexibility to use external compute tools (e.g., Databricks, Synapse).
 - Long-term archival or low-cost storage.
- **Choose Fabric if you need:**
 - An all-in-one platform for analytics (storage + compute + BI).
 - Simplified collaboration (OneLake, shared datasets).
 - AI/ML integration (e.g., Synapse Data Science).

4. Integration

- **Fabric uses ADLS Gen2 under the hood** (via OneLake).
- You can **connect ADLS to Fabric** for existing data lakes.

Example Workflow:

1. **ADLS:** Ingest raw data.
2. **Fabric:** Transform (Spark), model (Warehouse), and visualize (Power BI).

Question

You have 1000 JSON files stored in Azure Data Lake Storage (ADLS) Gen2 that you want to bring into Fabric.

The ADLS Gen2 storage is secured using a virtual network.

Which of these actions would you need to perform FIRST?

- A) In Fabric, go to 'Manage connections and gateways' and then click 'Create a virtual network data gateway'.
- B) Create a shortcut to the ADLS Gen2 storage account.
- C) **In Azure, register a new Resource Provider, create a Private Endpoint and subnet.**
- D) Install an on-premise data gateway on an Azure Virtual Machine in the same virtual network.
- E) Enable public access in the storage account network settings

To securely bring 1000 JSON files from **Azure Data Lake Storage (ADLS) Gen2 (secured via a virtual network)** into **Microsoft Fabric**, the first step is:

Fabric Analytics Concepts Notes

Correct Answer: C)

"In Azure, register a new Resource Provider, create a Private Endpoint and subnet."

Why?

1. **ADLS Gen2 is secured via a virtual network** (private, no public access).
2. To access it from Fabric, you need a **Private Endpoint** to establish a secure, private connection between Fabric and ADLS *without* exposing data to the public internet.
3. This requires:
 - o Registering the **Microsoft.Synapse** Resource Provider (if not already done).
 - o Creating a **Private Endpoint** in Azure to link Fabric to the ADLS virtual network.
 - o Configuring a **subnet** for the endpoint.

Why Not the Others?

- **A)** Virtual Network Data Gateway is for **on-premises** data sources, not Azure-native services like ADLS.
- **B)** Shortcuts require connectivity *first* (which the Private Endpoint enables).
- **D)** On-premises gateways are unnecessary for Azure-to-Azure connections.
- **E)** Enabling public access violates security best practices (question states the storage is private).

Steps to Implement (After Option C):

1. Create the Private Endpoint in Azure.
2. In Fabric, create a **shortcut** to ADLS Gen2 (using the private connection).
3. Ingest JSON files via Spark or Data Pipeline.

This ensures secure, network-compliant data transfer.

Fabric Analytics Concepts Notes

Question

You have data stored in tables in Snowflake.

Which of the following CANNOT be used to bring the data into Fabric?

- A) Use the Data Pipeline CopyData activity.
- B) Create a Shortcut to the Snowflake tables from your Lakehouse.
- C) Use a Dataflow Gen2 with a Snowflake connection.
- D) Use database mirroring to create a mirrored Snowflake database in Fabric.

In Microsoft Fabric, **shortcuts** are currently supported only for **OneLake data storage locations** like Azure Data Lake or other Lakehouses.

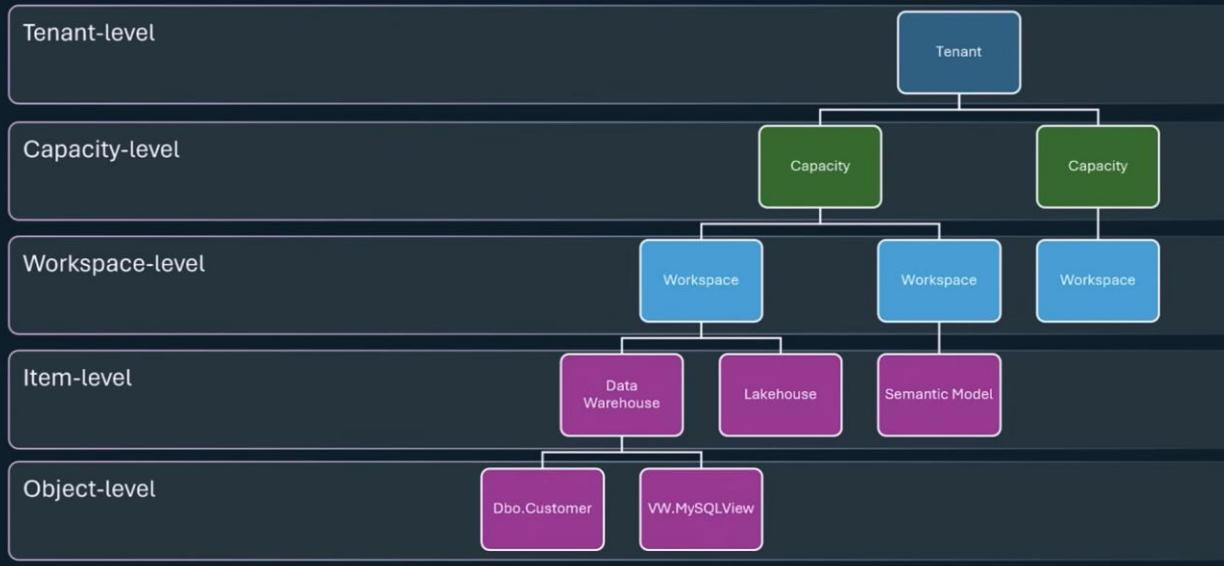
Snowflake is not a supported source for shortcuts — you cannot directly create a shortcut to a Snowflake table from a Lakehouse.

- A) Use the Data Pipeline CopyData activity
 - ✓ This is a standard way to copy data from Snowflake into Fabric via data pipelines.
- C) Use a Dataflow Gen2 with a Snowflake connection
 - ✓ Dataflows Gen2 support Snowflake connectors, enabling data ingestion and transformation.
- D) Use database mirroring to create a mirrored Snowflake database in Fabric
 - ✓ Fabric supports **database mirroring** for Snowflake as one of its real-time data integration methods.

Fabric Analytics Concepts Notes

Ch3. Implement & Manage a Data Analytics Env

The structure of a Fabric implementation



Fabric administration

Level	Administrator	Role required	Where admin happens	
Tenant-level Tenant	Tenant Admin	Entra ID role of Global administrator, Power Platform administrator or Fabric administrator	Fabric Admin Portal (Tenant Settings)	→ Last lesson
Capacity-level Capacity	Capacity Admin	Assigned in Azure when provisioning the capacity (must be person or Service Principal in Entra ID tenant) – can also be updated in Azure Portal	Azure Portal & Fabric Admin Portal (Capacity Settings)	→ This lesson
Workspace-level Workspace	Workspace Admin	Person or group with the Workspace Role of 'Admin'	Workspace Settings (in the workspace)	→ This lesson

Fabric Analytics Concepts Notes

Capacity administrator settings

Capacity administration tasks in Azure

- 1. Creation of a new capacity
- 2. Deleting a capacity
- 3. Changing the size of a capacity
- 4. Changing the capacity administrator

* plus pausing and resuming a capacity

Capacity administration tasks in Fabric Capacity Settings

- 1. Enable Disaster Recovery
- 2. View capacity usage report
- 3. Define who can create workspaces
- 4. Define who is a Capacity Administrator
- Workspace creation permissions
- 5. Update Power BI connection settings from/to this capacity
- 6. Permit workspace admins to size their own custom Spark pools based on workspace compute requirements.
- 7. Assign workspaces to the Capacity

Capacity settings in Azure Portal

The screenshot shows the Microsoft Azure portal interface. At the top, there's a navigation bar with icons for search, notifications, and account settings. Below the bar, the 'Azure services' section is visible, featuring icons for creating a resource, Microsoft Fabric, Microsoft Entra ID, Storage accounts, Subscriptions, Azure SQL, SQL databases, Quickstart Center, Virtual machines, and More services. The main area is titled 'Resources' and shows a list of recent resources. The 'Recent' tab is selected, displaying the following table:

Name	Type	Last Viewed
outputtfabric	Storage account	2 weeks ago
fabric-sql	Resource group	2 weeks ago
Azure subscription 1	Subscription	3 weeks ago
learnmsf	SQL server	4 weeks ago
fabrictest	SQL database	4 weeks ago
FabricIngestor	Resource group	2 months ago

At the bottom of the resource list, there's a link 'See all'.

Fabric Analytics Concepts Notes

Home >

Create Fabric capacity

Project details
Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize your resources.

Subscription * ⓘ Azure subscription 1
Resource group * ⓘ FabricIngeste Create new

Capacity details
Name your Capacity and select a location.

Capacity name * ⓘ fabric2learn
Region * UK South

Size ⓘ F64
64 Capacity units Change size

Fabric capacity administrator * ⓘ william Select

Review + create < Previous Next: Tags >

Select

Select the resource size

SKU	Capacity Units
F2	2
F4	4
F8	8
F16	16
F32	32
F64	64
F128	128
F256	256

Prices presented here are estimates in your local currency that include only A infrastructure costs and any subscription or location discounts. Final charges provided in your local currency, in cost analysis and billing views. [View the Az calculator.](#)

Costs for F2



Product Details

Microsoft Fabric
by Microsoft
[Terms of use](#) | [Privacy policy](#)

COST (ESTIMATED/MONTH) ⓘ
US\$306.60

Capacity settings in Fabric Portal

Fabric Analytics Concepts Notes

The Microsoft Fabric Home page displays the following tiles:

- 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.
- 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**: Rapidly ingest, transform, and analyze data from 1 GB to 1 PB.

The Microsoft Fabric Admin portal shows the following sidebar navigation:

- Tenant settings (New)
- Usage metrics
- Users
- Premium Per User
- Audit logs
- Domains (New)
- Capacity settings** (selected)
- Refresh summary
- Embed Codes
- Organizational visuals
- Azure connections
- Workspaces
- Custom branding
- Protection metrics
- Featured content
- Help + support

The Capacity settings page details the Delegated tenant settings, including:

- Disaster Recovery
- Capacity usage report
- Notifications
- Contributor permissions (Disabled for the entire organization)
- Admin permissions
- Power BI workloads
- Preferred capacity for My workspace
- Data Engineering/Science Settings
- Workspaces assigned to this capacity

Fabric Analytics Concepts Notes

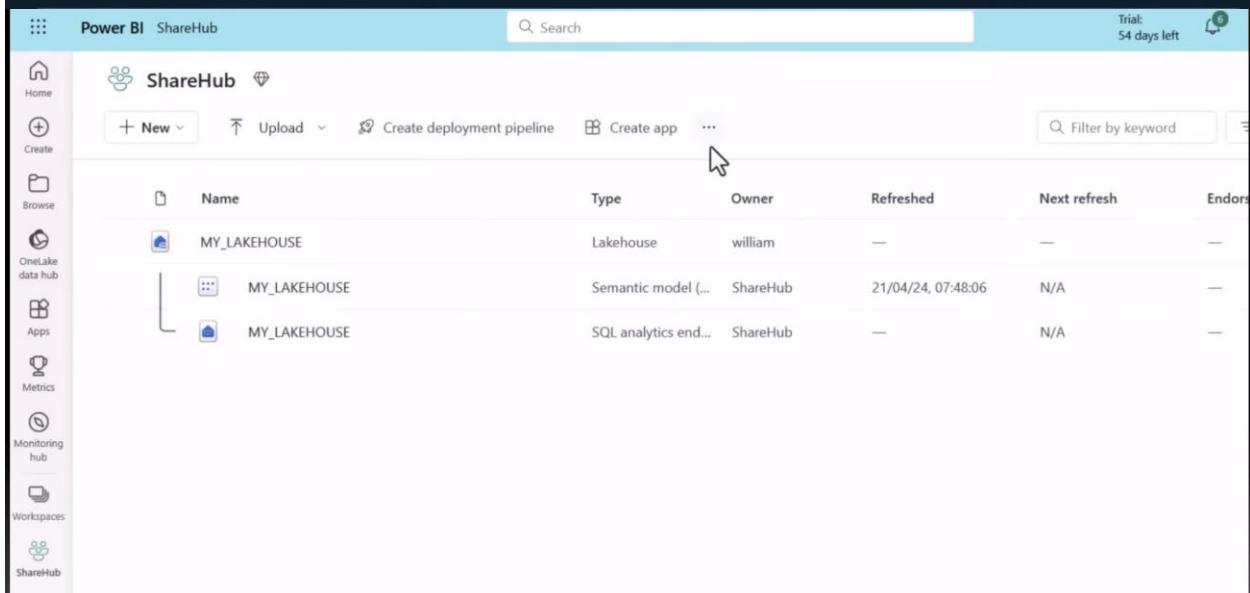
Workspace administrator settings

Workspace administration tasks in Workspace Settings

1. Edit license for the workspace (e.g. Pro, PPU, Fabric, Trial etc)
2. Configure Azure connections
3. Configure Azure DevOps connection (Git)
4. Setup workspace identity
5. Power BI Settings
6. Spark Settings

Note: managing access is done through 'Manage Access', not Workspace Settings.

And remember, it's not just Workspace Admins that can give people access to a Workspace – we'll explore this in more detail shortly.



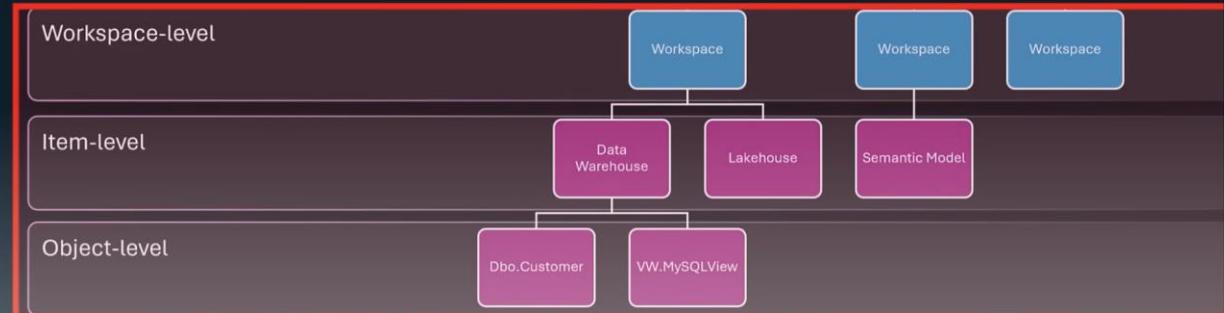
The screenshot shows the Power BI ShareHub interface. The left sidebar contains navigation links: Home, Create, Browse, OneLake data hub, Apps, Metrics, Monitoring hub, Workspaces, and ShareHub (which is selected). The main area is titled "ShareHub" and displays a list of datasets. The columns are: Name, Type, Owner, Refreshed, Next refresh, and Endorsed. There are three entries, all named "MY_LAKEHOUSE". The first entry is a Lakehouse type owned by "william". The second is a Semantic model owned by "ShareHub" with a refresh timestamp of "21/04/24, 07:48:06" and "N/A" for next refresh. The third is a SQL analytics endpoint owned by "ShareHub" with "N/A" for both refresh fields.

Name	Type	Owner	Refreshed	Next refresh	Endorsed
MY_LAKEHOUSE	Lakehouse	william	—	—	—
MY_LAKEHOUSE	Semantic model (...)	ShareHub	21/04/24, 07:48:06	N/A	—
MY_LAKEHOUSE	SQL analytics end...	ShareHub	—	N/A	—

The structure of a Fabric implementation

Sharing things in Fabric can be done at one of these three levels.

(Object-level sharing is not assessed as part of the exam)



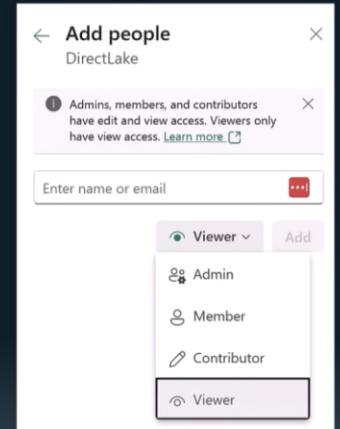
Workspace-level sharing

People or groups can be given workspace-level access.

When sharing, the person or group is assigned a workspace role:

- Admin
- Member
- Contributor
- Viewer

This role applies to **all items in the workspace**. For example a Viewer in the workspace will be able to View all items in the workspace.



Fabric Analytics Concepts Notes

Capability	Admin	Member	Contributor	Viewer
Update and delete the workspace.	✓			
Add or remove people, including other admins.	✓			
Add members or others with lower permissions.	✓	✓		
Allow others to reshare items. ¹	✓	✓		
Create or modify database mirroring items.	✓	✓		
Create or modify warehouse items.	✓	✓		
Create or modify SQL database items.	✓	✓		
View and read content of data pipelines, notebooks, Spark job definitions, ML models and experiments, and eventstreams.	✓	✓	✓	✓
View and read content of KQL databases, KQL query-sets, and real-time dashboards.	✓	✓	✓	✓
Connect to SQL analytics endpoint of Lakehouse or the Warehouse	✓	✓	✓	✓
Read Lakehouse and Data warehouse data and shortcuts ² with T-SQL through TDS endpoint (ReadData).	✓	✓	✓	✓
Read Lakehouse and Data warehouse data and shortcuts ² through OneLake APIs and Spark (ReadAll).	✓	✓	✓	
Read Lakehouse data through Lakehouse explorer (ReadAll).	✓	✓	✓	
Subscribe to OneLake events.	✓	✓	✓	
Write or delete data pipelines, notebooks, Spark job definitions, ML models, and experiments, and eventstreams.	✓	✓	✓	
Write or delete Eventhouses ³ , KQL Querysets, Real-Time Dashboards, and schema and data of KQL Databases, Lakehouses, data warehouses, and shortcuts.	✓	✓	✓	
Execute or cancel execution of notebooks, Spark job definitions, ML models, and experiments.	✓	✓	✓	
Execute or cancel execution of data pipelines.	✓	✓	✓	
View execution output of data pipelines, notebooks, ML models and experiments.	✓	✓	✓	✓
Schedule data refreshes via the on-premises gateway. ⁴	✓	✓	✓	
Modify gateway connection settings. ⁴	✓	✓	✓	

Fabric Analytics Concepts Notes

Workspace-level access example

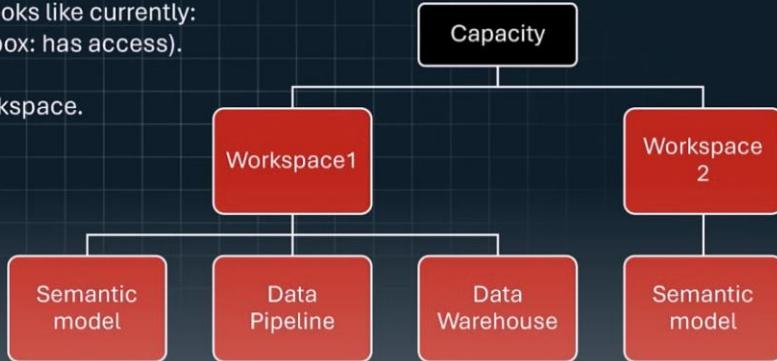


This is John.
He is a business analyst for Camilla.
Camilla has asked you to give him Contributor access to Workspace 1.

No access
Has access

This is what John's access looks like currently:
(Red box: no access, green box: has access).

You are an Admin in the Workspace.



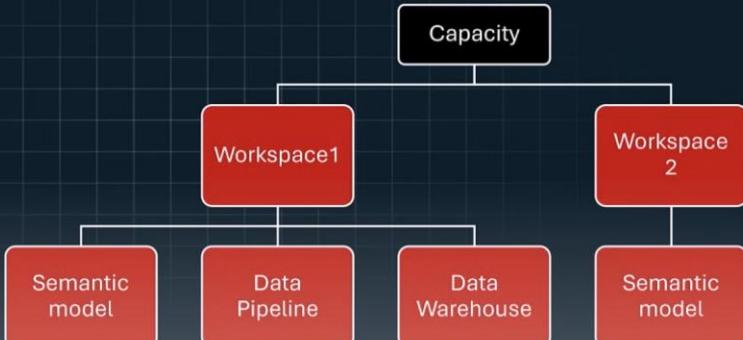
Workspace-level access example

What steps would you take?

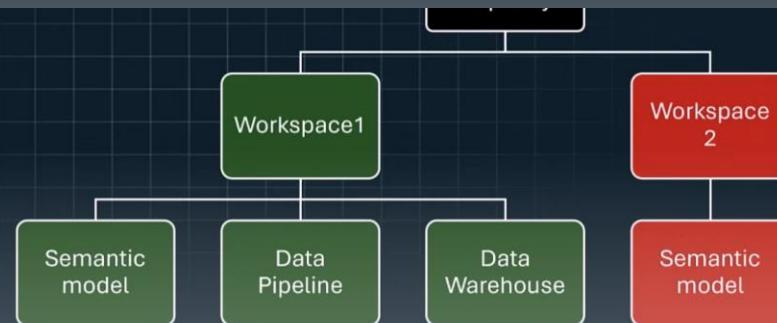
1. Does John fit into an existing security group that currently has Contributor access to the workspace?
2. If no, could you create an 'Analysts' security group, and add John to it? This would future-proof his access (and others like him).

No access
Has access

You create an Analysts security group, add John to the security group, and give the group Contributor access to Workspace1. What changes in the picture on the right?



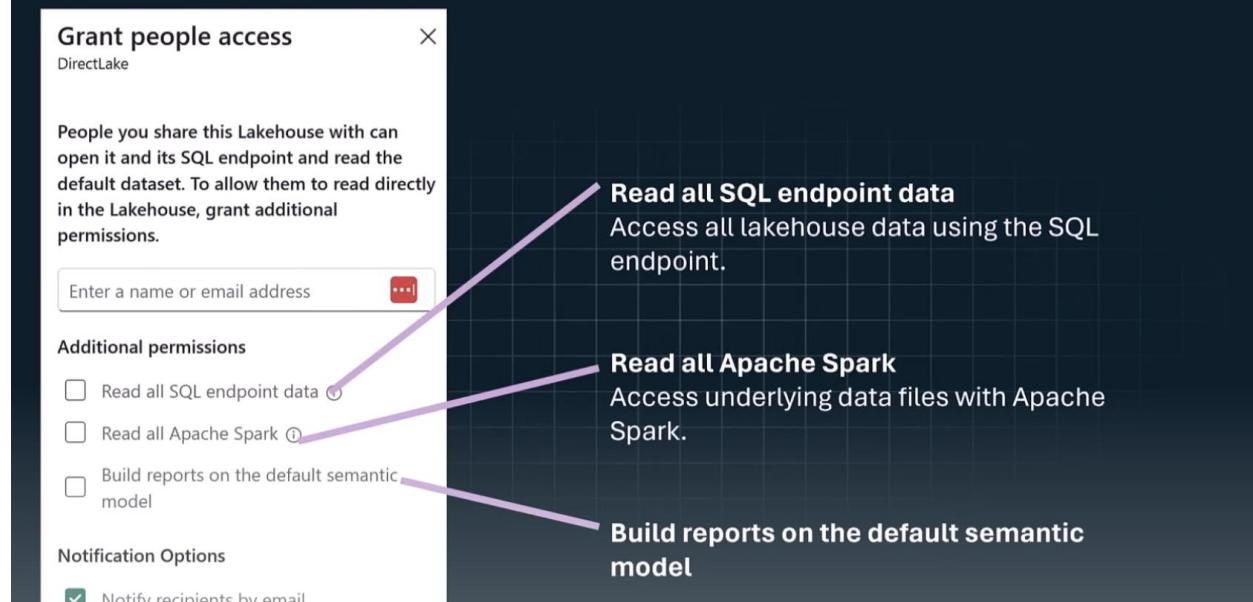
John's access after the changes have been implemented.



Fabric Analytics Concepts Notes



LAKEHOUSE: Additional permissions when sharing



Fabric Analytics Concepts Notes

OneLake Data Access Model (RBAC) – preview

A SQL analytics endpoint for SQL querying and a default Power BI semantic model for faster reporting were created and will be updated with any...

Manage OneLake data access (preview)

Explorer

businesses

Showing 1000 rows

	ABC business_id	ABC name	ABC address	ABC city	ABC state
1	kqZNMEhQEsg...	Basements Love ...	3653 Salmon St	Philadelphia	NJ
2	AcDQ4C7Plc0Ioj...	Herbiary	51 N 12th St	Philadelphia	PA
3	vi5NzBtUaBoYf...	Pennsylvania Gen...	51 N 12th St	Philadelphia	PA
4	DS-aX6GzBVwRp...	Condiment	51 N 12th St	Philadelphia	PA
5	AMWELaiyJGf5W...	Mueller Chocolat...	51 N 12th St	Philadelphia	PA
6	GqE-dUBnMmvjz...	The Head Nut	51 N 12th St	Philadelphia	PA
7	mSyqu7A0TrGMe...	Golden Fish Mar...	51 N 12th St	Philadelphia	PA
8	PtFZJdyXzxVO_V...	Profi's Crêperie	51 N 12th St	Philadelphia	PA
9	RLXT0560avopH...	Fair Food Farmst...	51 N 12th St	Philadelphia	PA
10	n5Fdf9RzooS7ob...	Four Seasons Jui...	51 N 12th St	Philadelphia	PA
11	xVExbX_AWD3jlz...	Downtown Chees...	51 N 12th St	Philadelphia	PA
12	hsl36MhSVSpH3...	Martin's Quality ...	51 N 12th St	Philadelphia	PA

Question

Toby creates a new workspace with some Fabric items to be used by Data Analysts.

A) Viewer

Toby creates a new security group called 'Data Analysts'. He includes himself as a member of this security group.

B) Member

Toby gives the Data Analysts security group a **Viewer role** within the workspace.

C) Admin

D) Contributor

What workspace role does Toby have?

The Admin role supersedes and includes all permissions of the Viewer role.

Fabric Analytics Concepts Notes

Question

[..case study continued...]

Toby wants to delegate some of the management responsibilities in the workspace.

He wants to give this person the ability to share content within the workspace, invite new Contributors to the workspace, but not add new Admins to the workspace.

Which role should Toby give this person?

A) Admin

B) Member

C) Contributor

D) Viewer

Question

You have Admin role in a workspace.

Shiela is a data engineer in your team. Currently she has no access to the workspace.

Sheila needs to update a data transformation script in a PySpark notebook. The script gets data from a Lakehouse table, cleans it and then writes it to a table in the same Lakehouse.

You want to adhere to the principle of least privilege.

What actions should you take to enable this?

A) Give Sheila the Contributor role in the workspace.

B) Share the Lakehouse item with Read All Spark Data permission

C) Give Sheila the Admin role in the workspace

D) Share the Lakehouse item with Read All Spark Data permission and share the Notebook with Edit permission

Fabric Analytics Concepts Notes

Question

You have Admin role in a workspace.

You want to pre-install some useful Python packages to be used across all notebooks in your workspace.

How do you achieve this?

- A) In the Fabric Admin Portal, go to Spark Settings and install the libraries
- B) Go to Workspace Settings > Spark Settings > Library Management
- C) Create an Environment, install the packages in the Environment, then go to Workspace Settings > Spark Settings and Set the Default Environment
- D) Go to Capacity Settings > Default Libraries

Ch4. Manage analytics dev lifecycle

The screenshot shows the Azure DevOps interface for a repository named 'dp-600-practice'. The left sidebar includes links for Overview, Boards, Repos (selected), Files, Commits, Pushes, Branches, Tags, Pull requests, Advanced Security, Pipelines, Test Plans, and Artifacts. The main content area displays the message 'dp-600-practice is empty. Add some code!'. Below this, there's a 'Clone to your computer' section with 'HTTPS' selected, showing the URL 'https://fabric-university@dev.azure.com/fabric-university/dp-600-practice'. There's also an option to 'Generate Git Credentials'. A note at the bottom of this section says: 'Having problems authenticating in Git? Be sure to get the latest version [Git for Windows](#) or our plugins for [IntelliJ](#), [Eclipse](#), [Android Studio](#) or [Windows command line](#)'. Below this is a 'Push an existing repository from command line' section with 'HTTPS' selected, showing the command 'git remote add origin https://fabric-university@dev.azure.com/fabric-university/dp-600-practice/_git/dp-600-practice'. At the bottom, there's an 'Import a repository' section with an 'Import' button and a video thumbnail of a man speaking. A note next to the video says 'Initialize & main branch with a README or'. The top navigation bar shows 'Azure DevOps fabric-university / dp-600-practice / Repos / Files / dp-600-practice' and a search bar.

Restrict who can make changes to main branch of repo

Fabric Analytics Concepts Notes

The screenshot shows the Azure DevOps interface for managing project settings. The main window displays 'All Repositories' under the 'Policies' tab, listing various repository policies such as Commit author email validation, File path validation, Case enforcement, Reserved names, Maximum path length, and Maximum file size. Each policy has a toggle switch set to 'Off'. A video overlay of a person speaking is visible in the background of the main pane. Below this, another section shows 'Branch Policies' for protecting branch namespaces across all repositories in the project.

Project Settings
dp-600-practice

All Repositories

Repository Policies

- Commit author email validation**: Off
- File path validation**: Off
- Case enforcement**: Off
- Reserved names**: Off
- Maximum path length**: Off
- Maximum file size**: Off

Add branch protection

Branches to protect

- Protect the default branch of each repository
- Protect current and future branches matching a specified pattern

Git/ version control summary of key points (for the exam)

Version control with Git allows:

- track changes made to Fabric items
- revert to older versions of an item
- Multiple users can collaborate on the same Fabric item (for example a Power BI report), at the same time, and their changes can be merged together (assuming no conflicts).
- implement a check and approval process for approving changes made to Fabric items.



The following items are currently supported:

- [Data pipelines](#)
- [Lakehouse](#)
- [Notebooks](#)
- [Paginated reports](#)
- Reports (* excluding AAS/SSAS connected reports)
- Semantic models (except push datasets, live connections, model v1, and semantic models [created from](#) the Data warehouse/lakehouse.)

Well actually, let's start with 'deployment'

Because for many coming from the world of analytics, 'deployment' is a new concept.

Rather than having just one copy of a Power BI report (the production copy), instead we have (typically) three:

- a development version (the version you use when you are making changes to a Live report)
- A TEST version, which you might give to colleagues/ a client to test/ review (or do some automated testing)
- The production report (which is public and shared with stakeholders).

Microsoft released a feature called deployment pipelines to help manage these three (or more!) environments. Let's [redacted] what you can do today...



You worked in prod env in ergo & edeka projects

Setup deployment pipeline in MS Fabric

Testing: Integration test, unit tests, data validation tests

Fabric Analytics Concepts Notes

Deployment pipelines

Manage your workspace content through deployment stages

The screenshot shows the Fabric Analytics interface. A modal window titled "Customize your stages" is open, prompting the user to define stages for a pipeline. The stages listed are "Development", "Test", and "Production". Each stage has a trash icon to its right. Below the modal, three deployment pipeline stages are visible: "dp600-dev", "dp600-test", and "dp600-prod". Each stage has sections for "Warehouses", "Lakehouses", and "Data pipelines". A "Deploy" button is present at the bottom of each stage section. The background shows a video feed of a person.

Summary of deployment pipelines (for the exam)

- The overall goal is add layers of control when developing and deploying new Fabric items (or making changes to existing items).
- Ultimately to ensure new things you develop are not going to break your existing analytics solutions.
- Normally includes Development, Test/Staging, Production stages.
- 'Deployment' using deployment pipelines involves copying items from one workspace to another.
- Deployment rules can be implemented to the Default Lakehouse (for a notebook) like stages.

As well as the Deployment Pipelines functionality there are other ways to manage deployment:

- Can be managed through Branching
- Can be managed through Azure DevOps Pipelines (YAML templates)
- For semantic models, you can do it using the XMLA endpoint (see below for more info).

Other options for deployment

CI/CD for Microsoft Fabric Data Warehouses using YAML Pipelines

Published by [Kevin Chant](#) on October 25, 2023

Reading Time: 6 minutes

In this post I cover how to perform [CI/CD for Microsoft Fabric Data Warehouses using YAML pipelines](#). Which can now be done gracefully with the new target platform thanks to a new `SqlPackage` update.

CI/CD for Microsoft Fabric Data Warehouses using Azure DevOps

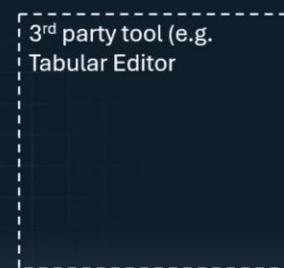
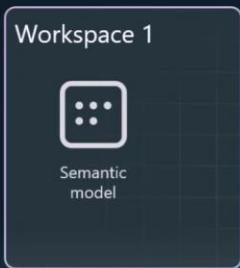
Published by [Kevin Chant](#) on October 23, 2023

Reading Time: 5 minutes

In this post I want to cover [CI/CD for Microsoft Fabric Data Warehouses using Azure DevOps](#). Which can now be done gracefully with the new target platform thanks to a new `SqlPackage` update.

Deploying semantic models using XMLA endpoint

Broadly, there's two ways to create and manage semantic models



XLMA (XML for Analysis) Endpoint is a legacy SOAP-based protocol used for interacting with SQL Server Analysis Services (SSAS) and other OLAP (Online Analytical Processing) data sources. It allows client applications to execute queries (MDX, DAX, or XMLA commands) against multidimensional or tabular models.

- Protocol: Uses XMLA over HTTP/HTTPS (SOAP-based).
- Supported SSAS Models: Multidimensional (OLAP) models & Tabular models (compatibility level 1200+)

When is XLMA Endpoint Used Today?

- Power BI Premium/PPU: Power BI uses XMLA endpoints for:
 - Advanced dataset management (via SSMS/Tabular Editor).
 - Automated deployments (using PowerShell or DevOps).
 - Use Cases: Executing MDX (Multidimensional Expressions) **or DAX (Data Analysis Expressions) queries**. Performing administrative tasks (e.g., processing cubes, deploying models). Used by tools like Excel, Power BI, SQL Server Management Studio (SSMS), and custom apps.
- Legacy SSAS Systems: Older on-premises SSAS deployments still rely on it.

Fabric Analytics Concepts Notes

- Third-Party Tools: Some BI tools connect via XMLA for live queries.

Connecting to your workspace (XMLA)

Workspace settings

Connection link
Use this link to connect third-party software to the workspace. Copy the link and add it to your third-party software.

powerbi://api.powerbi.com/v1.0/myorg/PBi%20to%20Fabric%20End-to-end

Feature	.pbix	.pbip	.pbit	.pbids
Full Name	Power BI Desktop File	Power BI Project File	Power BI Template File	Power BI Data Source File
Contains Data?	✓ Yes (embedded)	✗ No (Live/DQ only)	✗ No (template)	✗ No (connection only)
Git-Friendly?	✗ No (binary)	✓ Yes (folder-based)	✗ No	✓ Yes (lightweight)
Editable in PBI Desktop?	✓ Yes	✓ Yes	✓ Yes (prompts for data)	✓ Yes (opens as connection)
Used for Publishing?	✓ Yes	✗ No (must publish .pbix)	✗ No	✗ No
Primary Purpose	Standalone reports	Team projects (DevOps)	Reusable templates	Quick data source setup

- .pbix → Standard Power BI reports (single-user or small teams).
- .pbip → Team projects with Git/DevOps (future of Power BI development).
- .pbit → Distributing reusable templates (e.g., company-wide dashboards).
- .pbids = Just a connection (no report).

Fabric Analytics Concepts Notes

Impact Analysis Tool of Fabric (Lineage View)

The screenshot shows the Impact Analysis Tool interface for a "PBi to Fabric End-to-end" pipeline. On the left, a sidebar lists various fabric components: Home, Create, Browse, Workspaces, PBi to Fabric End-to-end (selected), ETL Notebook, datapipeline, MyDataPipeline, Competitive Marketing..., and The main area displays a lineage graph starting from "Azure Data Lake Storage Gen2" (https://outputfabric.dfs.core.windows.net/_william) and "Azure Data Lake Storage Gen2" (Output Fabric william). These feed into a central "LH_BRONZE" node labeled "Lakehouse". From this lakehouse, arrows point to "Notebook 1" and "Notebook 2". Further downstream, "Notebook 1" connects to "LH_BRONZE SQL analytics endpoint" and "LH_BRONZE Semantic model (default)". "Notebook 2" connects to "Warehouse" and "DF_LoadDataFromADLSToDWH" (Dataflow Gen2). A video overlay of a person speaking is visible in the bottom right corner.

Impacted by this Lakehouse

LH_BRONZE

Child Items All downstream items

3 Impacted child items ⓘ 1 Workspaces

Browse by item type

SQL analytics endpoint 1 LH_BRONZE Notebooks 2

Making changes? Notify the contact lists of workspaces that contain child items of LH_BRONZE.

Notify contacts

Fabric Analytics Concepts Notes

Question

You are looking to improve the efficiency and consistency of your Power BI Development team.

You want each report created by the team to always consist of three pages: Intro, Context, Analysis.

The reports should always align to the company branding guidelines.

Which of the following would help you achieve this

A) Create .pbip file

B) Create a .pbix file

C) **Create a .pbit file**

D) Create a .pbids file

E) Using a JSON custom report theme



Question

Which of the following most accurately describes Git?

A) To use Git, you must be using GitHub.

B) Git is a Microsoft product for tracking changes made to Fabric items.

C) **Git is a open-source version control system that tracks changes in any set of (text-based) files.**



It allows us to add deployment rules to Fabric deployment pipelines.

Question

In a Azure DevOps repo, the MAIN branch is 'protected' (needs approval before any changes are merged into it).

The repo contains one PBIP. You have to update the Title in the report, merging these changes to the Main branch.

In which order should you carry out the following tasks to achieve this?

1. **Clone the repository to your local machine.**

2. **Checkout a new feature branch from the MAIN branch**

3. **Make the required changes to the report.**

4. **Commit and push the feature branch.**

5. **Open a pull request in Azure Repos.**

6. **Wait for approval, then merge into the main**



Fabric Analytics Concepts Notes

Question

You want to deploy a semantic model using the XMLA endpoint.

Where can you do to find the XMLA endpoint to set up a connection with a third-party tool?

A) Go to the workspace settings for the workspace you want to deploy your model to

B) Go to the Fabric Admin Portal > Capacity Settings

C) In your workspace, find your semantic model and click on settings to get the XMLA endpoint address

D) In the Azure Portal, in your Fabric Capacity, go to the XMLA endpoint connection string settings.



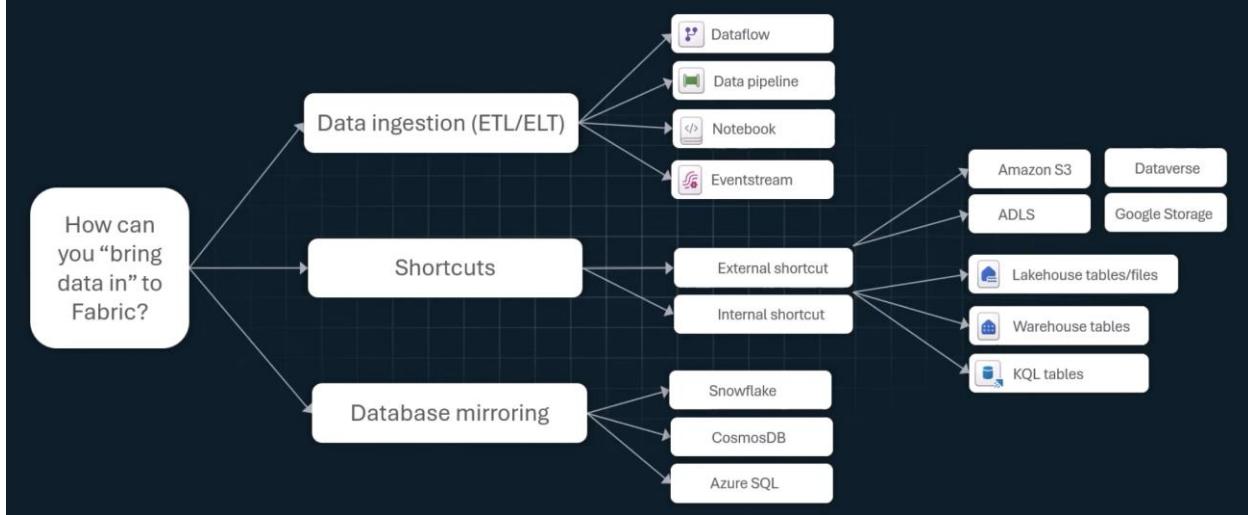
1. The XMLA endpoint is **workspace-specific** and tied to the underlying capacity.
 - o Open the **Power BI Service** (app.powerbi.com) & Navigate to the **workspace** containing your semantic model.
 - o Format: powerbi://api.powerbi.com/v1.0/[tenant-name]/[workspace-name]
2. To deploy a semantic model via XMLA (e.g., using **SSMS, Tabular Editor, or DevOps**), you need the endpoint address from:
 - o **Workspace Settings > Premium/ Fabric Capacity tab > XMLA Endpoint** field.
 - o Requires Premium/Fabric Capacity: XMLA endpoints are only available for workspaces on Premium Per User (PPU), Premium, or Fabric SKUs.
 - o Third-Party Tools: Use the endpoint to connect via: SQL Server Management Studio (SSMS) (or) Tabular Editor (or) PowerShell (for automation)

Fabric Analytics Concepts Notes

Ch5. Getting Data to Fabric

- Ingest data by using a data pipeline, dataflow, or notebook
- Copy data by using a data pipeline, dataflow, or notebook
- Choose an appropriate method for copying data from a Fabric data source to a lakehouse or warehouse
- Create and manage shortcuts

Ingest data by using a data pipeline, dataflow, or notebook



Ingest data with a dataflow

You can use 150+ connectors to external systems to bring data into a familiar Power Query low/no-code interface, perform data transformation and then write to Fabric data stores.

Pros (when you should consider using)	Cons (maybe when you shouldn't use)
<ul style="list-style-type: none">• To use 150+ external connectors• No/ low-code solution• Can do Extract, Transform AND Load• Accessing on-premise data (OPDG)• When you need to get more than one dataset at a time and combine them (although you might want to space this out to allow data validation).• Can upload raw local files (static files)	<ul style="list-style-type: none">• Can struggle with large datasets (although Fast Copy has recently been introduced which should speed up your ETL)!• Difficult to implement data validation• Currently can't pass in external parameters

Ingest data with a data pipeline

Primarily an orchestration tool (do this, then do that). Can also be used to get data into Fabric, using the Copy Data activity (and others!).

Pros (when you should consider using)	Cons (maybe when you shouldn't use)
<ul style="list-style-type: none">• Large datasets (although now Dataflow has Fast Copy, so performance should be comparable between the two).• Importing 'cloud' data (e.g. data in Azure)• When you need control flow logic• Triggering wide variety of actions in Fabric (and outside of Fabric), like Dataflows, Notebooks, Stored Procs, KQL Scripts, Webhooks, Azure Functions, Azure ML, Azure Databricks.	<ul style="list-style-type: none">• Can't do the Transform natively (but can embed notebooks or dataflow).• No ability to 'upload' local files• Does not work cross-workspace (currently - although this feature is now planned :))

Ingest data with a notebook

General purpose coding notebook which can be used to bring data into Fabric, via connecting to APIs or by using client Python libraries

Pros (when you should consider using)

- Extraction from APIs (using Python requests library, or similar!)
- To use client libraries (e.g. Azure libraries, or the Hubspot client library in Python to access Hubspot data).
- Good for code re-use (and can be parameterized)
- For data validation and data quality testing of incoming data
- The fastest in terms of data ingestion (and most efficient for CU spend - see here)

Cons (maybe when you shouldn't use)

- When you don't have a Python capability in your organisation
- When you want to write to a Data Warehouse (currently not possible)

Create and manage shortcuts

- A shortcut enables you to create a live link to data stored either in another part of Fabric (internal shortcut) or in external storage locations:
 - ADLS Gen2
 - Amazon S3
 - Other services that USE Amazon S3 for storage (like Cloudflare)
 - Google Cloud Storage
 - Dataverse
- Shortcuts can be setup for individual files, but more commonly to a folder. If you create a shortcut to a 'base' folder, it will monitor and sync all files in subfolders.
- Be careful of cross-region egress fees. If your Fabric capacity is in UK-South region (for example), but your ADLS storage account is in West-US, you will be charged for 'cross-region' egress fees (\$0.01 per GB)
- OneLake shortcuts can now be created using the [Fabric REST API](#)

Fabric Analytics Concepts Notes

The following table shows the shortcut-related permissions for each workspace role. For more information, see [Workspace roles](#).

 Expand table

Capability	Admin	Member	Contributor	Viewer
Create a shortcut	Yes ¹	Yes ¹	Yes ¹	-
Read file/folder content of shortcut	Yes ²	Yes ²	Yes ²	-
Write to shortcut target location	Yes ³	Yes ³	Yes ³	-
Read data from shortcuts in table section of the lakehouse via TDS endpoint	Yes	Yes	Yes	Yes

¹ Users must have a role that provides write permission the shortcut location and at least read permission in the target location.

² Users must have a role that provides read permission both in the shortcut location and the target location.

³ Users must have a role that provides write permission both in the shortcut location and the target location.

Feature	Lakehouse	Warehouse	Eventhouse
Data Types	All (structured/unstructured)	Structured only	Time-series/logs
Query Language	Spark SQL, T-SQL	T-SQL	KQL
Latency	Minutes-hours	Seconds-minutes	Sub-second
Best For	Raw data exploration	Business reports	Real-time streams

Some other decision criteria

- Need for 'real-time' data - if there is a requirement for near-real-time, you will prefer either database mirroring (for data in Azure SQL, Azure Cosmos DB, or Snowflake), or Shortcuts for Files and Folders.
- Skills in your team:
 - No/Low-code? Dataflow + data pipeline
 - SQL-based? Data pipeline + stored proc activity/ script activity
 - Python? Notebooks
- Cross-workspace limitations: data pipeline has quite a few limitations working 'cross-workspace'. This means that currently your data pipeline must be in the same workspace as your destination datastore (for data ingestion).
- Scalability / Size of your data / Cost / capacity usage - in general, the notebook is the most efficient method for data ingestion.

Question

You are trying to create a shortcut to a folder of CSV files in Azure Data Lake Storage Gen2.

Which of the following is a valid connection string you can connect to?

- A) <https://contosoadlscdm.database.windows.net>
B) <https://contosoadlscdm.blob.core.windows.net>
C) <https://contosoadlscdm.file.core.windows.net>
D) <https://contosoadlscdm.dfs.core.windows.net>

Key Difference: blob vs. dfs Endpoints

Feature	<code>blob.core.windows.net</code>	<code>dfs.core.windows.net</code>
Purpose	Blob Storage (flat namespace)	ADLS Gen2 (hierarchical folders)
Shortcuts	✗ No	✓ Yes
Folder Operations	Limited	Optimized (e.g., <code>ls</code> , <code>mkdir</code>)

Fabric Analytics Concepts Notes

- A) `database.windows.net` :
 - Used for **Azure SQL Database**, not ADLS Gen2.
- B) `blob.core.windows.net` :
 - Blob Storage endpoint (works for files, but lacks **hierarchical namespace** features like folders).
- C) `file.core.windows.net` :
 - For **Azure Files** (SMB/NFS shares), not ADLS Gen2.

3. Key Differences: HDFS vs. ADLS Gen2

Feature	HDFS (On-prem)	ADLS Gen2 (Cloud)
Namespace	Flat (emulated folders)	True hierarchical folders
Scalability	Manual cluster scaling	Automatic (unlimited storage)
Replication	3x (default)	Geo-redundancy (GRS, ZRS)
Access	<code>hdfs://</code> paths	<code>abfss://</code> (Azure Blob FS)

Why This Matters for Your Question

- The `dfs` endpoint in ADLS Gen2 (Option D) is the **cloud equivalent of HDFS** for Hadoop-based tools.
- When you use `dfs.core.windows.net`, you're leveraging:
 - HDFS-like paths (e.g., `/myfolder/data.csv`).
 - Compatibility with **Spark, Fabric, and Hadoop ecosystems**.



File Sharing Protocols	SMB (<i>Server Message Block</i>)	NFS (<i>Network File System</i>)
OS Support	Windows (native), Linux (Samba)	Linux/Unix (native), Windows (NFS client)
Security	Kerberos, AD integration/Entra id	IP-based, Kerberos (v4/v3)
Latency	Higher (chatty protocol)	Lower (lightweight)
Use Cases	Office docs, user home dirs	VMs, big data, media

Fabric Analytics Concepts Notes

Question

You are implementing the first stage in a medallion architecture.

Your goal is to retrieve data from a REST API (Get Request) and save the raw JSON response in the Files area of a BRONZE Lakehouse.

Which of the following methods can you use to achieve this?

(THREE CORRECT ANSWERS)

A) Data Pipeline (Copy Data Activity)

B) Dataflow (Web API connector)

C) Eventstream

D) Data Pipeline (Web Activity)

E) Fabric Notebook

Editor's Note: the Data Pipeline Web Activity can't perform the whole action on its own. It can perform the GET request and pass the output to a Copy Data activity to write the data into a Lakehouse Files area

Question

Your goal is to extract a CSV file in Azure Blob Storage write it to a Fabric Data Warehouse table.

Which of the following methods can you NOT use to achieve this?

A) Data Pipeline (Copy Data activity)

B) Dataflow with the Blob Storage connector

C) In a Fabric notebook, use the Azure Blob Storage client library for Python to get the file and write the Data Warehouse table

D) Use the COPY INTO statement (T-SQL) from within the Fabric Data Warehouse

Question

Workspace A contains Lakehouse A.

Workspace B contains Lakehouse B.

You want to create a shortcut in Lakehouse A (pointing to a table from Lakehouse B).

What is the minimum level of workspace permissions you need to achieve this?

A) Contributor in Workspace A and Viewer in Workspace B.

B) Contributor in Workspace A and Contributor in Workspace B.

C) Member in Workspace A and Contributor in Workspace B.

D) Viewer in Workspace A and Viewer in Workspace B

Question

One of your team is a superstar using M code for data extraction.

In which of following data extraction tools can you write M code?

A) Fabric notebook

B) In a T-SQL script in the Data Warehouse

C) Dataflow

D) Data pipeline (Mapping dataflow)

Fabric Analytics Concepts Notes

1. **A) Fabric Notebook**
 - **Primary Language:** Python, SQL, Scala, R.
 - **M Code Support?** ✗ No (M code is specific to Power Query, not notebooks).
2. **B) T-SQL Script in Data Warehouse**
 - **Primary Language:** T-SQL (Transact-SQL).
 - **M Code Support?** ✗ No (T-SQL is for querying relational data, not ETL transformations).
3. **C) Dataflow**
 - **Primary Language: M Code** (Power Query).
 - **M Code Support?** ✅ Yes (Dataflows are built on Power Query; M is the native language).
4. **D) Data Pipeline (Mapping Dataflow)**
 - **Primary Language:** Power Query (via UI) or Spark/SQL.
 - **M Code Support?** ✅ Partial (Mapping Dataflows use Power Query transformations, but M code is hidden behind the UI).

Ch6. SQL, DWH, Scheduling

View is read only (It can't update /insert just like function). To insert / update, use appropriate dml statement or stored procedure

```
CREATE VIEW dbo.vw_Employee_GetJack AS
SELECT
    EmployeeId
    , Name
    , Age
FROM dbo.Employees
WHERE Name like 'Jack%';

select * from dbo.vw_Employee_GetJack

-- KEY POINTS FOR A VIEW
-- [x] Transformed data is not stored in a view, the result is calculated at query time.
-- [x] You can create a View from another View (but be careful, this can lead to poor performance, and difficult to
-- [x] you cannot specify parameters for a View (this is possible with functions and stored procedures)
-- [x] Reading a SQL View into a semantic model will fallback to DirectQuery (Direct Lake not possible).
```

Fabric Analytics Concepts Notes

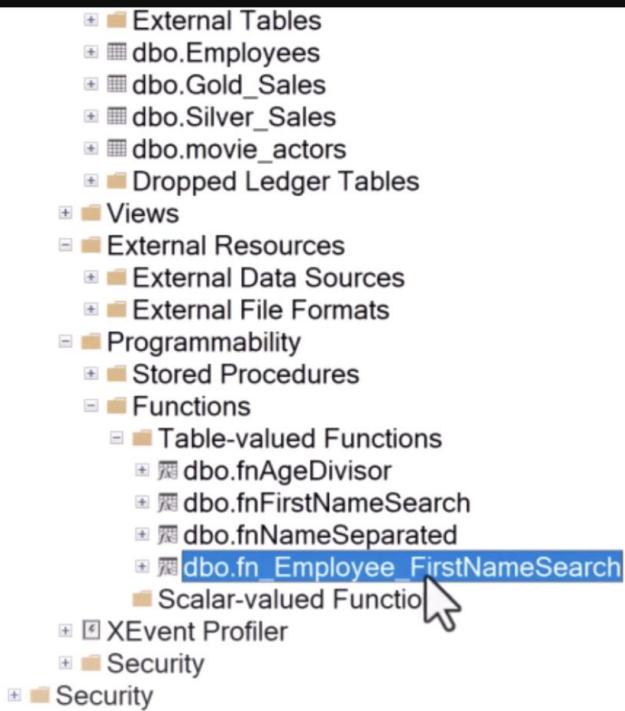
```
-- Create a Function that takes a first name as input
-- and Filters the Employees table
CREATE FUNCTION dbo.fn_Employee_FirstNameSearch ( @firstname varchar(20) = '' )
RETURNS TABLE
AS
RETURN
(
    SELECT
        EmployeeId
        , Name
        , Age
    from dbo.Employees
    WHERE Name like @firstname + '%'
) ;

-- Test the function
select * from dbo.fn_Employee_FirstNameSearch('Jack')
```

Function can take only input parameter & not output parameter (in mysql, sql server. Postgres function can take both though)

```
-- Key Points for a FUNCTION
-- [x] Useful for packaging up logic, which can be parameterized
-- [x] Only SELECT statements, can never edit the underlying data (using UPDATE/ INSERT etc)
-- [x] Can be called from DDL statements, or from within Stored Procedures, or a View
-- [x] Can't be orchestrated (directly) from a data pipeline (although can embed in a stored proc).
-- [x] Allows one or many input parameters
-- [x] Output type should also be a Table
```

Fabric Analytics Concepts Notes



- **DDL:**
 - Creating tables (CREATE TABLE).
 - Adding columns (ALTER TABLE).
 - Dropping databases (DROP DATABASE).
- **DML:**
 - Adding records (INSERT).
 - Updating salaries (UPDATE).
 - Querying data (SELECT).

1. DDL (Data Definition Language) Statements

DDL commands define/modify database structures.

CREATE TABLE

```
CREATE TABLE Employees (
    EmployeeID INT PRIMARY KEY,
    FirstName NVARCHAR(50) NOT NULL,
    LastName NVARCHAR(50) NOT NULL,
    HireDate DATE DEFAULT GETDATE(),
    Salary DECIMAL(10, 2) CHECK (Salary > 0)
);
```

ALTER TABLE

```
-- Add a column
ALTER TABLE Employees ADD DepartmentID INT;
```

Fabric Analytics Concepts Notes

```
-- Add a foreign key  
ALTER TABLE Employees  
ADD CONSTRAINT FK_Department  
FOREIGN KEY (DepartmentID) REFERENCES Departments(DepartmentID);  
CREATE INDEX  
CREATE INDEX IX_Employees_LastName ON Employees(LastName);  
DROP TABLE  
DROP TABLE Employees; -- Deletes the table and its data
```

2. Stored Procedures

Stored procedures encapsulate reusable SQL logic.

```
drop procedure if exists dbo.sp_Employee_GetByFirstName  
  
-- stored procedure with input parameter  
CREATE PROCEDURE dbo.sp_Employee_GetByFirstName @firstname varchar(20)  
as  
select * from dbo.Employees  
where Name like @firstname + '%'  
  
exec dbo.sp_Employee_GetByFirstName 'Jack'
```

Basic Stored Procedure

```
CREATE PROCEDURE GetEmployeeByID @EmployeeID INT  
AS  
BEGIN  
    SELECT * FROM Employees WHERE EmployeeID = @EmployeeID;  
END;
```

Execution:

```
EXEC GetEmployeeByID @EmployeeID = 101;
```

Procedure with Output Parameter

```
CREATE PROCEDURE GetEmployeeCount  
    @DepartmentID INT,  
    @Count INT OUTPUT  
AS  
BEGIN  
    SELECT @Count = COUNT(*)  
    FROM Employees  
    WHERE DepartmentID = @DepartmentID;  
END;  


#### Execution:



```
DECLARE @Result INT;
EXEC GetEmployeeCount @DepartmentID = 5, @Count = @Result OUTPUT;
PRINT @Result; -- Prints the count
```


```

Fabric Analytics Concepts Notes

Transaction Handling

```
CREATE PROCEDURE UpdateSalary
    @EmployeeID INT,
    @NewSalary DECIMAL(10, 2)
AS
BEGIN
    BEGIN TRY
        BEGIN TRANSACTION;
        UPDATE Employees
        SET Salary = @NewSalary
        WHERE EmployeeID = @EmployeeID;
        COMMIT TRANSACTION;
    END TRY
    BEGIN CATCH
        ROLLBACK TRANSACTION;
        THROW; -- Re-throws the error
    END CATCH
END;
```

3. Dynamic SQL in Stored Procedures

```
CREATE PROCEDURE SearchEmployees
    @ColumnName NVARCHAR(50),
    @SearchValue NVARCHAR(100)
AS
BEGIN
    DECLARE @SQL NVARCHAR(MAX);
    SET @SQL = N'SELECT * FROM Employees WHERE ' +
        QUOTENAME(@ColumnName) + ' = @Value';
    EXEC sp_executesql @SQL, N'@Value NVARCHAR(100)', @Value = @SearchValue;
END;
Execution:
EXEC SearchEmployees @ColumnName = 'LastName', @SearchValue = 'Smith';
```

- KEY POINTS
 - [x] Ability to define input AND OUTPUT parameters
 - [x] Are called using EXEC (not as part of a select statement)
 - [x] Can call other stored procedures
 - [x] Can also give access to people to the Stored Procedure
 - (and not the underlying dataset)
 - [x] Can be embedded in a data pipeline (with parameters passed)
 - [x] Ability to do INSERTS, UPDATES, DELETES

Fabric Analytics Concepts Notes

Azure DataFactory

The screenshot shows the Azure Data Factory interface for the workspace 'ETL-DP600'. The left sidebar includes options like Home, Create, Browse, OneLake data hub, Monitoring hub, and Workspaces. The main area displays a table of resources:

Name	Type	Owner	Refreshed	Next refresh	Endorsement	Sensitivity
dp600-dw	Warehouse	william	30/04/24, 19:22:02	N/A	—	—
dp600-dw	Semantic model (...)	ETL-DP600	29/04/24, 13:06:28	N/A	—	—
DP_Embroiding	Data pipeline	william	—	—	—	—
Load Data To DW Tables	Dataflow Gen2	william	01/05/24, 12:50:40	N/A	—	—
NB_Load_To_Silver	Notebook	william	—	—	—	—

Data Pipeline

In pipeline, u can define what to do when a dataflow fails, and link to notebook & store procedures and have retry of scripts and many settings (which any one individual item can't do)

The screenshot shows the Azure Data Factory interface for the workspace 'DP_Embroiding'. The left sidebar includes options like Home, Activities, Run, View, Create, Browse, OneLake data hub, Monitoring hub, and Workspaces. The main area shows the 'Activities' tab selected, with a 'Stored procedure' activity highlighted. The configuration pane at the bottom shows the following details:

General tab selected.

Name *: Stored procedure1

Fabric Analytics Concepts Notes

The screenshot shows two main sections of the Fabric Analytics interface.

Top Section (Stored Procedure Parameters):

- Data store type:** Workspace (selected)
- Warehouse:** dp600-dw
- Stored procedure name ***: [dbo].[sp_Employee_GetByFirstName]
- Stored procedure parameters** (dropdown menu):
 - [dbo].[StoredProcedureName]
 - [dbo].[get_employees]
 - [dbo].[spFirstNameSearch]
 - [dbo].[spGetAllEmployees]
 - [dbo].[sp_Employee_GetByFirstName] (highlighted with a red circle)
 - [dbo].[sp_Employees_GetAll]
- Treat as null** checkbox

Bottom Section (Notebook Settings):

- General** tab selected
- Workspace ***: ETL-DP600
- Notebook ***: NB_Load_To_Silver
- Base parameters** (dropdown menu):
 - [+ New]

If spark cluster in the notebook code is busy, below retry option retries the code 2 times after retry interval of 30sec

Fabric Analytics Concepts Notes

General Settings

Name * Learn more 

Description

Activity state  Activated  Deactivated

Timeout 

Retry  

Advanced

Retry interval (sec)  

 Dataflow1
   


General **Settings**

Workspace *  Refresh

Dataflow *  Refresh  Open  New

Notification Option *  No notification Mail on completion Mail on failure



Fabric Analytics Concepts Notes

The maximum number of scheduled refreshes allowed per day with a Dataflow Gen2 is...

- A) 12 refreshes per day
- B) 24 refreshes per day
- C) 48 refreshes per day**
- D) 96 refreshes per day
- E) Unlimited refreshes per day

Which of the following can you use to update a row in a Data Warehouse table?

- A) Stored Procedure**
- B) View
- C) Create Table
- D) Function

Which of the following statements is **FALSE**, when talking about operations in a Fabric Data Warehouse:

- A) You can call a function from a stored procedure
- B) You can call a function from a view
- C) You can query a view and another view
- D) You can call a stored procedure from a function**

You are orchestrating many Spark notebooks to perform data transformation activities at the same time.

Sometimes the notebook execution is failing because your cluster is busy.

What modification can you make to the data pipeline notebook activity to give the pipeline more chances to run successfully?

- A) Deactivate and reactivate the activity
- B) Set the number of retries to 2 or 3**
- C) Change the parameters you are passing into the notebook
- D) Add a Fail activity to handle the execution failure

Ch7. Transform Data

Data cleansing:

- Implement a data cleansing process
- Identify and resolve duplicate data, missing data, or null values
- Convert data types by using Dataflows or PySpark
- Filter data

Data enrichment

- Merge or join data
- Enrich data by adding new columns or tables

Data modelling

- Implement a star schema for a lakehouse or warehouse, including Type 1 and Type 2 slowly changing dimensions
- Implement bridge tables for a lakehouse or a warehouse
- Denormalize data
- Aggregate or de-aggregate data

The data cleansing process

Bronze

Silver

Gold

- Duplicate data
- Missing data
- Null values
- Converting data types
- Filter data

Many tools:
• T-SQL
• Spark
• Dataflow

Fabric Analytics Concepts Notes

Load csv files to lakehouse bronze layer

Name	Date modified	Type
branches.csv	5/3/2024 12:28:31 PM	CSV
countries.csv	5/2/2024 11:37:03 PM	CSV
date.csv	5/2/2024 11:37:03 PM	CSV
dealers.csv	5/2/2024 11:37:03 PM	CSV
products.csv	5/2/2024 11:43:57 PM	CSV
revenue.csv	5/2/2024 11:37:03 PM	CSV

Transformations in PowerQuery

Power Query Home Transform Add column View Help

Queries [1]

Dealer_ID	Model_ID	Branch_ID	Date_ID	Revenue
DLR0001	BMW-M1	BR0001	DT00001	13363978
DLR0002	BMW-M2	BR0011	DT00002	19979446
DLR0026	For-M26	BR0251	DT00030	6866014
DLR0030	Lin-M30	BR0291	DT00030	16295156
DLR0034	Cad-M34	BR0331	DT00034	19425790
DLR0036	Cad-M36	BR0351	DT00036	9737520
DLR0038	Cad-M38	BR0371	DT00038	7654362
DLR0039	Cad-M39	BR0381	DT00039	4722916
DLR0042	Che-M42	BR0411	DT00042	18106066
DLR0045	Che-M45	BR0441	DT00045	2907854
DLR0046	Che-M46	BR0451	DT00046	9805538
DLR0049	Che-M49	BR0481	DT00049	16954836
DLR0050	Che-M50	BR0491	DT00050	13123946
DLR0053	GMC-M53	BR0521	DT00053	12806852
DLR0061	Acu-M61	BR0601	DT00061	827032
DLR0062	Acu-M62	BR0611	DT00062	887358
DLR0064	Hon-M64	BR0631	DT00064	6012624
DLR0073	Hyu-M73	BR0721	DT00073	10618472
DLR0080	Nis-M80	BR0791	DT00080	1831088

Query settings

Properties

Name: revenue

Entity type: Custom

Applied steps

- Source
- Navigation 1
- Navigation 2
- Navigation 3**

Data destination

No data destination

Fabric Analytics Concepts Notes

Transformations in SQL Endpoint

```
-- identifying/ removing nulls / filtering with WHERE
```

```
with removed_nulls as (
    SELECT [Dealer_ID]
        ,[Model_ID]
        ,[Branch_ID]
        ,[Date_ID]
        ,[Units_Sold]
        ,[Revenue]
    FROM [LH_BRONZE].[dbo].[revenue]
    WHERE Revenue IS NOT NULL
)
select count(*) from removed_nulls
select count(*) from [LH_BRONZE].[dbo].[revenue]
```

```
-- type conversion / adding new columns
```

```
SELECT [Dealer_ID]
    ,[Revenue]
    ,CAST([Revenue] as FLOAT) as RevFloat
    ,[Revenue] / 2 as HalfRevenue
FROM [LH_BRONZE].[dbo].[revenue]
```

Fabric Analytics Concepts Notes

Transformations in PySpark

The screenshot shows a PySpark notebook interface. On the left, there's a sidebar with 'Lakehouses' expanded, showing 'LH_BRONZE' with 'Tables' like 'branches', 'countries', 'date', 'dealers', 'products', and 'revenue'. A cursor is hovering over the 'revenue' table. The main area has a code cell with the following content:

```
1 # Welcome to your new notebook
2 # Type here in the cell editor to add code!
3
4
5 # With Spark SQL, Please run the query onto the lakehouse which is from the same workspace as
6
7 df = spark.sql("SELECT * FROM LH_BRONZE.revenue")
8 display(df)
```

A note at the top says: "Other people in your organization may have access to notebooks and Spark job definitions in this workspace. Carefully review this item before running it." Below the code cell, a log message indicates: "✓ 2 sec - Command executed in 1 sec 607 ms by william on 4:25:27 PM, 5/06/24". Underneath the code cell, there are tabs for 'Table', 'Chart', and 'Showing rows 1 - 1000'. A 'Log' tab is also present. At the bottom, there's an 'Inspect' button. To the right of the code cell, a table is displayed with the following data:

	ABC Dealer_ID	ABC Model_ID	ABC Branch_ID	ABC Date_ID	I23 Units_Sold	I23 Revenue
33	DLR0150	Mar-M150	BR1491	DT00150	2	13805586
34	DLR0152	Mar-M152	BR1511	DT00152	2	6576128
35	DLR0153	Hyu-M153	BR1521	DT00153	2	11687542
36	DLR0155	Hyu-M155	BR1541	DT00155	2	12200662
37	DLR0157	Hyu-M157	BR1561	DT00157	2	1397660
38	DLR0163	Hyu-M163	BR1621	DT00163	2	NULL

The screenshot shows a PySpark notebook. A cursor is hovering over the variable 'df' in the following code cell:

```
1 # identifying duplicates
2 df \
3     .groupby(['Branch_ID','Date_ID']) \
4     .count() \
5     .where('count > 1') \
6     .show()
```

Below the code cell, a note indicates: "- Command executed in 1 sec 494 ms by william on 1:11:09 PM, 5/03/24".

dropDuplicates()

The screenshot shows a PySpark notebook. A cursor is hovering over the variable 'deduped' in the following code cell:

```
1 # identifying duplicate data
2 df.count()
3 deduped = df.dropDuplicates()
4 deduped.count()
5
6 print(f"Proces removed {df.count() - deduped.count()} rows from the dataset" )
```

Below the code cell, a note indicates: "- Command executed in 2 sec 432 ms by william on 1:14:20 PM, 5/03/24".

Missing data/ nulls values

Identifying missing values in a column

```
1 # option 1: using isNull()
2 nulls = df.filter(df.Revenue.isNull())
3 display(nulls)
4
5 # option 2, using .where and col
6 from pyspark.sql.functions import col
7 nulls2 = df.where(col("Revenue").isNull())
8 display(nulls2)
9
```

Dropping nulls values using dropna()

[See here for documentation](#)

Params explained:

- **how** - can be 'any' or 'all'. If 'any', drop a row if it contains any nulls. If 'all', drop a row only if all its values are null.
- **thresh** - drop rows that have less than thresh non-null values. This overwrites the how parameter.

```
1 no_nas = df.dropna(subset=['Revenue'])
2 print(f"Proces removed {df.count() - no_nas.count()} rows from the dataset" )
```

✓ 1 sec - Command executed in 796 ms by william on 4:32:36 PM, 5/06/24

> Spark jobs (5 of 5 succeeded) Resources Log

· Proces removed 6 rows from the dataset

Fabric Analytics Concepts Notes

Type conversion + adding new columns

```
1 # take a look at the schema (including data types)
2 df.printSchema()
3
4 # using cast, create a new column
5 type_conv = df.withColumn('UnitsSoldConverted', df.Units_Sold.cast("string"))
6 type_conv.printSchema()
7 display(type_conv)
8
9
```

✓ 1 sec - Command executed in 849 ms by william on 4:34:35 PM, 5/06/24

> Spark jobs (2 of 2 succeeded) Resources Log

```
root
|-- Dealer_ID: string (nullable = true)
|-- Model_ID: string (nullable = true)
|-- Branch_ID: string (nullable = true)
|-- Date_ID: string (nullable = true)
|-- Units_Sold: integer (nullable = true)
|-- Revenue: integer (nullable = true)

root
|-- Dealer_ID: string (nullable = true)
|-- Model_ID: string (nullable = true)
|-- Branch_ID: string (nullable = true)
|-- Date_ID: string (nullable = true)
|__ Units_Sold: integer (nullable = true)
```

Filter data

```
1 from pyspark.sql.functions import col
2
3 # df.filter()
4 filtered_df = df.where(col("Revenue") > 10000000)
5 print(f"Proces removed {df.count() - filtered_df.count()} rows from the dataset" )
```

- Command executed in 1 sec 457 ms by william on 1:56:35 PM, 5/03/24

Part 2) Data enrichment

Adding new columns (withColumn or withColumns)



```
1 # using withColumn to create new columns.  
2 enriched_df = df.withColumn('halfRevenue', df.Revenue/2)  
3 display(enriched_df)
```

- Command executed in 1 sec 502 ms by william on 3:01:10 PM, 5/03/24

Joining and merging

```
1 dealers_df = spark.sql("SELECT * FROM LH_BRONZE.dealers")  
2 display(dealers_df)|  
3  
4 countries_df = spark.sql("SELECT * FROM LH_BRONZE.countries")  
5 display(countries_df)  
6  
7
```

- Command executed in 4 sec 949 ms by william on 11:18:06 AM, 5/06/24

```
1 joined_df = (  
2     dealers_df  
3         .join(countries_df, dealers_df.Country_ID == countries_df.Country_ID)  
4         .select(dealers_df.Dealer_ID, dealers_df.Country_ID, countries_df.Country_Name)  
5     )  
6 display(joined_df)
```

Fabric Analytics Concepts Notes

Data Modeling

Data modelling

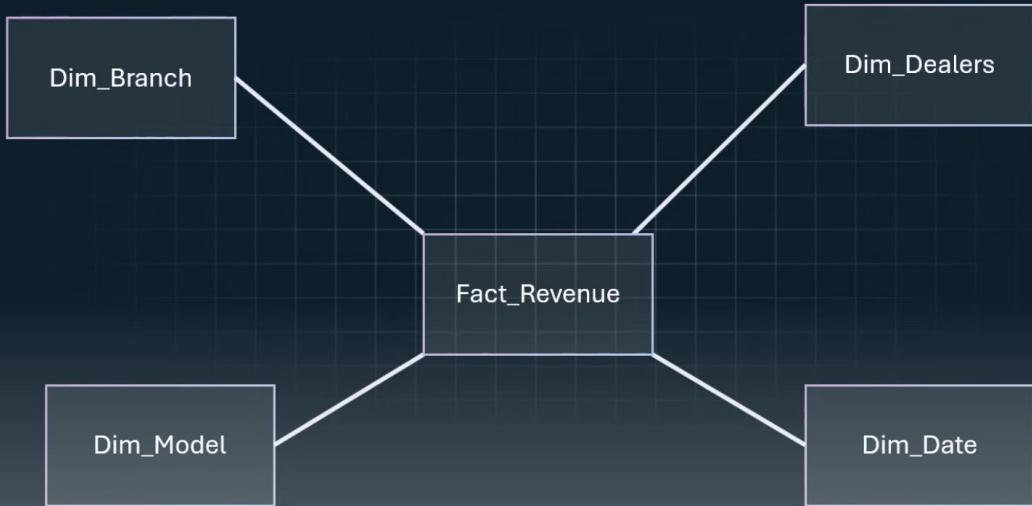
Bronze

Silver

Gold

- Star schema
- Type 1 and Type 2 SCDs
- Bridge tables
- Denormalize data
- Aggregate or de-aggregate data

Star schema



Fabric Analytics Concepts Notes

SCD

Type 1 SCD – No History (Overwrite)

From this:

EmployeeId	EmployeeName	Department
1	Sarah Jones	Consulting
2	Danny Walker	Marketing

To this:

EmployeeId	EmployeeName	Department
1	Sarah Jones	Consulting
2	Danny Walker	Sales

Danny
changes
departments

Implement with:

- OVERWRITE writing mode from either a data pipeline, dataflow or using pySpark

Note: if you're using a Lakehouse as your data store, the 'history' can be retrieved at a point in time using the delta log.

Type 2 SCD – ValidFrom, ValidTo and IsCurrent

From this:

EmployeeId	EmployeeName	Department	ValidFrom	ValidTo	IsCurrent
1	Sarah Jones	Consulting	2019-01-19	9999-01-01	1
2	Danny Walker	Marketing	2022-06-09	9999-01-01	1

To this:

EmployeeId	EmployeeName	Department	ValidFrom	ValidTo	IsCurrent
1	Sarah Jones	Consulting	2019-01-19	9999-01-01	1
2	Danny Walker	Marketing	2022-06-09	2023-11-02	0
2	Danny Walker	Sales	2023-11-03	9999-01-01	1

Implementing Type 2 SCD

- If using Lakehouse/ Spark
 - Bear in mind that delta logs store a history of all writes, so this might be a simpler option, depending on how you need to use the history
 - MERGE INTO – Delta update
- If using Data Warehouse T-SQL
 - Load to staging, then Stored Procedure to perform inserts, updates and ValidFrom, ValidTo calculations
 - MERGE operation in Fabric Data Warehouse is currently not supported (so you'll have to implement this manually)
 - Can implement row hashing to check for changes to rows

SCD Type 2 vs. Type 3: Key Differences

Feature	Type 2 SCD	Type 3 SCD
History Tracking	Full history (new row for each change).	Limited history (only current + previous value).
Storage Impact	High (adds rows over time).	Low (adds columns, fixed row count).
Implementation	Uses <code>ValidFrom</code> / <code>ValidTo</code> dates + <code>IsCurrent</code> flag.	Adds <code>PreviousValue</code> columns (e.g., <code>PreviousDept</code>).
Query Complexity	Requires time-based joins for historical data.	Simpler (current + previous value in one row).
Best For	Auditing, compliance, full historical analysis.	Short-term trend analysis (e.g., "last change").

Fabric Analytics Concepts Notes

Example: Employee Department Changes

SCD Type 2 (Full History)

EmployeeID	Name	Department	ValidFrom	ValidTo	IsCurrent
101	Alice	Sales	2023-01-01	2023-05-31	False
101	Alice	Marketing	2023-06-01	NULL	True

SCD Type 3 (Limited History)

EmployeeID	Name	CurrentDept	PreviousDept	LastChangeDate
101	Alice	Marketing	Sales	2023-06-01

Bridge tables (problem)

A company has many projects running, and employees are assigned to one or many projects. What does this look like as a data model?

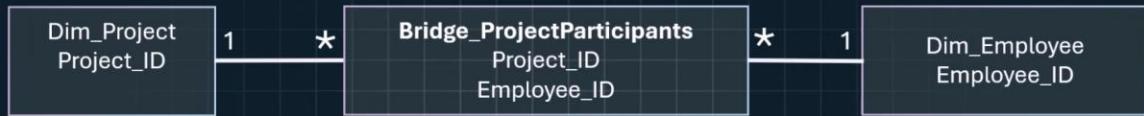


A many-to-many relationship exists between Dim_Projects and Dim_Employee

Fabric Analytics Concepts Notes

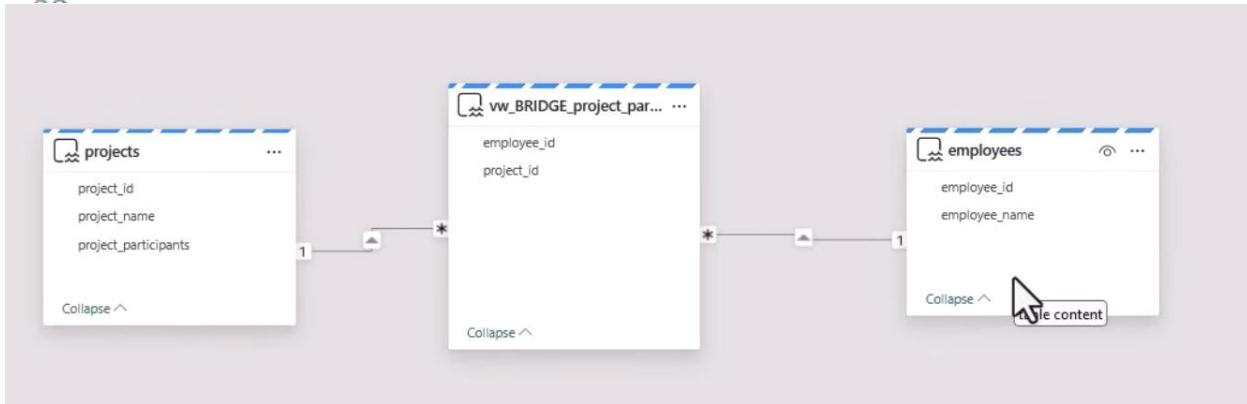
Bridge tables (solution)

A company has many projects running, and employees are assigned to one or many projects.
What does this look like as a data model?



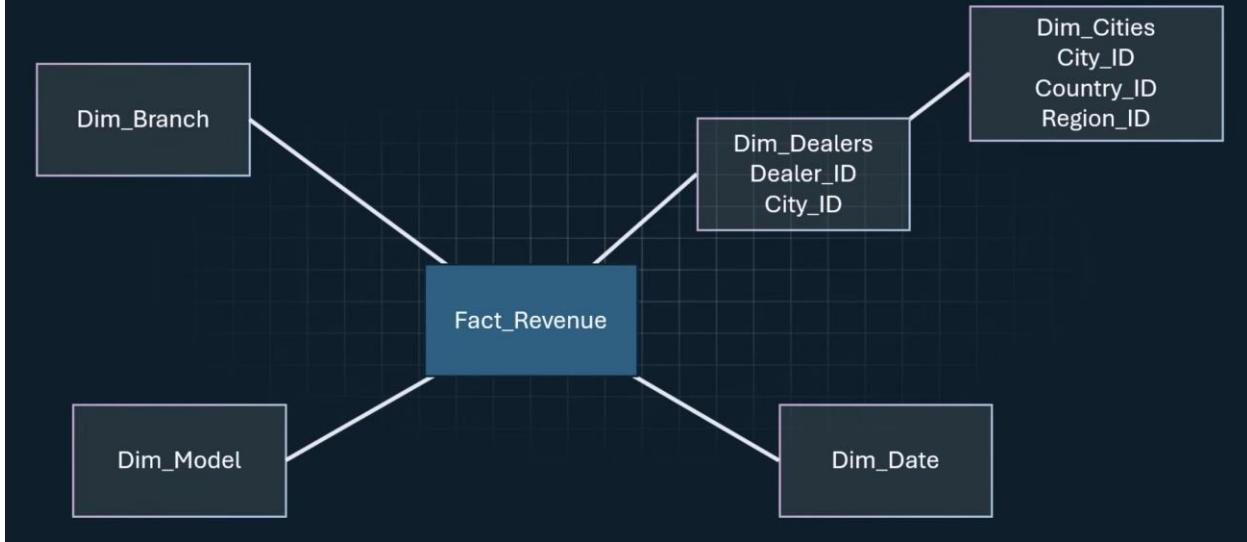
The bridge table 'resolves' the many-to-many relationship – 1 to 1 mapping of all Project / Project Member combinations

```
21
22
23  create view dbo.vw_BRIDGE_project_participants as
24    SELECT project_id, CAST(value as INT) as employee_id
25    FROM dbo.projects
26    |      CROSS APPLY STRING_SPLIT(project_participants, ',');
27
28
```



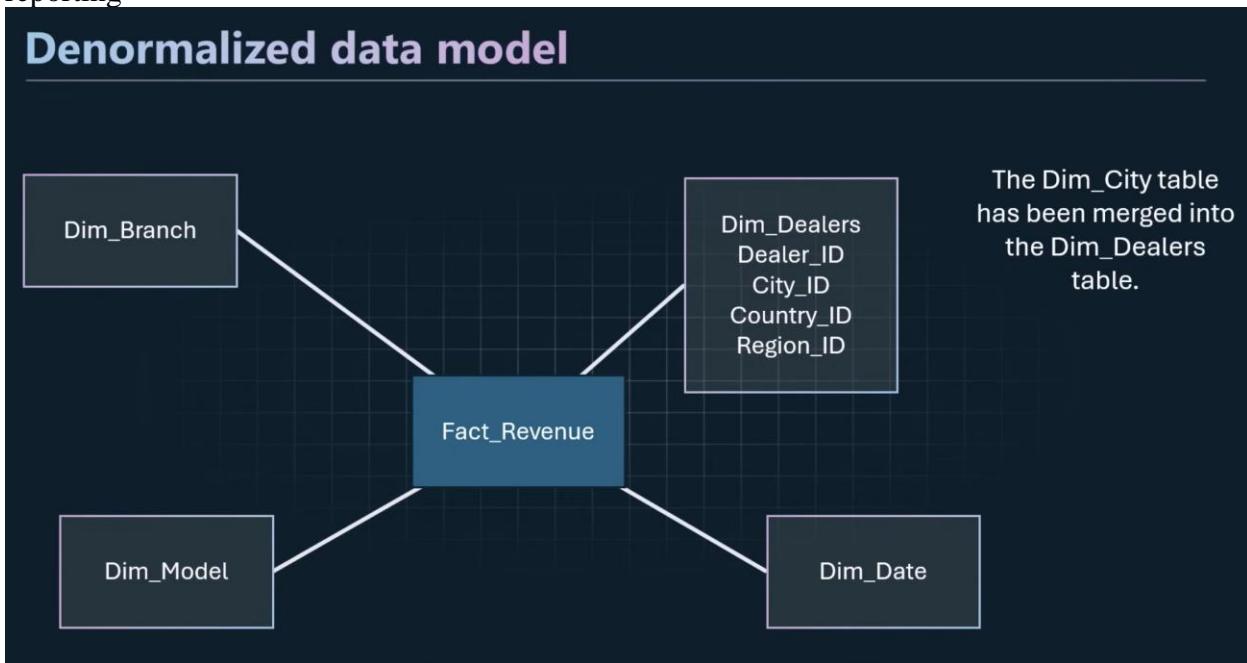
Fabric Analytics Concepts Notes

Normalized data model (Snowflake)

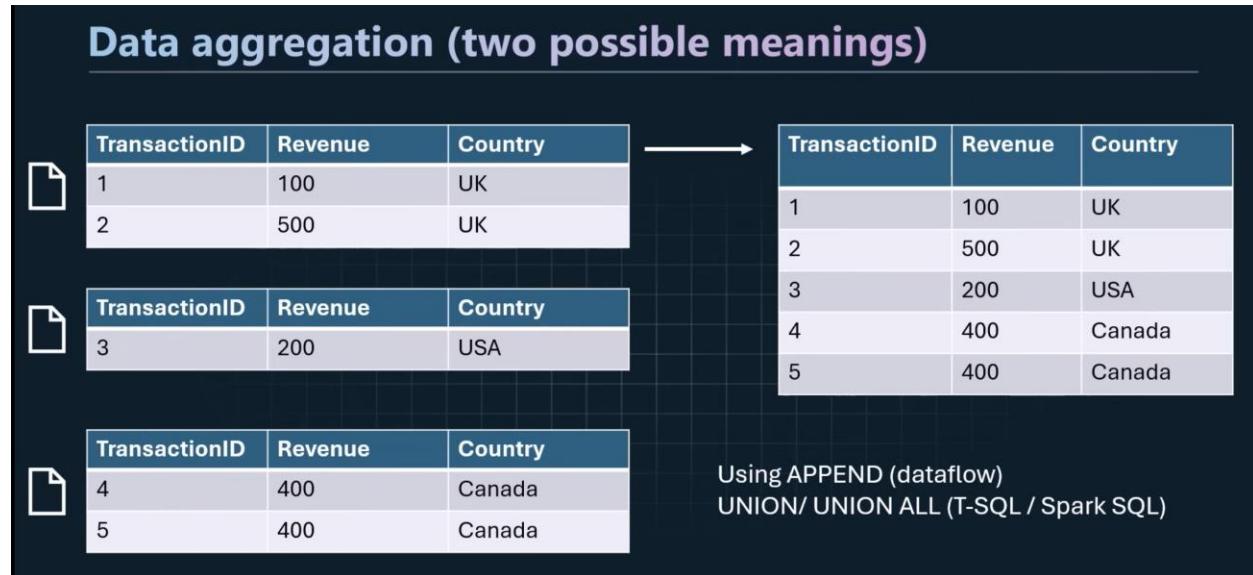


Denormalized from snowflask to star schema for efficiency (although redundant) in powerbi reporting

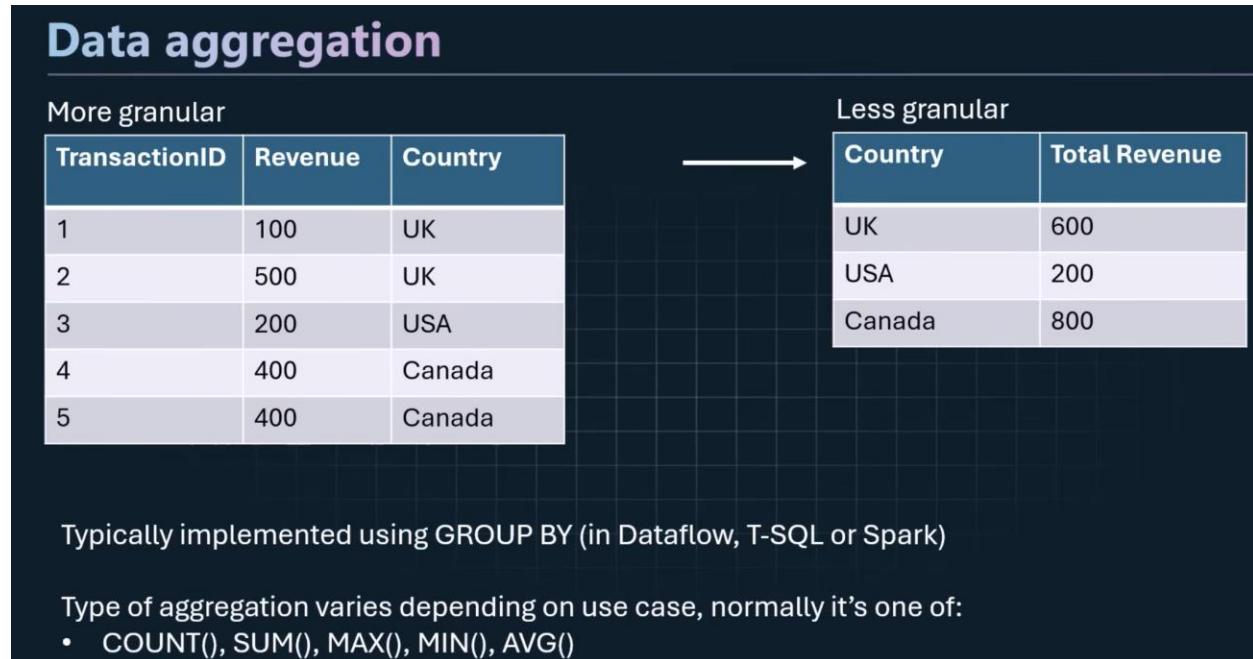
Denormalized data model



Fabric Analytics Concepts Notes



Union – appends all datasets (no chk for duplicates)
Union – removes any duplicate rows in result dataset



Fabric Analytics Concepts Notes

Question

When using dfjoin() in a pySpark notebook, the default join type

- A) FULL OUTER JOIN
- B) INNER JOIN
- C) LEFT JOIN
- D) RIGHT JOIN
- E) ANTI JOIN

JOIN Types

DataFrames

Left DF (Customers)	Right DF (Orders)
ID: A, Name: Alice	ID: B, Order: 01
ID: B, Name: Bob	ID: C, Order: 02
ID: C, Name: Carol	ID: D, Order: 03

LEFT JOIN Result

```
python
customers.join(orders, "ID", "left").show()
```

ID	Name	Order
A	Alice	NULL
B	Bob	01
C	Carol	02

Copy

▼

Fabric Analytics Concepts Notes

Left DF: [A, B, C] Right DF: [B, C, D]

INNER: [B, C]

LEFT: [A, B, C] + NULL for A

RIGHT: [B, C, D] + NULL for D

FULL: [A, B, C, D] + NULLs for A/D

ANTI: [A]

SEMI: [B, C] (no duplicates)

Question

You are looking to migrate a data transformation workload that currently is done using a Dataflow Gen2. You need to convert the same transformation into a T-SQL script.

A) LEFT JOIN

B) UNION ALL

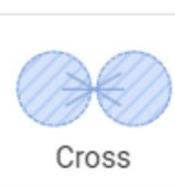
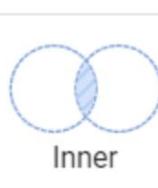
C) CONCAT

D) APPEND

E) UNION

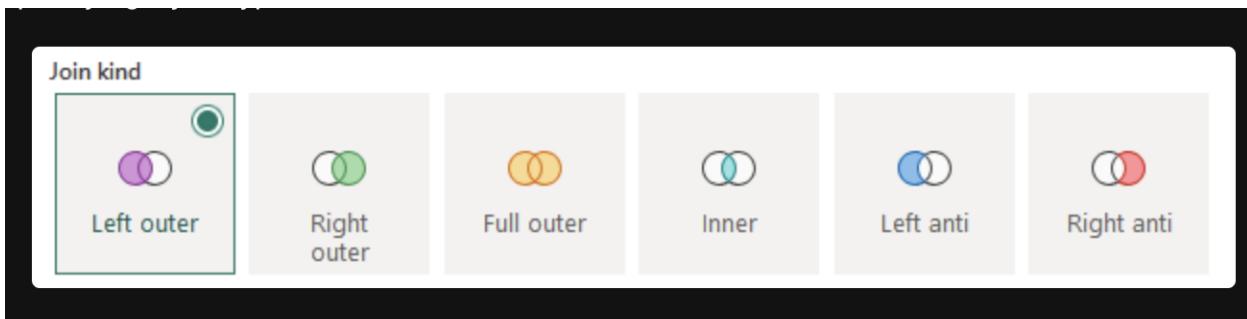
The Dataflow appends two datasets together and removes any duplicate rows.

Which T-SQL command can you use to implement this transformation?



Returns matching rows from the right table, plus non-matching rows from the left table

Fabric Analytics Concepts Notes



Question

You have a Spark Dataframe (df). Your goal is to remove rows that contain a null value in the TransactionDate column.

Which of the following will help you achieve this?

- A) `df.dropDuplicates()`
- B) `df.dropna(how='all')`
- C) `df.filter(df.TransactionDate.isNull())`
- D) `df.dropna(how='any')`
- E) `df.dropna(subset=['TransactionDate'])`

Why the Other Options Are Wrong

Option	Problem	What It Actually Does
A) <code>df.dropDuplicates()</code>	✖ Doesn't handle nulls.	Removes duplicate rows (ignores nulls).
B) <code>df.dropna(how='all')</code>	✖ Too broad.	Only drops rows where all columns are null (your goal is to target one column).
C) <code>df.filter(df.TransactionDate.isNull())</code>	✖ Opposite effect.	Keeps rows with nulls (instead of removing them). Use <code>isNotNull()</code> to fix.
D) <code>df.dropna(how='any')</code>	✖ Overkill.	Drops rows if any column has a null (not just <code>TransactionDate</code>).

Fabric Analytics Concepts Notes

Question

You inherit a data project and are inspecting the tables in the data warehouse.

One of the tables is a dimension table (Dim_Contacts), with the following columns:

- Contact_ID
- Contact_Name
- Contact_Address
- Effective_Date
- Effective_Until

A) Type 0 SCD

B) Type 1 SCD

C) Type 2 SCD

D) Type 3 SCD

Make an assumption about the type of data modelling that is being implemented in this dimension table

Ch8. Optimizing Performance

Framing performance optimization

	Identify performance issues	Possible methods to resolve issues
Dataflow	Refresh History Monitoring Hub Capacity Metrics App	Refactoring Staging Fast copy
SQL Data Warehouse	Query Insights DMVs Capacity Metrics App Statistics & Query Plan (in future)	Refactoring/ Query optimization
Notebook (Spark)	Spark History Server Monitoring Hub Capacity Metrics App	Refactoring/ Query optimization
Delta files	DESCRIBE	V-Order Optimization File partitioning VACUUM & OPTIMIZE

See logs

Fabric Analytics Concepts Notes

Power BI Monitoring hub

View and track the status of the activities across all the workspaces for which you have permissions within Microsoft Fabric.

Monitoring hub

To apply filters, select the values from the Filter dropdown menu.

Activity name	Status	Item type	Start time	Submitted by	Location
Data protection metrics (automatically generated)	Succeeded	Semantic model	4:30 AM, 5/10/24	william	My workspace
Fabric Capacity Metrics	Succeeded	Semantic model	12:00 AM, 5/10/24	william	Microsoft Fabric
Delta Optimization_c6bbebc2-8176...	Succeeded	Notebook	1:04 PM, 5/9/24	william	DP600-Transfo
Feature Usage and Adoption	Succeeded	Semantic model	12:40 PM, 5/9/24	Admin Monitoring	Admin monitor
Purview Hub	Succeeded	Semantic model	12:38 PM, 5/9/24	Admin Monitoring	Admin monitor
Delta Optimization_3ee8806e-1701-4e30-9055...	Succeeded	Notebook	12:13 PM, 5/9/24	william	DP600-Transfo
Notebook 8_4f6c4aa-f3c0b-464f-b2d3-e24db0a...	Succeeded	Notebook	10:50 AM, 5/7/24	william	My workspace
DP-600 Data Transformation Notebook _474803...	Succeeded	Notebook	4:23 PM, 5/6/24	william	DP600-Transfo
DP-600 Data Transformation Notebook _38bfd0...	Succeeded	Notebook	11:17 AM, 5/6/24	william	DP600-Transfo
LoadDimEmployeeToStar...	Success/Start	Dataflow Gen2	9:28 AM, 5/6/24	william	FTI - DP600

Monitoring hub

View and track the status of the activities across all the workspaces for which you have permissions within Microsoft Fabric.

Monitoring hub

To apply filters, select the values from the Filter dropdown menu.

Activity name	Status	Item type	Start time	Submitted by
Transformation Notebook _/8d/zb...	Succeeded	Notebook	2:41 PM, 5/3/24	william
Transformation Notebook _9d958c...	Succeeded	Notebook	1:08 PM, 5/3/24	william
Transformation Notebook _664f97...	Succeeded	Notebook	12:42 PM, 5/3/24	william
_TableLoad_159731f7-895f-4493-89...	Succeeded	Lakehouse	12:28 PM, 5/3/24	william
_TableLoad_75977776-d82f-4060-bf...	Succeeded	Lakehouse	11:37 PM, 5/2/24	william
to DW Tables	Failed	Dataflow Gen2	12:46 PM, 5/1/24	william
_ffbf39a8a-e8c5-46fe-82a2-f5f13eae...	Succeeded	Notebook	10:43 AM, 4/23/24	william
_2daa331e-c7ee-4ac2-a2cd-93d550...	Succeeded	Notebook	8:50 AM, 4/23/24	william
ouse_TableLoad_d949209d-ee5d-4b...	Succeeded	Lakehouse	4:24 PM, 4/22/24	william
OUSE_TableLoad_39e053b7-ea43-4f1...	Succeeded	Lakehouse	7:49 AM, 4/21/24	william

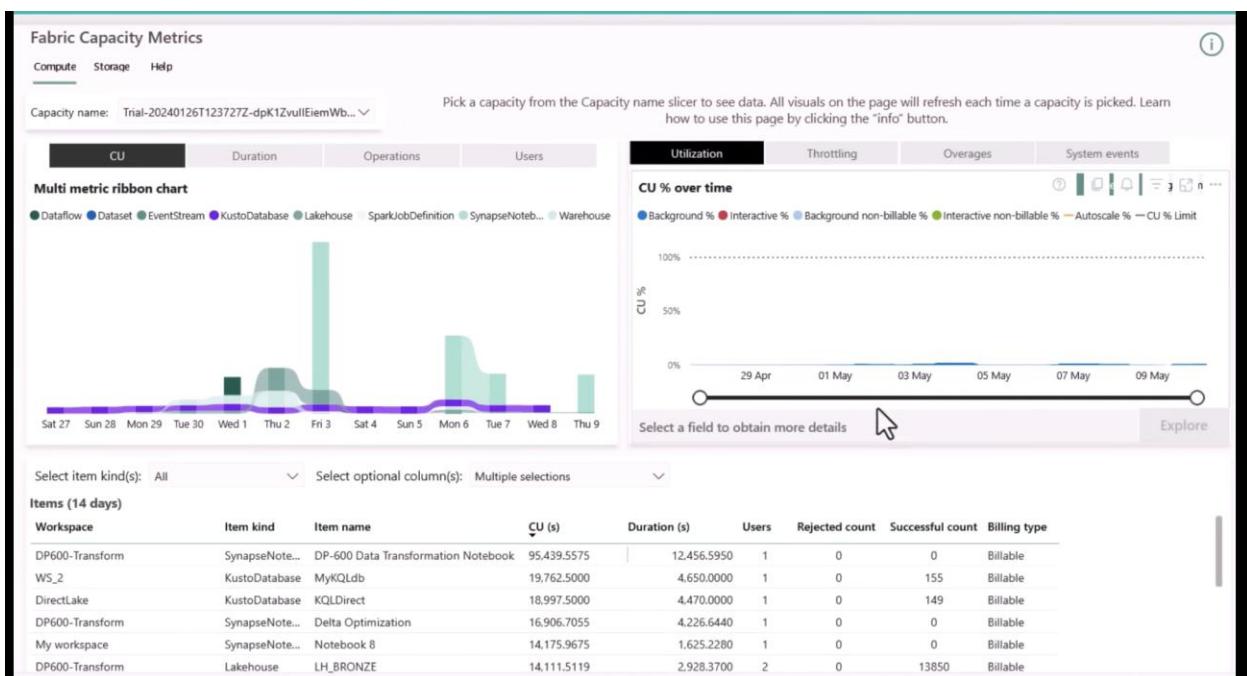
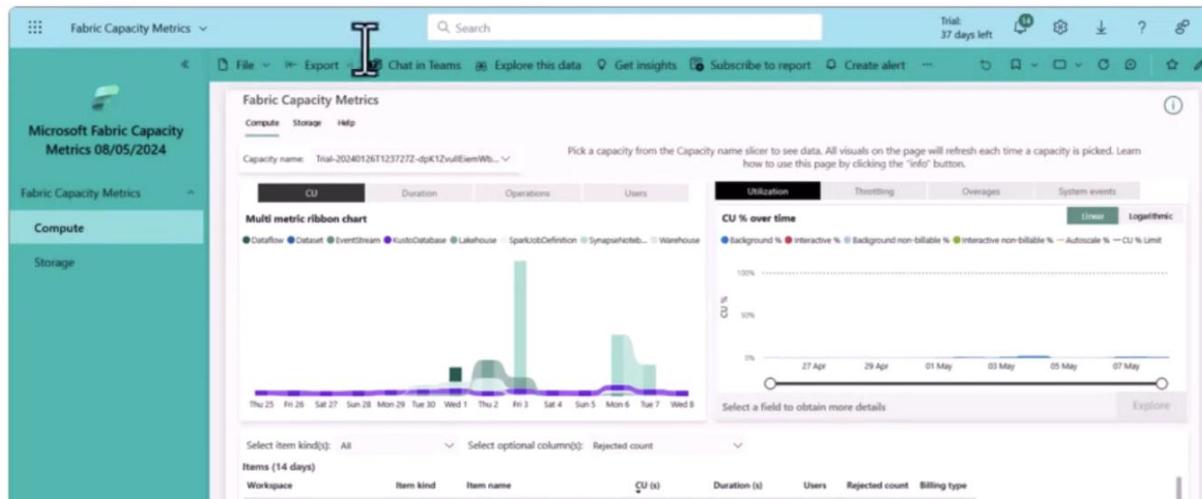
Item type: Dataflow Gen2
Status: Failed
Start time: 01/05/2024, 11:46:54
End time: 01/05/2024, 11:50:40
Duration: 3m 45s
Submitted by: william
Location: ETL-DP600
Capacity: Trial-20240126T123727Z-dpK1ZvullEiemWbdWzv95w
Average duration: 0s
Refreshes per day: 0

Fabric Analytics Concepts Notes

Capacity Metrics App

The Capacity Metrics App ([install instructions](#)), is a Power BI App that you can download (for free) and install onto your Fabric tenant.

You give it your Capacity ID (you can find this in the [Capacity Settings](#) part of the Admin Portal), and it retrieves capacity usage statistics like so:



Dataflows + performance

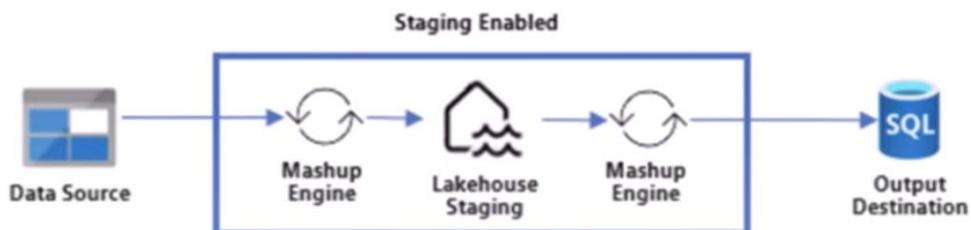
Performance monitoring (Dataflows)

- High-level summary of a dataflow run is available in the **Monitoring Hub**.
- Lower-level data and understanding of what is happening in a particular dataflow is available in the Refresh History of a particular dataflow.
 - You can inspect error messages,
 - Get a breakdown of the different sub-activities being performed by the Dataflow.

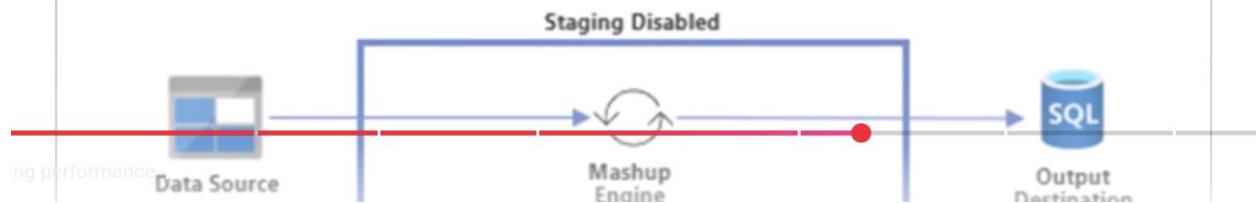
Performance optimization/ features:

Staging - what is staging? By default, transformations in a dataflow are carried out by the PowerQuery engine (also known as the Mashup engine), this is relatively slow, compared to the Spark engine. So staging involves using the Mashup engine to first get the data from source, then it writes the data out to a Staging Lakehouse, performs the transformation using the Spark engine (rather than the Mashup engine), then the Mashup engine reads the result and then writes to data to destination.

Behind the scenes, controlling the staging setting of a query will change orchestration from its default setting:



To orchestration that extracts, transforms, and loads data from source to destination in memory:



Fabric Analytics Concepts Notes

The added overhead of writing the data out to a Lakehouse (and then back) adds considerable overhead, so Staging is only ever more efficient when you've got large datasets, or lots of transformation steps (or both).

Fast Copy (preview) - a new feature to speed up getting new data from sources (but not on-premise). I believe it uses the same technology as the Data Pipeline CopyData activity. Still a preview feature, and relatively recent, so might not be on the exam yet.

More resources:

- [Spotlight blog on Dataflow \(including performance optimization\)](#)
- [Best practices for complex dataflows](#)

SQL Data Warehouse + Performance

So you've got a really long-running SQL query in Data Warehouse. What to do about it??
There's a few things we can do:

1) Monitor using the Capacity Metrics App ([see here](#))

The Capacity Metrics App can be used to see how much CU(s) has been used for different read/write queries you've executed in the Data Warehouse.

2) Dynamic Management Views (DMVs)

For the current version, there are three dynamic management views (DMVs) provided for you to receive live SQL query lifecycle insights.

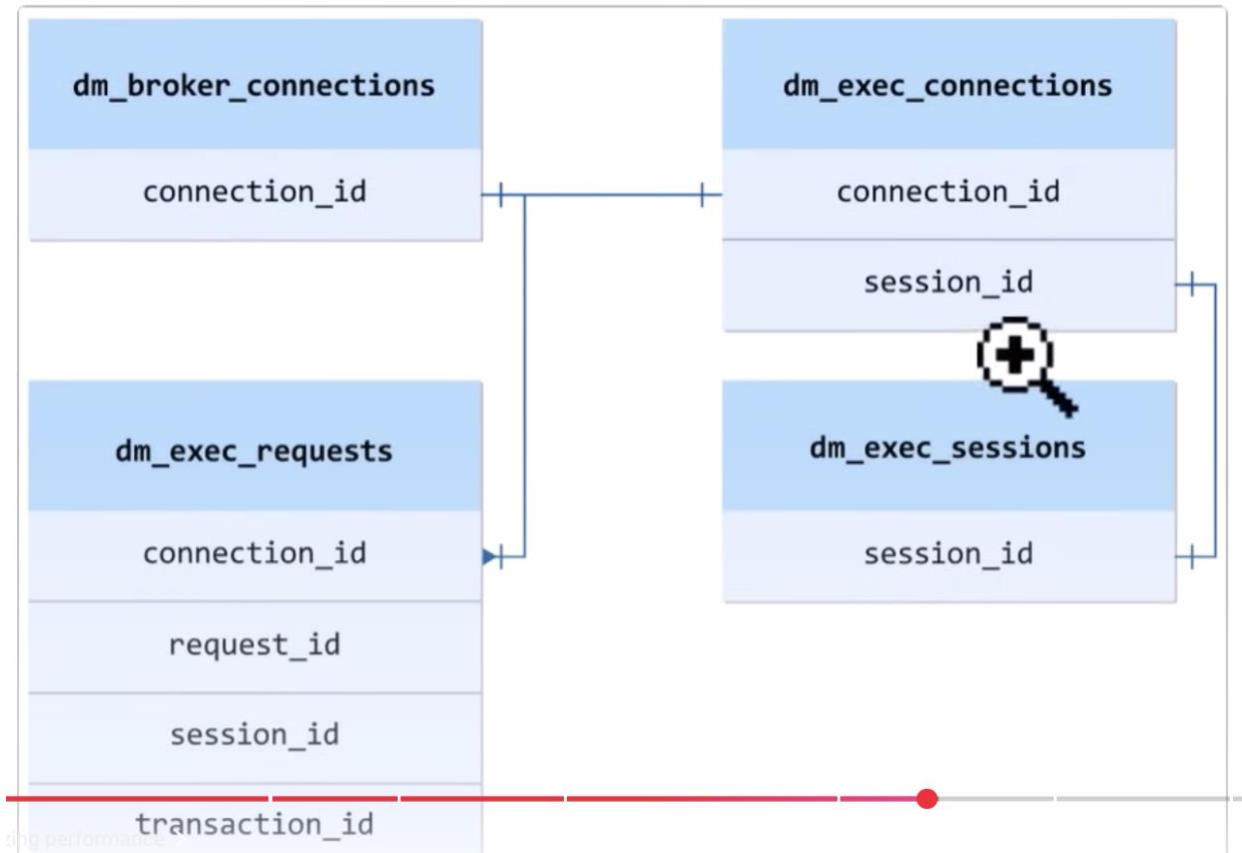
- [sys.dm_exec_connections](#)
 - Returns information about each connection established between the warehouse and the engine.
- [sys.dm_exec_sessions](#)
 - Returns information about each session authenticated between the item and engine.

Fabric Analytics Concepts Notes

- [sys.dm_exec_requests](#)

- Returns information about each active request in a session.

How these queries are related:



These three DMVs provide detailed insight on the following scenarios:

- Who is the user running the session?
- When was the session started by the user?
- What's the ID of the connection to the data Warehouse and the session that is running the request?
- How many queries are actively running?
- Which queries are long running?

DMV are low level logs

Fabric Analytics Concepts Notes

Queryinsights is a schema exposed in every Data Warehouse, it consists of four main views, all of which are useful for performance monitoring:

- [exec_requests_history](#): Returns information about each completed SQL request/query.
- [frequently_run_queries](#): Returns information about frequently run queries
- [long_running_queries](#): Returns the information about queries by query execution time.

Note: it's not currently possible to view the Execution/ Query Plan for SQL query in the Fabric Data Warehouse - although there are plans to support/ expose this in the future.

Notebook spark jobs detailed info with spark history server

DP600-Transform > Delta Optimization > Delta Optimization_c6bbebc3-8176-4a48-9ee7-0989c4ccf992								Run details			
Jobs	Resources (Preview)	Logs	Data	Item snapshots	Status	Stages	Tasks	Duration	Processed	Data	Status
> Job 10	getRowsInJsonString at Display.scala:452				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	353 ms	4,096 rows	1.62	Stopped
> Job 9	parquet at NativeMethodAccessorImpl.java:0				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	373 ms	0 rows	0 B	Application application_
> Job 8	takeAsList at <console>:44				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	534 ms	3 rows	1.13	Total duration 20 min 46 sec
> Job 7	\$anonfun\$recordDeltaOperationInternal\$1 at SynapseLoggingShim.scala:111				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	287 ms	3 rows	639 B	Queued duration 0 sec
> Job 6	\$anonfun\$recordDeltaOperationInternal\$1 at SynapseLoggingShim.scala:111				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	164 ms	1 row	353 B	Running duration 20 min 46 sec
> Job 5	\$anonfun\$recordDeltaOperationInternal\$1 at SynapseLoggingShim.scala:111				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	1 sec 392 ms	5 rows	1.2 K	Livy ID c6bbebc3-8
> Job 4	\$anonfun\$recordDeltaOperationInternal\$1 at SynapseLoggingShim.scala:111				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	461 ms	8 rows	1.27	Submitter william@fab
> Job 3	toString at String.java:2951				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	512 ms	4 rows	1.27	Submit time 5/9/24 1:04:11
> Job 2					✓ Succeeded	0/0	<div style="width: 100%;">0/0 succeeded</div>	< 1 ms	0 rows	0 B	
> Job 1	save at <console>:38				✓ Succeeded	1/1	<div style="width: 100%;">1/1 succeeded</div>	6 sec 84 ms	6 rows	292 B	
> Job 0	save at <console>:38				✓ Succeeded	1/1	<div style="width: 100%;">3/3 succeeded</div>	1 sec 848 ms	3 rows	0 B	

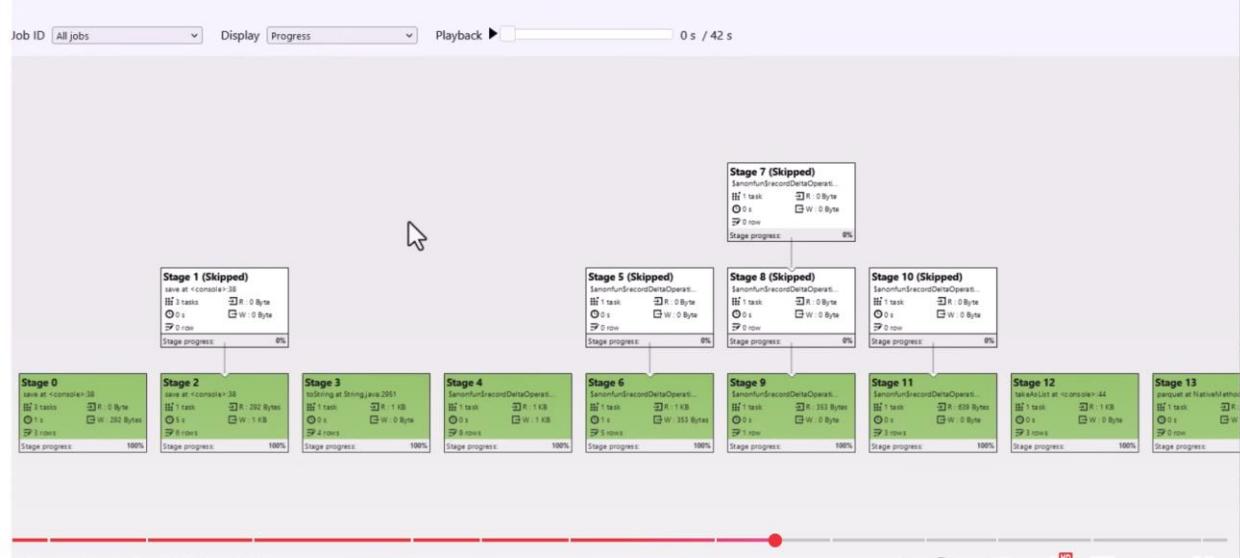
3:47:42 / 6:03:40 - Optimizing performance



5/9/24 1:04:11

Fabric Analytics Concepts Notes

Spark Application & Job Graph



Delta Table Optimization

Fabric is built on delta tables, so this is a pretty important topic. The delta file format is great, but it can lead to poor performance and bloated storage sizes, if not managed correctly.

This is a topic that can run very deep, but we'll just cover some of the basics for the exam.

Download and run the Spark notebook below

In this notebook, we will first look at file partitioning, then we will look at other methods for optimizing delta table (and Spark) performance.

This can improve:

- read/retrieval performance (no need to scan through millions of rows to find the record(s) you need)
- transformation performance (partitions can be transformed in parallel)

A client you're working with wants to reduce the SKU of their Fabric Capacity from an F16 to an F8 to save money.

They want to find the most resource-intensive workloads and optimize them to use less CU(s).

Where should they look to find this information?

A) Monitoring Hub

B) Capacity Metrics App

C) Query Insights

D) Spark History Server

E) OneLake Hub

Fabric Analytics Concepts Notes

Correct Answer: B) Capacity Metrics App

Why?

- The **Capacity Metrics App** provides detailed insights into **Capacity Unit (CU) consumption** across all workloads (Power BI, Data Factory, Spark, etc.).
- Key metrics:
 - **CU usage per item/workspace** (identify high-cost queries).
 - **Peak utilization times** (optimize scheduling).
 - **Resource-heavy operations** (e.g., long-running Spark jobs).

How to Access:

1. Go to your Fabric **admin portal**.
2. Navigate to **Capacity Settings** → **Metrics App**.

Why Not the Others?

Option	Why It's Not the Best Fit
A) Monitoring Hub	Shows activity logs (who did what) but not CU consumption .
C) Query Insights	Optimizes individual queries (Power BI), not overall capacity planning.
D) Spark History Server	Limited to Spark job logs (ignores other workloads like Dataflows).
E) OneLake Hub	Manages data storage , not compute resources.

Fabric Analytics Concepts Notes

Best practices

This article describes best practices when using Delta Lake.

Choose the right partition column

You can partition a Delta table by a column. The most commonly used partition column is `date`. Follow these two rules of thumb for deciding on what column to partition by:

- If the cardinality of a column will be very high, do not use that column for partitioning. For example, if you partition by a column `userId` and if there can be many distinct user IDs, then that is a bad partitioning strategy.
- Amount of data in each partition: You can partition by a column if you expect data in that partition to be at least 1 GB.

Compact files

If you continuously write data to a Delta table, it will over time accumulate a large number of files, especially if you add data in small batches. This can have a negative effect on the efficiency of table reads, and it can also affect the performance of your file system. Ideally, a large number of small files should be rewritten into a smaller number of larger files on a regular basis. This is known as compaction.

You can compact a table by repartitioning it to smaller number of files. In addition, you can specify the option `dataChange` to be `false` indicates that the operation is a compaction.

High cardinality refers to a column or field in a database that contains a large number of unique values (close to the total number of rows). Ex: Employee ID,

Write data to delta (three ways)

```
1 # write WITHOUT Partitions, first
2 transformed_df.write.mode("overwrite").format("delta").save("Tables/flights_not_partitioned")
3
4 # then write WITH partitions
5 transformed_df.write.mode("overwrite").format("delta").partitionBy("year", "month").save("Tables/flights_partitioned")
6
7 # then write WITH (more) partitions
8 transformed_df.write.mode("overwrite").format("delta").partitionBy("year", "month", "day").save("Tables/flights_partitioned_daily")
```

Fabric Analytics Concepts Notes

Inspecting our delta tables (and parquet files)

```
1 %%sql  
2 describe detail flights_partitioned  
3
```

* 4 sec - Running

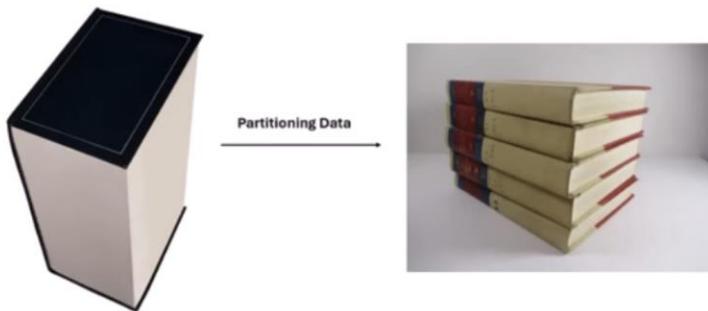
```
1 %%sql  
2 describe detail flights_not_partitioned
```

- Command executed in 1 sec 459 ms by william on 12:40:25 PM, 5/09/24

```
1 %%sql  
2 describe detail flights_partitioned_daily
```

If partitioned column has data that is less than 1gb, u come across small file pb. Partitioning is balancing act.

- As we know, Spark is a distributed processing engine, meaning it is capable of processing multiple sets of data 'in-parallel'.
- To make this possible, it helps to 'partition' the data. Rather than passing in one massive file to the Spark engine, multiple partitions allow the process/ transform the data in parallel, greatly improving performance.
- It is a balancing act though, because if you go too far, then you will suffer from the 'small file problem'.



Fabric Analytics Concepts Notes

1 / 2

V-Order optimization

What is V-Order?

V-Order (not to be confused with [Z-order](#)) is a Microsoft proprietary algorithm. It works by providing some special-sau~~T~~ sorting/ compaction. compression files to reduce file size, and optimize read performance across all engines within Fabric. The outputted Parquet file is fully parquet-compliant.

Checking if V-Order is enabled

In a notebook you can check if V-Order is enabled (it should be ENABLED by default):

```
1 # Check if v-order optimization is Enabled
2 spark.conf.get('spark.sql.parquet.vorder.enabled')
```

- Command executed in 277 ms by william on 12:35:05 PM, 5/09/24

PySp

Disabling/ re-enabling it:

```
1 # disable v-order
2 spark.conf.set('spark.sql.parquet.vorder.enabled', 'false')
3
4 # enable v-order
5 spark.conf.set('spark.sql.parquet.vorder.enabled', 'true')
```

Delta table maintenance & optimization techniques

The following section describes a few delta table maintenance and optimization techniques:

OPTIMIZE & VACUUM (Delta)

- OPTIMIZE is a delta lake method that performs bin-compaction (read more [here](#)). This means is can tidy up lots of small files. It's also idempotent.
- VACUUM ([read more here](#)) involves the removal of files no longer referenced by a Delta table.

Coalesce & Repartition (Spark)

- `coalesce()` ([read more here](#)) is a Spark method for reducing the amount of partitions in a delta table. Say you have 100 partitions, you can 'coalesce' them into 10 partitions. Importantly, this is quite an efficient operation, because it doesn't require a shuffle of your data.
- `repartition()` ([read more here](#)) involves the breaking of existing partitions to create new partitions, these can either be more or less than the original partitions. Repartitioning is an expensive operation because it involves shuffling (unlike `coalesce`)

Fabric Analytics Concepts Notes

Question

You notice one of your Dataflow Gen2 runs failed to refresh last night.

Where would you go to find out why a particular Dataflow might have failed a particular run?

- A) Capacity Metrics App
- B) Monitoring Hub
- C) Power Query Error Hub
- D) Dataflow Refresh History**
- E) Data Pipeline Run History

Question

When talking about delta table optimization, which of the following operations removes old files no longer referenced by a Delta table log?

- A) V-order optimization
- B) Z-ordering
- C) VACUUM**
- D) OPTIMIZE
- E) Bin compaction

Z-Order Optimization & Bin Compaction

These are advanced data organization techniques used in **distributed computing** and **big data systems** (like Apache Spark, Delta Lake, or Iceberg) to improve query performance, reduce I/O, and optimize storage.

1. Z-Order Optimization (Morton Order)

What It Is

Z-Order is a **multi-dimensional clustering** method that maps data with multiple columns into a single **Z-curve** (space-filling curve), keeping similar values close on disk.

Why Use It?

- Improves **query performance** for multi-column filters (e.g., `WHERE date = X AND region = Y`).
- Reduces **data skipping** (fewer files scanned).
- Works well in **columnar formats** (Parquet, ORC).



Feature	Z-Order Optimization	Bin Compaction
Purpose	Optimize query performance	Reduce small-file overhead
Method	Multi-column clustering	Merging small files
Impact	Faster reads (data skipping)	Faster metadata ops & scans
Used in	Delta Lake, Iceberg, databases	Data lakes, streaming systems

Fabric Analytics Concepts Notes

2. Bin Compaction (Small File Problem Fix)

What It Is

A technique to **merge small files** (e.g., in data lakes) into larger, optimized "bins" to:

- Reduce **metadata overhead** (e.g., too many Parquet footers).
- Improve **read parallelism** (larger files = better throughput).
- Avoid the "**small files problem**" (common in streaming/CDC systems).

How It Works

1. **Group small files** (e.g., <128MB) into logical "bins."
2. **Compact them** into larger files (e.g., 1GB).
3. **Update metadata** (e.g., Hive/Delta Lake transactions).

When to Use Both?

- First compact small files (Bin Compaction).
- Then apply Z-Order for query optimization.

Fabric Analytics Concepts Notes

V-Order (short for "Vertex Order") is a **write-time optimization** technique developed by **Databricks** to enhance the performance of **Delta Lake and Parquet files** in Apache Spark. It reorganizes data on disk to improve **compression, query speed, and cost efficiency** in cloud storage (e.g., S3, ADLS).

Key Features of V-Order

1. Automatic Sorting & Clustering

- Unlike **Z-Order** (which requires explicit `ZORDER BY`), **V-Order automatically sorts data** during writes.
- It groups similar values together, improving **dictionary encoding** and **columnar compression**.

2. Better Compression

- Achieves **smaller file sizes** (up to **2–4x better compression** vs. unsorted data).
- Reduces **storage costs** and **I/O operations** in cloud environments.

3. Faster Queries

- Enables **data skipping** (fewer files scanned).
- Works well with **Delta Lake's metadata** for optimized reads.

4. No Extra Maintenance

- Unlike **Z-Order**, which requires manual optimization (`OPTIMIZE ... ZORDER BY`), **V-Order works at write time** without additional commands.

Question

Which of the following statements about V-Order optimization is FALSE?

- A) V-Order optimization is enabled by default in the Fabric Spark Runtime.
- B) V-Order can be enabled during table creation using TBLPROPERTIES.
- C) A table can be both V-ordered and Z-ordered.
- D) V-Order improves the read performance for parquet files.
- E) **V-Order speeds up the write time of a Parquet file.**

Vorder works at write time (hence increase write time than decrease / speed up the write time). Vorder takes longer time to write files but it massively improves the read performance.

Fabric Analytics Concepts Notes

Question

You want to analyse long-running queries in a Fabric Data Warehouse.

What is the minimum workspace role you need to run the following query:

```
SELECT * FROM [queryinsights].[long_running_queries]
```

A) Admin

B) Member

C) Contributor

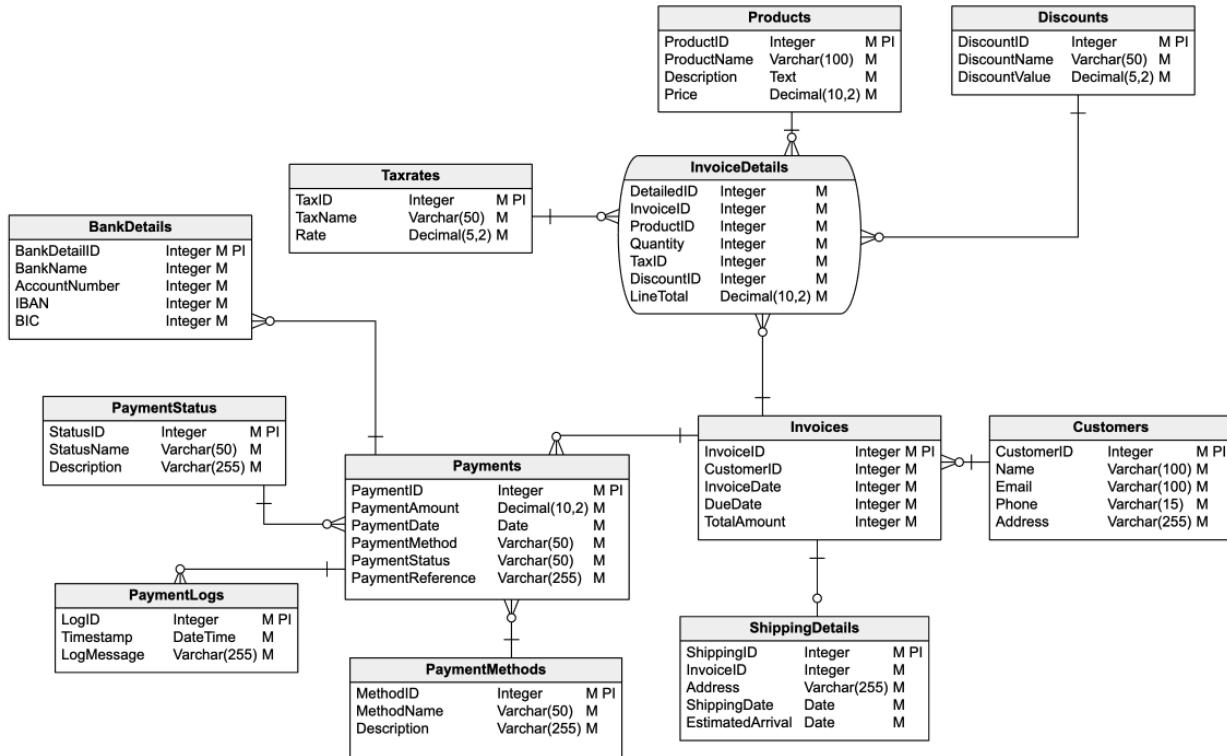
D) Viewer

Ch9. Design & build semantic models

Normalized vs de-normalized data

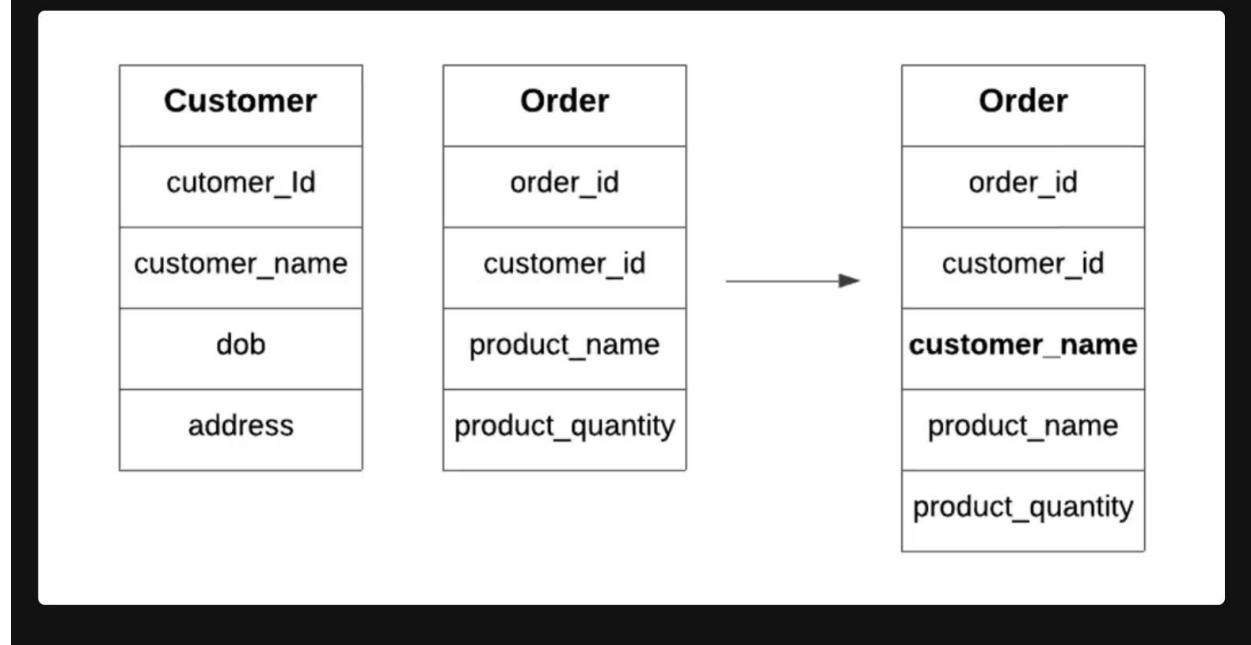
OLTP systems are built for to handle lots of transactions, very fast. The backend database tables are designed to be fully normalized. This means that entities in an OLTP application each have their own table in the backend. It might look something like this, which is an example of an entity-relationship diagram for a billing system. This structure minimizes redundancy - every value is only stored in one place, making it efficient, from a storage perspective. But this kind of heavily normalized structure is less optimal for analytics use cases. OLTP systems are commonly used as a data source for OLAP analytical systems, however, we must de-normalize how the data is structured to make it more useful for analytical workloads.

Fabric Analytics Concepts Notes



Denormalization of data

Put simply, denormalizing data looks like this (below), two or more tables are 'denormalized' into a single table for analytics use cases. This can be done in a Data Warehouse using JOINs.



Fabric Analytics Concepts Notes

KQL

Returning a preview sample of 5 records from `StormEvents` table, returning the State, EventType and DamageProperty columns:

```
StormEvents  
| take 5  
| project State, EventType, DamageProperty
```

Filtering using the WHERE operator:

```
StormEvents  
| where State == 'TEXAS' and EventType == 'Flood'  
| project StartTime, EndTime, State, EventType, DamageProperty
```

Simply aggregation using the summarise operator:

```
StormEvents  
| summarize TotalStorms = count() by State
```

Grouping data into bins:

```
StormEvents  
| where StartTime between (datetime(2007-01-01) .. datetime(2007-12-31))  
    and DamageCrops > 0  
| summarize EventCount = count() by bin(StartTime, 7d)
```

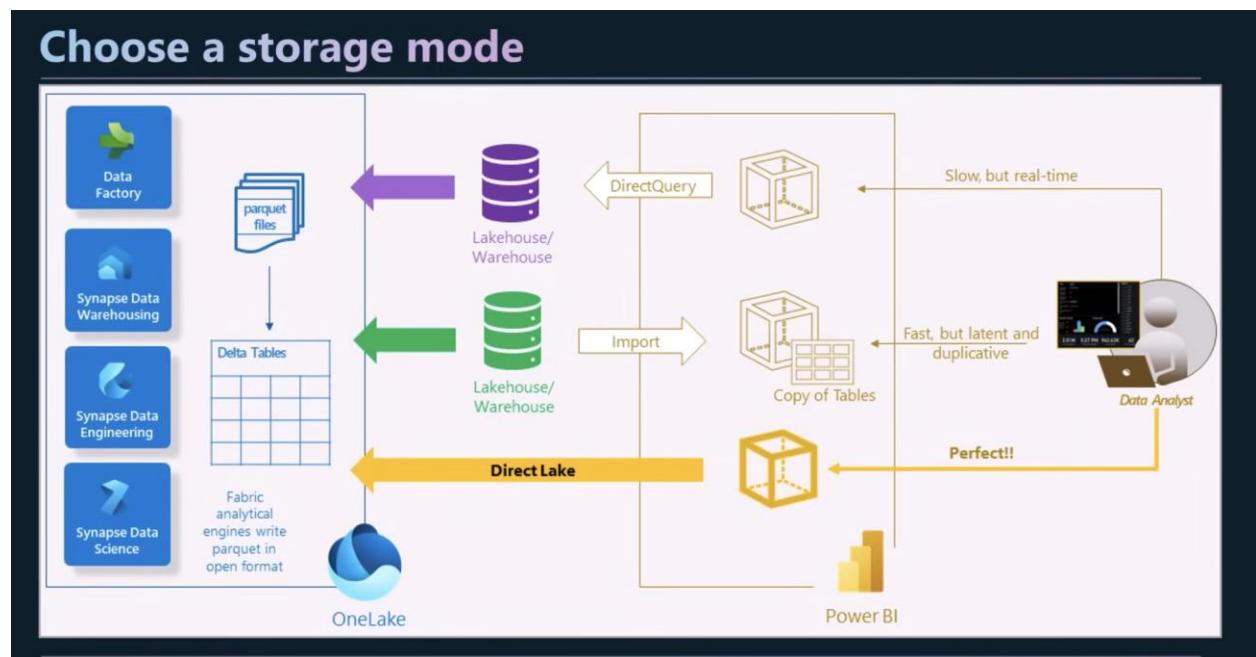
Calculating the max, min, avg and sum:

Fabric Analytics Concepts Notes

StormEvents

```
| where DamageCrops > 0  
| summarize  
    MaxCropDamage=max(DamageCrops),  
    MinCropDamage=min(DamageCrops),  
    AvgCropDamage=avg(DamageCrops)  
    by EventType  
| sort by AvgCropDamage
```

Storage modes



Choose a storage mode

Choose IMPORT MODE when...

- your data size is small enough to fit in memory
- you prioritise read performance
- you have no requirement for near real-time updates
- you want to use Calculated Columns/ Tables
- You want to combine data from multiple data sources

Choose DIRECT LAKE MODE when...

- you need near real-time updates
- your data is stored in one Fabric data store
- Your dataset size is large (tens or hundreds of GB)
- You have data modelling capability 'upstream', in a Lakehouse or Data Warehouse

Choose DIRECT QUERY mode when...

- you need near real-time updates
- You can do data transformation 'upstream'
- building a composite models

Composite models

- A composite model combines one more than connection mode.
- Commonly, this is a Direct Query Fact Table, and Import Mode Dimensions
- Composite models also provide a way to model many-to-many relationships, without the need for bridge tables.

Fact table has lots of data (million to billions of rows). Fact table is updated very often (hundreds of records every second in some oltp transactions). Direct mode is good for fact. As data is changing very fast in fact, direct mode gives near real time access to that changing / fresh data.

Dim tables might change a lot slower & they also have fewer rows than fact. So Dim table's can use import mode. Let's say products dim table is updated once per day.

Aggregations (Power BI feature)

- Aggregations are a specific feature within Power BI for aggregating large datasets and caching the aggregation.
- The goal is to improve performance in large data models.
- Aggregations can also be used in composite models.
- Aggregations can either be created in your data source, or you can use Import Mode and create the aggregation in the Power Query engine.
- User-defined aggregations – using the ‘Manage aggregations’ dialog in Power BI Desktop to define aggregations for aggregation columns with summarization, detail table, and detail column properties
- Automatic aggregations (Premium subscriptions only) - use state-of-the-art machine learning (ML) to continuously optimize DirectQuery semantic models for maximum report query performance

Large format semantic models

- Large format semantic models provide a highly compressed in-memory cache for optimized query performance, enabling fast user interactivity
- Normally used in your semantic models are > 10GB in size.
- Available on Premium P SKUs, Embedded A SKUs, and with Premium Per User (PPU)
- Used commonly when connecting third-party tools via the XMLA endpoint
- On-demand loading (same feature as Direct Lake)

Fabric Analytics Concepts Notes

DAX Variables, functions & parameters

AREA OF FOCUS: Variables																																																																																																															
NOTES:		IMPLEMENTATION:																																																																																																													
<p>DAX variables can help avoid code repetition and improve potentially performance too.</p> <p>You can create a DAX variable with the keyword VAR VariableName = {Your Expression}</p> <p>You will need to return a value, using the RETURN keyword.</p>		<table><thead><tr><th>year</th><th>Month</th><th>Count of Date</th><th>Sum of Revenue</th><th>Average Revenue per Day</th></tr></thead><tbody><tr><td>2017</td><td>1</td><td>31</td><td>629202447</td><td>20,296,853.13</td></tr><tr><td>2017</td><td>2</td><td>28</td><td>503865119</td><td>17,995,182.82</td></tr><tr><td>2017</td><td>3</td><td>31</td><td>524612179</td><td>16,922,973.52</td></tr><tr><td>2017</td><td>4</td><td>30</td><td>465174983</td><td>15,505,832.77</td></tr><tr><td>2017</td><td>5</td><td>31</td><td>506889924</td><td>16,351,287.87</td></tr><tr><td>2017</td><td>6</td><td>30</td><td>453424738</td><td>15,114,157.93</td></tr><tr><td>2017</td><td>7</td><td>31</td><td>446893848</td><td>14,415,930.58</td></tr><tr><td>2017</td><td>8</td><td>31</td><td>455363932</td><td>14,689,159.10</td></tr><tr><td>2017</td><td>9</td><td>30</td><td>363034438</td><td>12,101,147.93</td></tr><tr><td>2017</td><td>10</td><td>31</td><td>476420417</td><td>15,368,400.55</td></tr><tr><td>2017</td><td>11</td><td>30</td><td>397166455</td><td>13,238,881.83</td></tr><tr><td>2017</td><td>12</td><td>31</td><td>374976719</td><td>12,096,023.19</td></tr><tr><td>2018</td><td>1</td><td>31</td><td>393405848</td><td>12,690,511.23</td></tr><tr><td>2018</td><td>2</td><td>28</td><td>459268244</td><td>16,402,437.29</td></tr><tr><td>2018</td><td>3</td><td>31</td><td>512634612</td><td>16,536,600.39</td></tr><tr><td>2018</td><td>4</td><td>30</td><td>528785573</td><td>17,626,185.77</td></tr><tr><td>2018</td><td>5</td><td>31</td><td>448659487</td><td>14,472,886.68</td></tr><tr><td>2018</td><td>6</td><td>30</td><td>469669010</td><td>15,655,633.67</td></tr><tr><td>2018</td><td>7</td><td>31</td><td>468266332</td><td>15,105,365.55</td></tr><tr><td>Total</td><td></td><td>1247</td><td>19302539848</td><td>15,479,181.91</td><td></td></tr></tbody></table>				year	Month	Count of Date	Sum of Revenue	Average Revenue per Day	2017	1	31	629202447	20,296,853.13	2017	2	28	503865119	17,995,182.82	2017	3	31	524612179	16,922,973.52	2017	4	30	465174983	15,505,832.77	2017	5	31	506889924	16,351,287.87	2017	6	30	453424738	15,114,157.93	2017	7	31	446893848	14,415,930.58	2017	8	31	455363932	14,689,159.10	2017	9	30	363034438	12,101,147.93	2017	10	31	476420417	15,368,400.55	2017	11	30	397166455	13,238,881.83	2017	12	31	374976719	12,096,023.19	2018	1	31	393405848	12,690,511.23	2018	2	28	459268244	16,402,437.29	2018	3	31	512634612	16,536,600.39	2018	4	30	528785573	17,626,185.77	2018	5	31	448659487	14,472,886.68	2018	6	30	469669010	15,655,633.67	2018	7	31	468266332	15,105,365.55	Total		1247	19302539848	15,479,181.91	
year	Month	Count of Date	Sum of Revenue	Average Revenue per Day																																																																																																											
2017	1	31	629202447	20,296,853.13																																																																																																											
2017	2	28	503865119	17,995,182.82																																																																																																											
2017	3	31	524612179	16,922,973.52																																																																																																											
2017	4	30	465174983	15,505,832.77																																																																																																											
2017	5	31	506889924	16,351,287.87																																																																																																											
2017	6	30	453424738	15,114,157.93																																																																																																											
2017	7	31	446893848	14,415,930.58																																																																																																											
2017	8	31	455363932	14,689,159.10																																																																																																											
2017	9	30	363034438	12,101,147.93																																																																																																											
2017	10	31	476420417	15,368,400.55																																																																																																											
2017	11	30	397166455	13,238,881.83																																																																																																											
2017	12	31	374976719	12,096,023.19																																																																																																											
2018	1	31	393405848	12,690,511.23																																																																																																											
2018	2	28	459268244	16,402,437.29																																																																																																											
2018	3	31	512634612	16,536,600.39																																																																																																											
2018	4	30	528785573	17,626,185.77																																																																																																											
2018	5	31	448659487	14,472,886.68																																																																																																											
2018	6	30	469669010	15,655,633.67																																																																																																											
2018	7	31	468266332	15,105,365.55																																																																																																											
Total		1247	19302539848	15,479,181.91																																																																																																											

```
1 Average Revenue per Day =
2 VAR TotalRevenue = SUM('revenue'[Revenue])
3 VAR TotalDays = DISTINCTCOUNT('date'[Date])
4 RETURN
5 IF(TotalDays > 0, TotalRevenue / TotalDays, BLANK())
```

Fabric Analytics Concepts Notes

AREA OF FOCUS: Iterators	
<p>NOTES:</p> <p>Iterator functions in Power BI enumerate through all rows in a table, performing some calculation (depending on the iterator function you choose) and aggregating the result.</p> <p>Examples include SUMX, COUNTX.</p> <p>Here we're using it to perform a cumulative calculation</p>	<p>IMPLEMENTATION:</p> <p>Cumulative Revenue by Date</p> <p>19bn Sum Revenue</p>

1 Cumulative Revenue =

```
2 SUMX(  
3     FILTER(  
4         ALL('date'[Date]),  
5         'date'[Date] <= MAX('date'[Date]))  
6     ),  
7     [Sum Revenue]  
8 )
```

Fabric Analytics Concepts Notes

AREA OF FOCUS: Table filtering																																									
<p>NOTES: The FILTER function returns a table that represents a subset of another table or expression.</p> <p>For example: FILTER('InternetSales_USD', RELATED('SalesTerritory'[SalesTerritoryCountry])<>"United States")</p>	<p>IMPLEMENTATION:</p> <p>Total Revenue by Category by Product_Name</p> <table border="1"><thead><tr><th>Product_Name</th><th>Total Revenue by Category</th></tr></thead><tbody><tr><td>Toyota</td><td>1.5bn</td></tr><tr><td>Nissan</td><td>1.4bn</td></tr><tr><td>Ford</td><td>1.3bn</td></tr><tr><td>Honda</td><td>1.2bn</td></tr><tr><td>Hyundai</td><td>1.1bn</td></tr><tr><td>Maruti Suzuki</td><td>1.0bn</td></tr><tr><td>Tata</td><td>1.0bn</td></tr><tr><td>Audi</td><td>1.0bn</td></tr><tr><td>Mahindra</td><td>0.9bn</td></tr><tr><td>Chevrolet</td><td>0.9bn</td></tr><tr><td>BMW</td><td>0.8bn</td></tr><tr><td>Renault</td><td>0.7bn</td></tr><tr><td>Volkswagen</td><td>0.6bn</td></tr><tr><td>GMC</td><td>0.5bn</td></tr><tr><td>Jeep</td><td>0.5bn</td></tr><tr><td>Cadillac</td><td>0.5bn</td></tr><tr><td>Kia</td><td>0.4bn</td></tr><tr><td>Mercedes-Benz</td><td>0.4bn</td></tr><tr><td>Lincoln</td><td>0.3bn</td></tr></tbody></table>	Product_Name	Total Revenue by Category	Toyota	1.5bn	Nissan	1.4bn	Ford	1.3bn	Honda	1.2bn	Hyundai	1.1bn	Maruti Suzuki	1.0bn	Tata	1.0bn	Audi	1.0bn	Mahindra	0.9bn	Chevrolet	0.9bn	BMW	0.8bn	Renault	0.7bn	Volkswagen	0.6bn	GMC	0.5bn	Jeep	0.5bn	Cadillac	0.5bn	Kia	0.4bn	Mercedes-Benz	0.4bn	Lincoln	0.3bn
Product_Name	Total Revenue by Category																																								
Toyota	1.5bn																																								
Nissan	1.4bn																																								
Ford	1.3bn																																								
Honda	1.2bn																																								
Hyundai	1.1bn																																								
Maruti Suzuki	1.0bn																																								
Tata	1.0bn																																								
Audi	1.0bn																																								
Mahindra	0.9bn																																								
Chevrolet	0.9bn																																								
BMW	0.8bn																																								
Renault	0.7bn																																								
Volkswagen	0.6bn																																								
GMC	0.5bn																																								
Jeep	0.5bn																																								
Cadillac	0.5bn																																								
Kia	0.4bn																																								
Mercedes-Benz	0.4bn																																								
Lincoln	0.3bn																																								

```
1 Total Revenue by Category =  
2 CALCULATE(  
3     SUM('revenue'[Revenue]),  
4     FILTER(  
5         'products',  
6         'products'[Product_Name] = SELECTEDVALUE('products'[Product_Name])  
7     )  
8 )
```

AREA OF FOCUS: **Window functions**

NOTES :

Window functions allow you to perform calculations within a specific window of your table data.

Some use cases for a window might be:

- time-based (e.g. a 3 month moving average),
- or a window of a categorical variable (e.g. average revenue for each department in a company).

Fabric Analytics Concepts Notes

Applies to: Calculated column Calculated table Measure Visual calculation

Returns multiple rows which are positioned within the given interval.

Syntax

DAX

Copy

WINDOW (`from[, from_type]`, `to[, to_type]`[, `<relation>` or `<axis>`][, `<options>`]

Parameters



Expand table

Term	Definition
from	Indicates where the window starts. It can be any DAX expression that returns a scalar value. The behavior depends on the <code><from_type></code> parameter: - If <code><from_type></code> is REL, the number of rows to go back (negative value) or forward (positive value) from the current row to get the first row in the window.
from_type	Modifies behavior of the <code><from></code> parameter. Possible values are ABS (absolute) and REL (relative). Default is REL.
to	Same as <code><from></code> , but indicates the end of the window. The last row is included in the window.
to_type	Same as <code><from_type></code> , but modifies the behavior of <code><to></code> .
relation	(Optional) A table expression from which the output rows are returned. <small>If specified, all columns in <code><partitionBy></code> must come from it or a related table.</small>

Example 1 - measure

The following measure:

```
DAX Copy  
  
3-day Average Price =  
AVERAGEX(  
    WINDOW(  
        -2, REL, 0, REL,  
        SUMMARIZE(ALLSELECTED('Sales'), 'Date'[Date], 'Product'[Product]),  
        ORDERBY('Date'[Date]),  
        KEEP,  
        PARTITIONBY('Product'[Product])  
    ),  
    CALCULATE(AVERAGE(Sales[Unit Price]))  
)
```

Fabric Analytics Concepts Notes

AREA OF FOCUS: Information functions

NOTES:

DAX information functions look at the cell or row that is provided as an argument and tells you whether the value matches the expected type.

For example, the ISERROR function returns TRUE if the value that you reference contains an error.

Useful examples include:

CONTAINS

CONTAINSSTRING

HASONEVALUE

ISBLANK

ISERROR

SELECTEDMEASURE

USERPRINCIPALNAME

IMPLEMENTATION:

TransactionId	Revenue	IsRevenueBlank
1		No revenue recorded
2	800	Revenue recorded
3	200	Revenue recorded
4		No revenue recorded
5	400	Revenue recorded
8	230	Revenue recorded

```
1 IsRevenueBlank =
2 IF(
3     ISBLANK(SUM('SampleTransactions'[Revenue])),
4     "No revenue recorded",
5     "Revenue recorded"
6 )
```

an argument and tells you whether

Fabric Analytics Concepts Notes

AREA OF FOCUS: Calculation groups

NOTES:

Calculation groups are a simple way to reduce the number of measures in a model by grouping common measure expressions.

Calculation groups work with existing explicit DAX measures by automating repetitive patterns.

Calculation Groups can be created in Power BI Desktop & also Tabular Editor.

IMPLEMENTATION:

year	Total	Daily average	Monthly average
2017	5597025199	16657813	466418767
2018	5922609476	17267083	493550790
2019	5329967167	15815926	444163931
2020	2452938006	17274211	490587601

year	Total	Daily average	Monthly average
2017	336	1	28
2018	343	1	28.58333333333332
2019	337	1	28.08333333333332
2020	142	1	28.4

AREA OF FOCUS: Dynamic string formatting

NOTES:

Dynamic format strings - With *dynamic format strings for measures*, you can determine how measures appear in visuals by conditionally applying a format string with a separate DAX expression.

The difference between this and the modelling tab Format string is that the data type is maintained.

IMPLEMENTATION:

year	Month	Dynamic Format Measure
2017	1	629.2M
2017	2	503.9M
2017	3	524.6M
2017	4	465.2M
2017	5	506.9M
2017	6	453.4M
2017	7	446.9M
2017	8	455.4M
2017	9	363.0M
2017	10	476.4M
2017	11	397.2M
2017	12	375.0M
2018	1	393.4M
2018	2	459.3M
2018	3	512.6M
2018	4	528.8M
2018	5	448.7M
2018	6	469.7M
2018	7	468.3M
Total		19.3B

19.3B
Dynamic Format Measure

Fabric Analytics Concepts Notes

```
1 Dynamic Format Measure =  
2 VAR SumRev = [Sum Revenue]  
3 RETURN  
4     SWITCH(  
5         TRUE(),  
6             SumRev < 1000, FORMAT(SumRev, "0.00"),  
7             SumRev < 1000000, FORMAT(SumRev/1000, "0.0K"),  
8             SumRev < 1000000000, FORMAT(SumRev/1000000, "0.0M"),  
9             FORMAT([SumRev/1000000000, "0.0B"]  
10        )
```

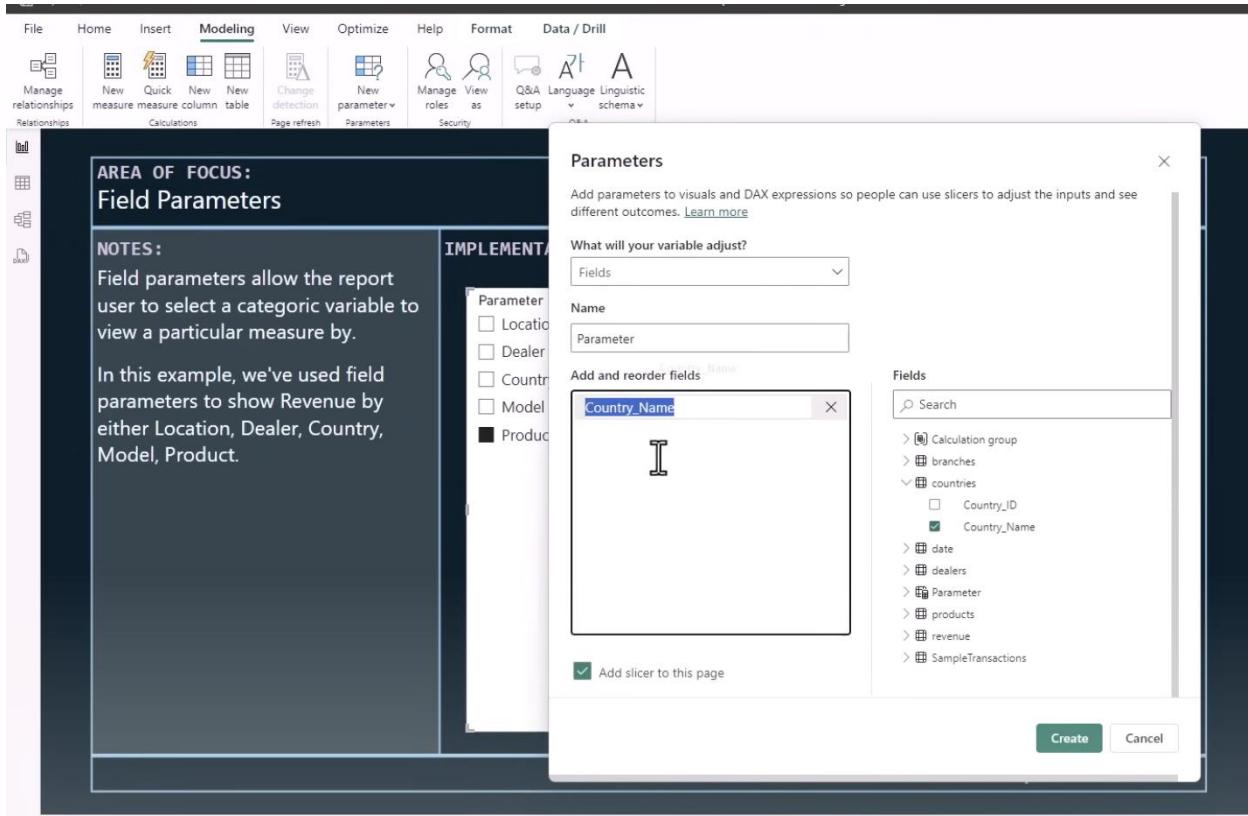
AREA OF FOCUS: Field Parameters																																					
<p>NOTES: Field parameters allow the report user to select a categoric variable to view a particular measure by. In this example, we've used field parameters to show Revenue by either Location, Dealer, Country, Model, Product.</p>	<p>IMPLEMENTATION:</p> <p>Parameter</p> <ul style="list-style-type: none"><input type="checkbox"/> Location<input type="checkbox"/> Dealer<input type="checkbox"/> Country<input type="checkbox"/> Model<input type="checkbox"/> Product <p>Sum Revenue by Location</p> <table border="1"><caption>Data for 'Sum Revenue by Location' chart</caption><thead><tr><th>Location</th><th>Sum Revenue</th></tr></thead><tbody><tr><td>(Blank)</td><td>~0.45bn</td></tr><tr><td>Chittagong</td><td>~0.35bn</td></tr><tr><td>Gyumri</td><td>~0.30bn</td></tr><tr><td>Santa Ana</td><td>~0.25bn</td></tr><tr><td>Spanish Town</td><td>~0.20bn</td></tr><tr><td>Bimbo</td><td>~0.18bn</td></tr><tr><td>Anse Etoile</td><td>~0.15bn</td></tr><tr><td>Santiago de Cuba</td><td>~0.12bn</td></tr><tr><td>Alajuela</td><td>~0.10bn</td></tr><tr><td>Sisimiut</td><td>~0.08bn</td></tr><tr><td>Freeport</td><td>~0.06bn</td></tr><tr><td>Cul De Sac</td><td>~0.05bn</td></tr><tr><td>Lelydorp</td><td>~0.04bn</td></tr><tr><td>Lubumbashi</td><td>~0.03bn</td></tr><tr><td>León</td><td>~0.02bn</td></tr><tr><td>Tartu</td><td>~0.01bn</td></tr><tr><td>Portsmouth</td><td>~0.005bn</td></tr></tbody></table>	Location	Sum Revenue	(Blank)	~0.45bn	Chittagong	~0.35bn	Gyumri	~0.30bn	Santa Ana	~0.25bn	Spanish Town	~0.20bn	Bimbo	~0.18bn	Anse Etoile	~0.15bn	Santiago de Cuba	~0.12bn	Alajuela	~0.10bn	Sisimiut	~0.08bn	Freeport	~0.06bn	Cul De Sac	~0.05bn	Lelydorp	~0.04bn	Lubumbashi	~0.03bn	León	~0.02bn	Tartu	~0.01bn	Portsmouth	~0.005bn
Location	Sum Revenue																																				
(Blank)	~0.45bn																																				
Chittagong	~0.35bn																																				
Gyumri	~0.30bn																																				
Santa Ana	~0.25bn																																				
Spanish Town	~0.20bn																																				
Bimbo	~0.18bn																																				
Anse Etoile	~0.15bn																																				
Santiago de Cuba	~0.12bn																																				
Alajuela	~0.10bn																																				
Sisimiut	~0.08bn																																				
Freeport	~0.06bn																																				
Cul De Sac	~0.05bn																																				
Lelydorp	~0.04bn																																				
Lubumbashi	~0.03bn																																				
León	~0.02bn																																				
Tartu	~0.01bn																																				
Portsmouth	~0.005bn																																				

DP-600 Exam Preparation Course

```
1 Parameter = {  
2     ("Location", NAMEOF('dealers'[Location_NM]), 0),  
3     ("Dealer", NAMEOF('dealers'[Dealer_NM]), 1),  
4     ("Country", NAMEOF('countries'[Country_Name]), 2),  
5     ("Model", NAMEOF('products'[Model_Name]), 3),  
6     ("Product", NAMEOF('products'[Product_Name]), 4)  
7 }
```

This example, we've used field parameters to select a location for the report.

Fabric Analytics Concepts Notes



The DAX expression AVERAGEX is an example of:

A) An information function

B) A calculation group

C) Table filtering

D) A window function

E) An iterator function

AVGX iterates over each row of a specified table, computes the expression, and then averages the results.

Fabric Analytics Concepts Notes

Question

On-demand loading (loading only the data that is needed for a query) is a feature of which TWO of the following:

- A) Import mode
- B) Direct Lake mode
- C) Direct Query mode
- D) Large format semantic models
- E) The XMLA endpoint

Question

Dynamic format strings overcome which significant limitation that comes from using the DAX FORMAT function?

- A) The Format function is slow on large datasets
- B) The Format function returns a string value, so the values can't be used in chart visuals.
- C) The Format function can't handle date locale whilst formatting.
- D) The Format function can't be used with Field Parameters

The DAX expressions SELECTEDMEASURE() is most likely found in the construction of which of the following:

- A) Calculation item (Calculation Groups)
- B) Field parameters
- C) An iterator function
- D) Large format semantic models
- E) A window function

-- Calculation Item: "YTD"
`CALCULATE(SELECTEDMEASURE(), DATESYTD(DimDate[Date]))`

Calculation Groups with SELECTEDMEASURE() are ideal for:

- Time intelligence (e.g., YTD, MoM).
- Dynamic formatting (e.g., "%" vs. "\$").
- Scenario analysis (e.g., "Actuals" vs. "Forecast").

Fabric Analytics Concepts Notes

Which of the following is an irreversible operations (can't be changed afterwards)?

- A) Changing the cross-filtering of a relationship to bidirectional
- B) Changing the Storage mode of a table to Import**
- C) Naming a Calculation Group
- D) Converting a semantic model into a large-format semantic model
- E) Creating a window function

When a table's storage mode is changed to Import in Power BI or Analysis Services, the data is physically copied into the model. If it was previously in Direct Query or Direct Lake mode, switching to Import discards the real-time query connection behavior. While technically you can switch it back in some tools, doing so may lose configurations or break dependencies, and in many enterprise-level environments, it's considered a one-way change due to downstream impacts.

Ch10. Secure & Optimize Semantic Models

Dynamic row-level security (in semantic models)

- In general, row-level security restricts who can see what data (at the row-level) in specific tables in your Power BI report.
- Dynamic RLS is a method of applying Row-Level Security using the UserPrincipalName() Information Function.
- You can configure row-level security in the semantic model, the data warehouse and the T-SQL endpoint of the Lakehouse in Fabric.
 - But if you're using Direct Lake mode, you need to configure RLS in the semantic model (otherwise it will fallback to Direct Query).
- RLS only restricts data access for users with Viewer permissions. It doesn't apply to Admins, Members, or Contributors.

High-level steps to implement RLS in Power BI Desktop:
1. Create a Role
2. Select the table you want to apply RLS to.
3. Enter a table filter DAX expression to configure when/who the RLS is applied.
4. Validate that RLS has been applied correctly.

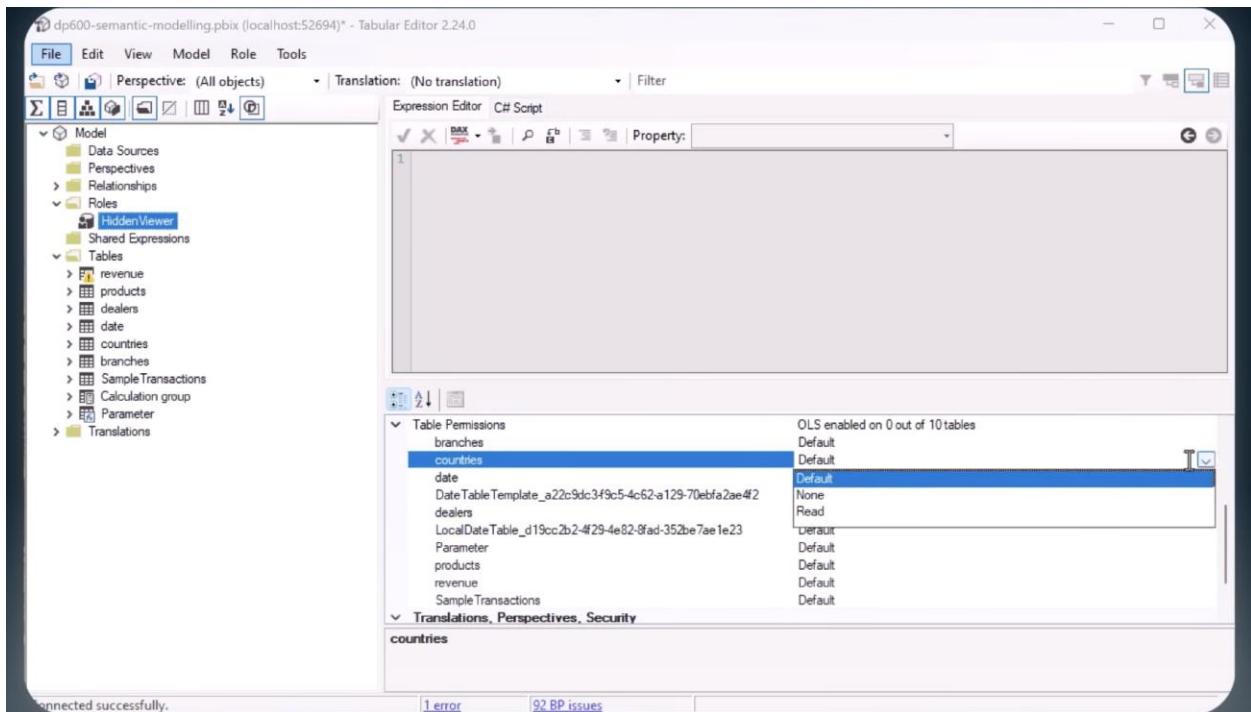
Fabric Analytics Concepts Notes

Object-level security (using Tabular Editor)

- With object-level security, you (as a model author) can configure who can view tables and specific columns within a Power BI report.
- To configure object-level security, you need to use an external tool, such as Tabular Editor.
- Similarly to RLS, OLS only restricts data access for users with Viewer permissions. It doesn't apply to Admins, Members, or Contributors.

High-level steps to implement OLS in Tabular Editor:

- Create a Role (in Power BI Desktop), if you haven't already
- Open Tabular Editor, find the Role, click on Table Properties for that Role
- Set the permissions for the table to *None* or *Read*
- Publish the report to the Service, and add people/groups to the Role.

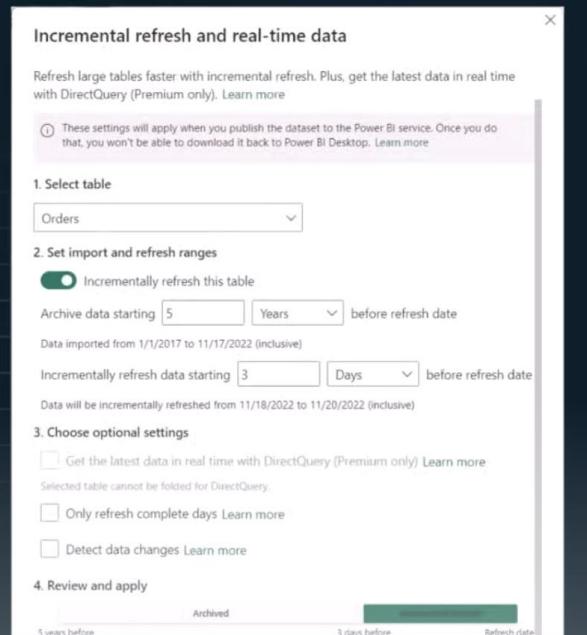


Incremental refresh (Power BI)

- Typically used on large Fact tables, incremental refresh allows you to pull in only the data that has changed in a given time period, rather than a full load of the data from source. This has the following benefits:
 - Fewer traditional refreshes are required
 - Refreshes are a lot quicker (only pulling in what is new)
 - Resource consumption is reduced
 - Refreshes can be more reliable (no long-running connections needed to your data source)
- Currently, incremental refresh is only possible within Power BI (not in any of the Fabric ETL items like a Dataflow Gen2 or Data Pipeline).
- Incremental refresh in Power BI is available for Power BI Premium licenses only (PPU or a Premium capacity subscription)
- Incremental refresh policies are defined in Power BI Desktop

Implementing incremental refresh (Power BI)

- Create RangeStart and RangeEnd parameters
- In Power Query, apply custom date filters on your table's date column, using the RangeStart and RangeEnd parameters
- Define your incremental refresh policy



Monitoring semantic model performance

For analyzing Power Query performance:

- You can use the Query Analyzer tool

For analyzing visual and query performance

- You can use Performance Analyzer in Power BI Desktop

For DAX performance:

- You can use DAX Studio

For semantic model performance:

- You can use the Best Practice Analyzer in Tabular Editor



Let's look at these three in more detail

DAX Studio

Optimizing DAX performance using DAX Studio

- You can load Power BI Performance Analyzer data into DAX Studio for further analysis
 - Better filtering and sorting than in Power BI Desktop, plus you can view the queries behind each visual load.
- Use View Metrics to look at the VertiPaq Analyzer:
 - Review Table and Column Sizes (and look for reasons WHY - cardinality, datatypes, other).
 - Look for referential integrity violations (mismatch in unique keys on two sides of a relationships).
 - Much more...
- Trace analysis:
 - All Queries trace: captures query events from client tools (like Power BI)
 - Query Plan trace: query plan trace events from a SSAS Tabular server
 - Server Timing trace: query timing from the server perspective

Fabric Analytics Concepts Notes

DAX Studio - 3.0.11

Advanced

File Home Advanced Help

Import Metrics Export Metrics View Metrics Export Data View As Benchmark SQL Profiler Analyze in Excel Swap Delimiters

Metrics Export Security Performance External Tools Utilities

Metadata Functions DMV

dp600-semantic-modelling Model

Search branches Calculation group countries date DateTableTemplate_a22c9dc3 (9...) dealers LocalDateTable_d19cc2b2-4f29... Parameter products revenue SampleTransactions

```

1 //=====
2 //| Operation : 1
3 //| Visual   : Slicer
4 //| Query Start : 18/05/2024 11:45:05
5 //| Query End  : 18/05/2024 11:45:05
6 //| Render Start : 18/05/2024 11:45:05
7 //| Render End  : 18/05/2024 11:45:05
8 //| Query Duration : 3 ms
9 //| Render Duration : 9 ms
10 //| Total Duration : 12 ms
11 //| Row Count   : 5
12 //=====
13 DEFINE
14     VAR __DSOCore =
15         SUMMARIZE(
16             VALUES('Parameter'),
17             'Parameter'[Parameter Fields],
18             'Parameter'[Parameter Order],
19             'Parameter'[Parameter]
20         )
21
22     VAR __DSOPrimaryWindowed =
23         TOPN(
24             101,
25             _____
26         )

```

Log Results History PBI Performance **VertiPaq Analyzer**

Tables	Name	Cardinality	Total Size ↓	Data	Dictionary	Hier Size	Encoding	Data Type	RI Violations	User Hier Size	Ref Size	% Tab
Columns	revenue	1,861	195,508	18,576	135,628	35,840	Many	-	-	0	5,464	
Relationships	date	1,247	171,246	7,544	130,798	30,160	Many	-	-	736	2,008	
Partitions	branches	2,470	148,493	9,416	115,349	23,728	Many	-	-	0	0	
Summary	dealers	262	132,876	2,168	120,284	10,176	Many	-	1	0	248	
	LocalDateTable...	1,461	131,904	6,096	89,008	12,320	Many	-	-	24,480	0	
	products	277	89,746	1,728	83,010	5,008	Many	-	-	0	0	
	countries	246	53,460	880	48,644	3,936	Many	-	1	0	0	
	DateTableTemp...	1	36,252	1,080	35,004	72	Many	-	-	96	0	
	Parameter	5	35,496	536	34,840	120	Many	-	-	0	0	
	Calculation group	3	17,896	400	17,448	48	Many	-	-	0	0	
	SampleTransact...	6	840	400	384	56	VALUE	-	-	0	0	

Log Results History PBI Performance **VertiPaq Analyzer**

Tables	Name	Cardinality	Total Size ↓	Data	Dictionary	Hier Size	Encoding	Data Type	RI Violations	User Hier Size	Ref Size	% Tab
Columns	revenue	1,861	195,508	18,576	135,628	35,840	Many	-	-	0	5,464	
Relationships	Branch_ID	1,843	72,653	3,112	54,789	14,752	HASH	String	-	-	-	3
Partitions	Date_ID	1,158	44,061	3,112	31,669	9,280	HASH	String	-	-	-	2
Summary	Model_ID	277	28,558	2,256	24,078	2,224	HASH	String	-	-	-	1
	Dealer_ID	267	27,904	2,256	23,504	2,144	HASH	String	-	-	-	1
	Revenue	1,851	15,112	7,576	128	7,408	VALUE	Int64	-	-	-	
	Units_Sold	3	1,500	136	1,332	32	HASH	Int64	-	-	-	
Tables	date	1,247	171,246	7,544	130,798	30,160	Many	-	-	736	2,008	
Relationships	Date	1,247	57,288	2,128	45,176	9,984	HASH	DateTime	-	-	-	3
Partitions	Date_ID5	1,247	44,931	2,128	32,819	9,984	HASH	String	-	-	-	2
Summary	Date_ID0	1,247	44,931	2,128	32,819	9,984	HASH	String	-	-	-	2
	Quarter	4	17,344	144	17,152	48	HASH	String	-	-	-	1
	Month	12	2,232	752	1,368	112	HASH	Int64	-	-	-	
	year	4	1,520	136	1,336	48	HASH	Int64	-	-	-	
Columns	branches	2,470	148,493	9,416	115,349	23,728	Many	-	-	0	0	
Relationships	dealers	262	132,876	2,168	120,284	10,176	Many	-	1	0	248	
Partitions	LocalDateTable...	1,461	131,904	6,096	89,008	12,320	Many	-	-	24,480	0	
Summary	products	277	89,746	1,728	83,010	5,008	Many	-	-	0	0	
	countries	246	53,460	880	48,644	3,936	Many	-	1	0	0	
	DateTableTemp...	1	36,252	1,080	35,004	72	Many	-	-	96	0	
	Parameter	5	35,496	536	34,840	120	Many	-	-	0	0	
	Calculation group	3	17,896	400	17,448	48	Many	-	-	0	0	
	SampleTransact...	6	840	400	384	56	VALUE	-	-	0	0	

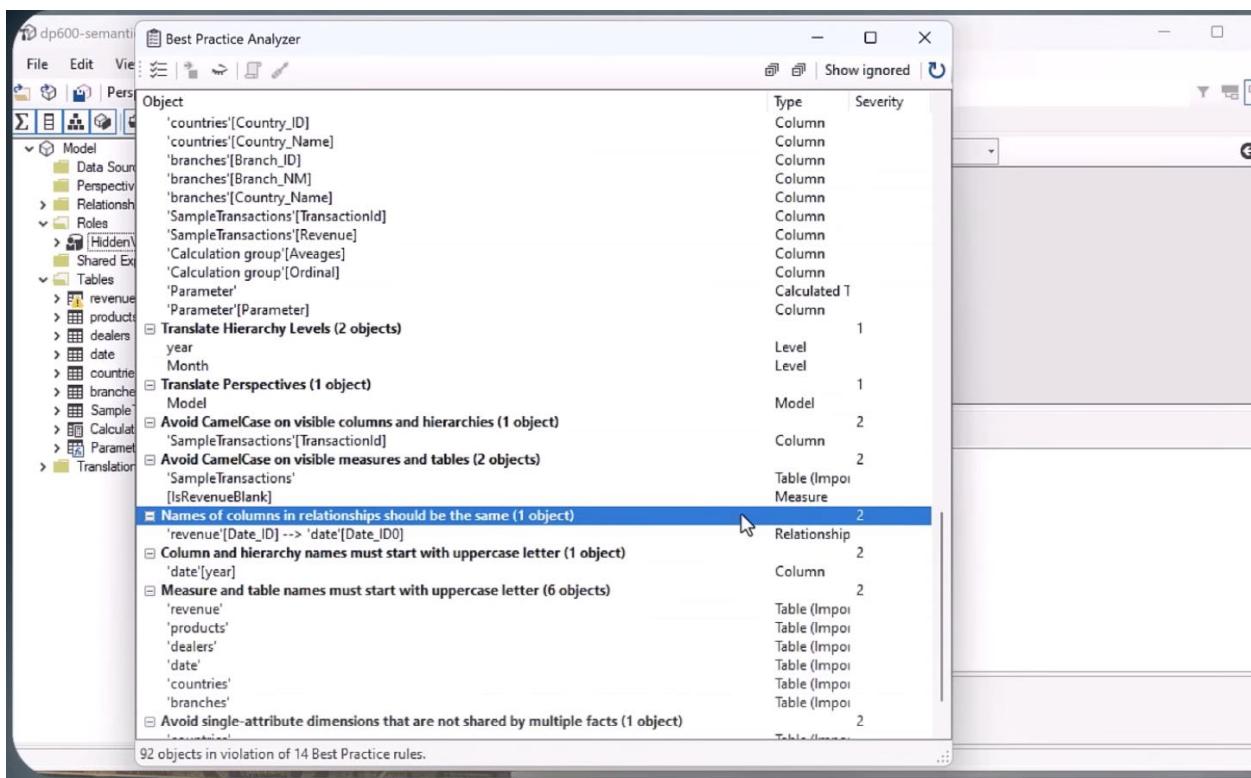
Fabric Analytics Concepts Notes

Tabular Editor

Optimizing semantic models using Tabular Editor

Best Practice Analyzer:

- This is a tool that performs a scan of your semantic model and checks for common issues.
 - The list of rules can be downloaded from GitHub, and are organized into the following categories:
 - Performance
 - DAX Expressions
 - Error Prevention
 - Formatting
 - Maintenance
 - The checks can also be run from the Tabular Editor CLI as part of a CI/CD process.



Use cases for DAX Studio and Tabular Editor

DAX Studio use cases

- Write, execute, and debug DAX queries (need to be manually copied over to Power BI Desktop, as DAX Studio is 'Read-Only').
- Use the VertiPaq Analyzer to understand the size of your semantic model (as well as individual tables and columns, within the model).
- Bring Power BI Desktop Performance Analyzer data into DAX Studio for further analysis.
- Perform trace analysis

Tabular Editor use cases

- Quickly edit data models - create measures, perspectives, calculation groups from the DAX editor (and publish them into a semantic model)
- Automate repetitive tasks using scripting
- Incorporate DevOps with tabular models
- Use the Best Practice Analyzer to identify common issues in Power BI semantic models.
- Implement object-level security

DAX Studio: The DAX Query & Performance Tool

Primary Purpose: Debugging, optimizing, and analyzing DAX queries/data models.

Key Use Cases:

1. DAX Query Development

- Write, test, and debug complex DAX queries in a dedicated IDE.
- *Limitation:* Queries must be manually copied to Power BI Desktop (read-only connection).

2. Performance Analysis

- **VertiPaq Analyzer:** Inspect semantic model size (tables/columns) to optimize storage (e.g., identify high-cardinality columns).
- **Trace Analysis:** Capture query execution plans to diagnose performance bottlenecks.

3. Integration with Power BI

- Import Performance Analyzer data for deeper investigation.

Tabular Editor: The Model Management & Automation Tool

Primary Purpose: Advanced modeling, automation, and DevOps for semantic models.

Key Use Cases:

1. Model Editing

- Rapidly create/edit measures, calculation groups, and perspectives without waiting for Power BI refreshes.
- *Example:* Build time-intelligence calculation groups (YTD, QoQ%) efficiently.

2. Automation & Scripting

- Use C# scripts to automate repetitive tasks (e.g., bulk measure edits, metadata updates).

3. DevOps & Best Practices

- **Best Practice Analyzer:** Flag issues (e.g., unused columns, missing descriptions).
- **Version Control:** Integrate with Git for collaborative model development.

4. Security

Fabric Analytics Concepts Notes

- Configure object-level security (OLS) for row-level security (RLS) models.

When to Use Which Tool?

Scenario	DAX Studio	Tabular Editor
Debugging a slow DAX measure	✓	⚠ (Limited)
Creating calculation groups	✗	✓
Analyzing data model size	✓	✗
Automating measure creation	✗	✓
Testing query performance	✓	✗

Pro Tip:

- **Combine Both:** Use Tabular Editor to build/optimize your model, then DAX Studio to fine-tune query performance.

Your goal is to analyze Performance Analyzer data in DAX Studio to find the visual in your report with the longest total load time (ms).

Put the following steps in the correct order to achieve this:

1. In Power BI Performance Analyzer, click Start Recording
2. Click 'Refresh Visuals' or interact with the report. Then click Stop Recording.
3. Export the Performance Data JSON file
4. Import the Performance Data JSON in DAX Studio
5. On the PBI Performance tab, sort by Total Ms (descending) to find the longest refresh time.

Step-by-Step Process

1. **Start Recording in Power BI Performance Analyzer**
 - In Power BI Desktop, go to the *View* tab → *Performance Analyzer* → Click **Start Recording**.
 - *Purpose:* Begins capturing metrics for all report interactions/refreshes.
2. **Refresh Visuals or Interact with the Report**
 - Click **Refresh Visuals** (or interact with filters/slicers).
 - After completing actions, click **Stop Recording**.
 - *Why?:* Generates performance data for all visuals during the session.
3. **Export the Performance Data as JSON**
 - In Performance Analyzer, click the **Export** button to save the data as a .json file.
 - *Key Point:* This file contains timing details (e.g., DAX query duration, visual rendering).

Fabric Analytics Concepts Notes

4. Import the JSON File into DAX Studio

- Open DAX Studio → Navigate to the *Tools* tab → Select **Import Performance Data**.
- Load the exported .json file.
- *Outcome:* DAX Studio parses the data into a queryable format.

5. Sort by Total Ms (Descending) in DAX Studio

- Go to the *PBI Performance* tab in DAX Studio.
- Sort the **Total Ms** column in descending order.
- *Result:* The visual with the highest value is your performance bottleneck.

You want to use the Best Practice Analyzer tool within Tabular Editor to assess your DAX performance.

Which of the following Severity codes for BPA rule violations indicates an Error:

A. Level 0

B. Level 1

C. Level 2 and above

D. Level 2

E. Level 3 and above

[Play with Tabular Editor 3](#) / Improve code quality with the Best Practice Analyzer

When creating a new rule, you must specify the following details:

- **Name:** The name of the rule, which will be displayed to users of Tabular Editor
- **ID:** An internal ID of the rule. Must be unique within a rule collection. If multiple rules have identical IDs across different collections, only the rule within the collection of the highest precedence is applied.
- **Severity:** The severity is not used within Tabular Editor's UI, but when running a Best Practice Analysis through [Tabular Editor's command line interface](#), the number determines how "severe" a rule violation is.
 - 1 = Information only
 - 2 = Warning
 - 3 (or above) = Error
- **Category:** This is used for logically grouping rules together to make management of rules easier

Fabric Analytics Concepts Notes

When implementing Dynamic Row-Level Security, which DAX Information Function should you use to filter the data in a specific table based on the logged-in user's email address?

A. USER()

B. USEROBJECTID()

C. USEREMAIL()

D. USERPRINCIPALNAME()

E. USERNAME()

In DAX Studio, which of the following records queries that are generated by a client tool like Power BI Desktop?

A. Query Plan trace

B. SQL Profiler

C. All Queries trace

D. Server Timings trace

E. VertiPaq Analyzer

1. All Queries Trace:

- **Purpose:** Captures *every* query sent to the semantic model by client tools (e.g., Power BI Desktop, Excel).
- **Use Case:** Identify which visuals/filters generate DAX queries and analyze their frequency/performance.

2. Why Other Options Are Incorrect:

- **A) Query Plan Trace:** Shows the *execution plan* for a *single* query (useful for optimization, not logging all queries).
- **B) SQL Profiler:** Legacy SQL Server tool (not native to DAX Studio).
- **D) Server Timings Trace:** Measures *storage engine* vs. *formula engine* processing times (doesn't log queries).
- **E) VertiPaq Analyzer:** Inspects data model structure (size/compression), not query logging.

Fabric Analytics Concepts Notes

The first step in implementing Incremental Refresh in Power BI is to:

- A. Add a RangeStart and a RangeEnd column to your dataset.
- B. Add RefreshStart and RefreshEnd parameters to your Power BI Desktop project
- C. Add a RefreshStart and a RefreshEnd column to your dataset.
- D. **Add RangeStart and RangeEnd parameters to your Power BI Desktop project**

For analyzing Power Query performance:

- You can use the Query Diagnostics tool

For analyzing visual and query performance

- You can use Performance Analyzer in Power BI Desktop

For DAX performance:

- You can use DAX Studio

For semantic model performance:

- You can use Tabular Editor

Fabric Analytics Concepts Notes

There are a number of methods you can use for improving DAX performance, we will look at using DAX Studio.

DAX Studio use cases:

- Use the VertiPaq Analyzer to understand the size of your semantic model (as well as individual tables and columns, within the model).
- Bring Power BI Desktop Performance Analyzer data into DAX Studio for further analysis.
- Trace analysis:
 - All Queries trace: captures query events from client tools (like Power BI)
 - Query Plan trace: query plan trace events from a SSAS Tabular server
 - Server Timing trace:
- Analyze DAX performance using Query Profiling
- Write, execute, and debug DAX queries (need to be manually copied over to Power BI Desktop, as DAX Studio is 'Read-Only').
- Charting and Pivot table for visualization

DAX Studio for DAX performance improvements

- You can load Power BI Performance Analyzer data into DAX Studio for further analysis

Ch11. Perform exploratory analytics

'Types' of analytics (in Microsoft's words)

1. **Descriptive analytics:** These analytics interpret past data and KPIs to identify trends and patterns. → **Describes**
2. **Diagnostic analytics:** By focusing on past performance data, these analytics decide which data element will influence specific trends and the possibility of any future events—created from techniques like data mining and correlation. → **Diagnoses**
3. **Predictive analytics:** By using statistics to forecast future outcomes with statistical models and machine learning techniques, these analytics provide context and clarity for future decisions. → **Predicts**
4. **Prescriptive analytics:** These analytics build on descriptive and predictive analytics to recommend specific actions that ensure the best or most profitable customer reactions and business outcomes possible. → **Prescribes**

PBI Visuals

Power BI visuals (and when to choose them)

Most of the content for this section of the exam, is from the PL-300 exam for Power BI Data Analyst, so that's what we'll focus on.

It's less about:

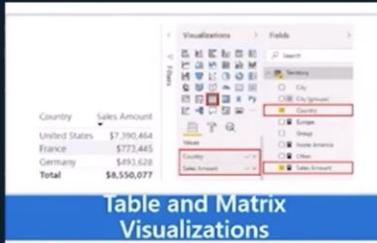
- descriptive,
- diagnostic,
- predictive, and
- prescriptive analytics,

...and more about:

- Power BI,
- it's visuals (and when to use them)
- specific Power BI features (and when to use them)

Fabric Analytics Concepts Notes

Visuals (when to choose/ not to choose)



- ✓ good for visualising fine-grained details, or for drill-throughs.
- ✓ can be used to display aggregate information (example: Monthly Revenue)
- ✗ difficult to spot long-term trends

- ✓ when you have one categoric variable and one numeric variable (example: Revenue by Region)
- ✓ can be 'stacked' to show more than one categoric variable
- ✗ not well suited to time-series information

- ✓ well-suited for time-series information
- ✓ legend can be used to show how a metric changes within each category
- ✗ area chart can be difficult to interpret



- ✓ good for KPI metrics
- ✓ can also include % change metrics, which give context
- ✗ don't give longer-term trends (can be coupled with a spark line to show 'momentum')

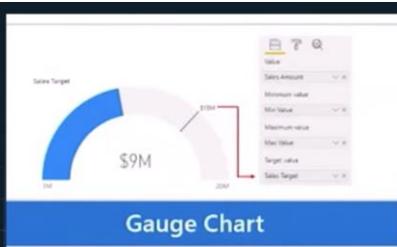
- ✓ can be used to show ratios
- ✗ difficult for human brain to compare between categories (area vs linear)
- ✗ presenting ratios obscures overall numbers
- ✗ cannot see trends over time

- ✓ typically used to visualize more than one metric on Y-axis
- ✓ ✗ allows user to come to conclusion about 'correlation' (good or bad?)
- ✗ can be difficult for user to interpret (which axis shows which metric), often needs a Legend to explain

Fabric Analytics Concepts Notes



Funnel Visualization



Gauge Chart



Waterfall Visualization

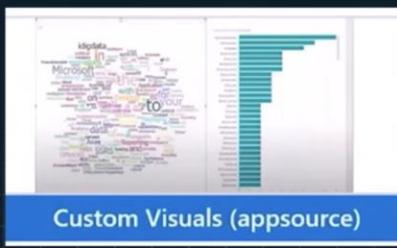
- can show movement through a linear process (example: tracking website conversion)
- difficult for user to grasp scale of difference

- can show progress of a particular metric towards a goal

- show a running total over a time period (or categories)
- can see which periods or categories contributed to change the most
- can be difficult for user to interpret



Scatter Charts



Custom Visuals (appsource)



Q&A Visualization

- visualising two numeric variables (and therefore *implied* correlation)
- Dots can be color-coded to deepen the analysis
- can sometimes lead to correlation conclusions (which != causation)

- can go beyond the out-of-the-box solutions
- normally involve some licensing/ paywall considerations

- allow the user to ask natural language questions about your data
- questions need to be carefully articulated to give best chance of good results
- simple queries only
- only in English and Spanish

Fabric Analytics Concepts Notes

1. Correlation

- **Definition:** A statistical relationship where two variables change together.
- **Key Points:**
 - **Does not imply causation** (e.g., ice cream sales and drowning incidents both rise in summer)
 - Measured by correlation coefficients (e.g., Pearson's r : -1 to $+1$).
- **Example:**
 - *Observation:* Higher social media usage correlates with higher anxiety rates.
 - *Reality:* Both may be driven by a third factor (e.g., stress).

2. Causation

- **Definition:** A direct cause-and-effect relationship where one variable *changes* another.
- **Key Points:**
 - Requires **experimental control** (e.g., randomized trials) or rigorous evidence.
 - Often summarized as " X causes Y ."
- **Example:**
 - *Observation:* Smoking increases lung cancer risk.
 - *Evidence:* Decades of controlled studies linking tobacco to cancer.

PBI Features

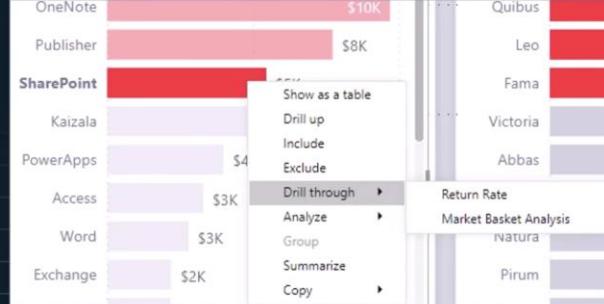
Fabric Analytics Concepts Notes

Drilldown



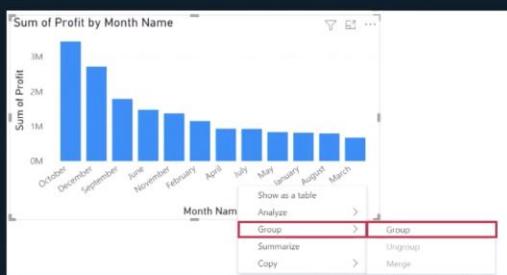
- allows user to explore data through layers of a hierarchy (either explicit or implicit)
- requires end user knowledge of the feature

Drillthrough



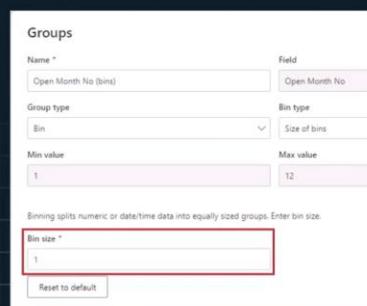
- allows user to drill-through to a separate report page, to show more detailed information for that given category.
- user journey/ navigation needs to be carefully thought
- requires end user knowledge of the feature

Grouping



- allows report author to group two or more categories

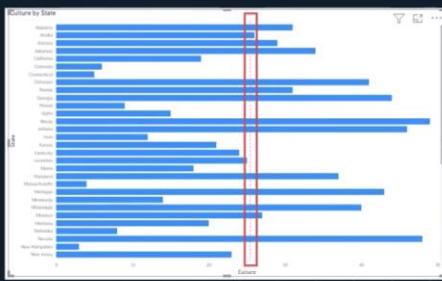
Binning



- Allows report author to create 'bins' for continuous variables
- Example: from a continuous 'Salary' variable to Salary Ranges (from \$0-\$30k, \$30k-60k etc)

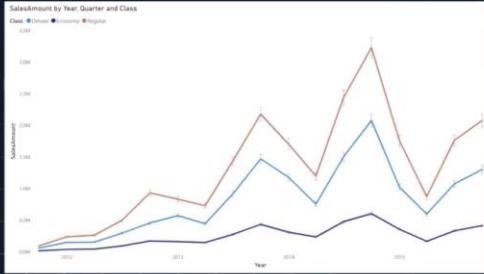
Fabric Analytics Concepts Notes

Other features



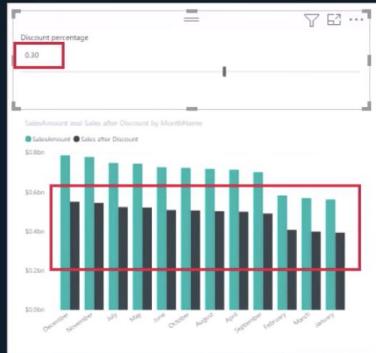
Reference lines

- Allows report authors to provide static reference lines across X or Y axis to give user some context.
- Example: average sales for all salespeople



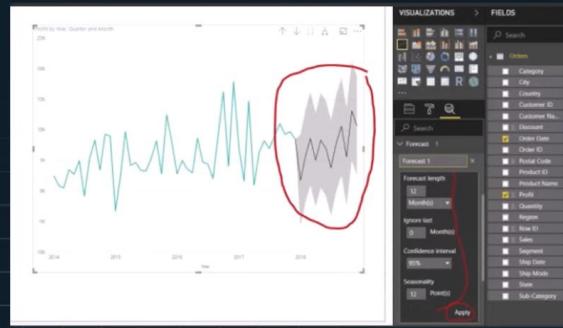
Visualising errors

- An important tool to visualise uncertainty
- Useful for showing uncertainty for a measurement or a prediction
- Can use markers, lines or shaded areas
- Requires data on errors



'What-if' parameters

- Rudimentary scenario analysis
- Requires some data engineering and/or prediction algorithm to prepare a what-if analysis



Time-series forecasting

- Gives report authors the ability to add a forecast line to time-series, continuous variables (with error bounds)
- Can be added in the Analytics pane
- Careful!

Data Profiling with Power Query

Fabric Analytics Concepts Notes

Power Query

Search (Alt + Q)

Home Transform Add column View Help

Data view Schema view Script view Diagram view Query settings

Always allow Advanced editor

Columns Parameters Advanced

Enable column profile

Show column quality details

Show column value distribution

Show column profile in details pane

Enable details pane

Monospaced

Show whitespace

✓ fx Table.SelectRows(#"Added custom", each ([Units_Sold] = 1))

Dealer_ID	Model_ID	Branch_ID	Date_ID
DLR0048	Che-M48	BR0471	DT00048
DLR0171	Hyu-M161	BR1902	DT00438
DLR0183	Hyu-M163	BR2223	DT00717
DLR0214	Tat-M194	BR0064	DT00748
DLR0039	Dod-M9	BR0984	DT00840
DLR0010	BMW-M247	BR0895	DT01078
7 DLR0242	Toy-M202	BR4446	DT00063
8 DLR0151	Toy-M101	BR0037	DT00239
9 DLR0022	Aud-M229	BR1618	DT00644
10 DLR0079	Dod-M9	BR2188	DT00701
11 DLR0250	Mah-M170	BR1630	DT01139

✓ fx Table.SelectRows(#"Added custom", each ([Units_Sold] = 1))

Model_ID	Branch_ID	Date_ID	Units_Sold	RevenueOverIM
Hon-M219	BR0615	DT0092	260 distinct, 66 unique	100% Valid, 0% Error, 0% Empty
Sko-M272	BR1145	DT00402	602 distinct, 66 unique	100% Valid, 0% Error, 0% Empty
Aud-M231	BR6186	DT00829	1 distinct, 0 unique	100% ..100% Valid, 0% ... 0% Error, 0% ... 0% Empty
Nis-M264	BR1667	DT00834	604 distinct, 66 unique	100% Valid, 0% Error, 0% Empty
Mar-M137	BR0999	DT01084	1 distinct, 0 unique	100% ..100% Valid, 0% ... 0% Error, 0% ... 0% Empty
Mar-M142	BR1049		1 distinct, 0 unique	100% ..100% Valid, 0% ... 0% Error, 0% ... 0% Empty
Cit-M115	BR1080		1 distinct, 0 unique	100% ..100% Valid, 0% ... 0% Error, 0% ... 0% Empty

Keep duplicates
Keep errors
Remove duplicates
Remove errors
Replace errors...

Fabric Analytics Concepts Notes

A company annual report shows a net profit of \$34 million dollars.

The company has 12 business units, each with their own net profit (or loss) amount.

Which of the following visual types could (best) be used to visually show how each business unit contributed to the overall net profit metric?

- A) Line chart
- B) Scatter chart
- C) Matrix visual
- D) Waterfall chart**

- E) Q & A visual

Why a Waterfall Chart?

1. Purpose-Built for Contributions:

- Clearly shows **cumulative impact** of each business unit (positive or negative) on the total net profit (\$34M).
- Starts with an initial value (often \$0), adds/subtracts each unit's profit/loss, and ends at the final total.

2. Example Visualization:

- **First bar:** Starting point (\$0).
- **Middle bars:** Each business unit's profit (green) or loss (red).
- **Final bar:** Total net profit (\$34M).

You want to add measurement error bars on a time-series line chart, to show potential error in each measurement.

Where would you go add this information to your visual?

- A) Analytics pane**
- B) Format visual pane
- C) View tab > Show error bars
- D) Build visual pane
- E) In the Model view

Fabric Analytics Concepts Notes

You want to visually compare two continuous variables (Age and Height of Survey respondents) in one chart.

Which of the following visual types would (best) represent this data?

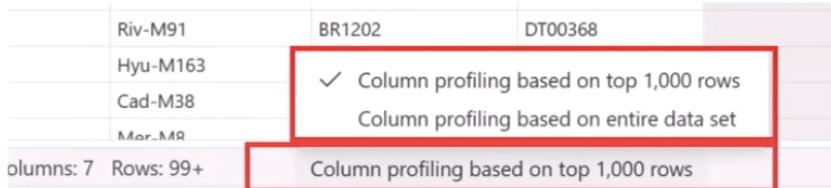
- A) Line chart
- B) Stacked bar visual
- C) Scatter visual**
- D) Q & A visual
- E) Matrix visual

By default, the Data Profiling tool reviews the top N rows of your dataset to show you potential data quality issues.

What is N ?

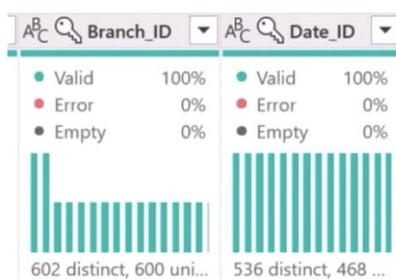
- A) 10
- B) 100
- C) 500
- D) 1000**

- E) 10000**



Which of the following features of the Data Profiling tool can help you identify duplicate values in a column that you plan to use as a key to join on?

- A) Column quality details
- B) Column value distribution**
- C) Column duplicate analysis
- D) Column key constraints
- E) Group By



Fabric Analytics Concepts Notes

Ch12. Query data using SQL

- Query a lakehouse in Fabric by using SQL queries or the visual query editor
- Query a warehouse in Fabric by using SQL queries or the visual query editor
- Connect to and query datasets by using the XMLA endpoint

Accessing the T-SQL engine

The screenshot displays two side-by-side views of the T-SQL engine interface, both titled "Visual queries".

Lakehouse T-SQL Endpoint (read-only): This view shows a "Home" tab with options like "New SQL query" and "New visual query". Below it is an "Explorer" pane showing a tree structure of databases and tables, and a main area with a SQL query editor containing the following code:

```
1 SELECT [Dealer_ID]
2     ,[Model_ID]
3     ,[Branch_ID]
4     ,[Date_ID]
5     ,[Units_Sold]
6     ,[Revenue]
7 FROM [LH_BRONZE].[dbo].[revenue]
```

Fabric Data Warehouse (read & write): This view shows a "Home" tab with options like "Get data", "New SQL query", and "New visual query". Below it is an "Explorer" pane showing a tree structure of warehouses, schemas, and tables, and a main area with a SQL query editor containing the following code:

```
1 CREATE TABLE [DW_GOLD].[dbo].[Revenue] AS
2 Select * from [LH_BRONZE].[dbo].[revenue]
3
4 CREATE TABLE [DW_GOLD].[dbo].[Products] AS
5 Select * from [LH_BRONZE].[dbo].[products]
```

Querying the XMLA endpoint

The screenshot shows the Power BI workspace settings, the "Connect to Server" dialog for SQL Server, and the Object Explorer and SQL Editor panes in SSMS.

Power BI Workspace Settings: Shows the "Workspace settings" pane with "General", "License info", "Azure connections", "System storage", "Git integration", "Create", and "Network security" sections. A "Connection link" section contains a redacted URL: powerbi.com/api/powerbi.com/v1.0/myorg/OPD0-transform.

SQL Server Connect to Server: Shows the "SQL Server" dialog with the following connection parameters:

- Server type: Analysis Services
- Server name: powerbi://api.powerbi.com/v1.0/myorg/(workspace_name)
- Authentication: Microsoft Entra MFA
- User name: (email_address)

Object Explorer: Shows the database structure with nodes like "DW_GOLD" and "LH_BRONZE".

SQL Editor: Shows the following T-SQL code:

```
-- CTAS to get us started
CREATE TABLE dbo.FactRevenue AS
SELECT * FROM [LH_BRONZE].[dbo].[revenue]

CREATE TABLE dbo.DimDate AS
SELECT * FROM [LH_BRONZE].[dbo].[date]

CREATE TABLE dbo.DimBranch AS
SELECT * FROM [LH_BRONZE].[dbo].[branches]
```

T-SQL

As well as understanding where you can write T-SQL, you will also need a pretty good understanding of T-SQL, and how to write T-SQL queries.

I don't think there exists a definitive list of which T-SQL functions you need to familiar with, but I recommend being familiar with the following areas:

- Difference between WHERE and HAVING
- GROUP BY and summarizations
- Difference between UNION and UNION ALL
- Different types of JOINS and when to use them
- Common Table Expressions
- LEAD and LAG
- ROW_NUMBER()
- Subqueries
- Cross-warehouse queries

T-SQL Concept	Purpose	Example/Syntax	Key Notes
WHERE vs. HAVING	Filters rows (WHERE) vs. filters aggregates (HAVING).	SELECT ... WHERE salary > 50000 vs. GROUP BY ... HAVING AVG(salary) > 50000	WHERE operates before GROUP BY; HAVING after.
GROUP BY	Groups rows by columns for aggregations (SUM, AVG, COUNT).	SELECT dept, AVG(salary) FROM employees GROUP BY dept	Essential for summary reports.
UNION vs. UNION ALL	Combines result sets (UNION removes duplicates; UNION ALL retains them).	SELECT ... UNION SELECT ...	UNION ALL is faster but less clean.
JOIN Types	Combines tables (INNER, LEFT, RIGHT, FULL, CROSS).	SELECT * FROM A INNER JOIN B ON A.id = B.id	LEFT JOIN keeps all left table rows; FULL keeps all from both.

Fabric Analytics Concepts Notes

T-SQL Concept	Purpose	Example/Syntax	Key Notes
Common Table Expressions (CTEs)	Creates temporary result sets for complex queries.	WITH CTE AS (SELECT ...) SELECT * FROM CTE	Improves readability; reusable in the same query.
LEAD/LAG	Accesses next/previous row values without self-join.	SELECT name, LAG(salary) OVER (ORDER BY date)	Window functions; useful for trend analysis.
ROW_NUMBER()	Assigns sequential numbers to rows (e.g., ranking).	SELECT ROW_NUMBER() OVER (PARTITION BY dept ORDER BY salary)	Often used with pagination or deduplication.
Subqueries	Nested queries (can be in SELECT, FROM, WHERE).	SELECT name FROM employees WHERE salary > (SELECT AVG(salary) FROM employees)	Correlated subqueries reference outer query.
Cross-Warehouse Queries	Queries across multiple databases/warehouses (e.g., Azure Synapse).	SELECT * FROM db1.table1 JOIN db2.table2 ON ...	Requires proper permissions and linked servers/federated queries.

Key Takeaways:

1. **WHERE** filters raw data; **HAVING** filters grouped data.
2. **UNION ALL** > **UNION** for performance (if duplicates don't matter).
3. **CTEs** simplify complex logic; **LEAD/LAG** avoid self-joins.
4. **ROW_NUMBER()** is versatile for rankings/pagination.

Fabric Analytics Concepts Notes

```
-- change the above query to only return the TOP 3 Branches in Spain
SELECT TOP 3 rev.Branch_ID, B.Branch_NM, B.Country_Name, SUM(Revenue) AS TotalRevenue FROM dbo.FactRevenue rev
INNER JOIN dbo.DimBranch B ON rev.Branch_ID = B.Branch_ID
WHERE B.Country_Name = 'Spain'
GROUP BY rev.Branch_ID, B.Country_Name, B.Branch_NM
ORDER BY SUM(rev.Revenue) DESC
```

The screenshot shows a SQL query execution window with the following output:

Branch_ID	Branch_NM	Country_Name	TotalRevenue
BR2047	Volkswagen South Africa Motors	Spain	26346189
BR2050	Volkswagen South Africa Motors	Spain	24546006
BR2048	Volkswagen South Africa Motors	Spain	14758116

For aggregates like sum, we are use having instead of where clause

```
-- what about if we want to filter after the aggregate? I.e. Give me all the Branches that
-- has a Revenue > some amount X
SELECT B.Branch_ID, B.Branch_NM, SUM(Revenue) AS TotalRevenue FROM dbo.FactRevenue rev
INNER JOIN dbo.DimBranch B ON rev.Branch_ID = B.Branch_ID
GROUP BY B.Branch_ID, B.Branch_NM
HAVING SUM(Revenue) > 50000000
```

CTE

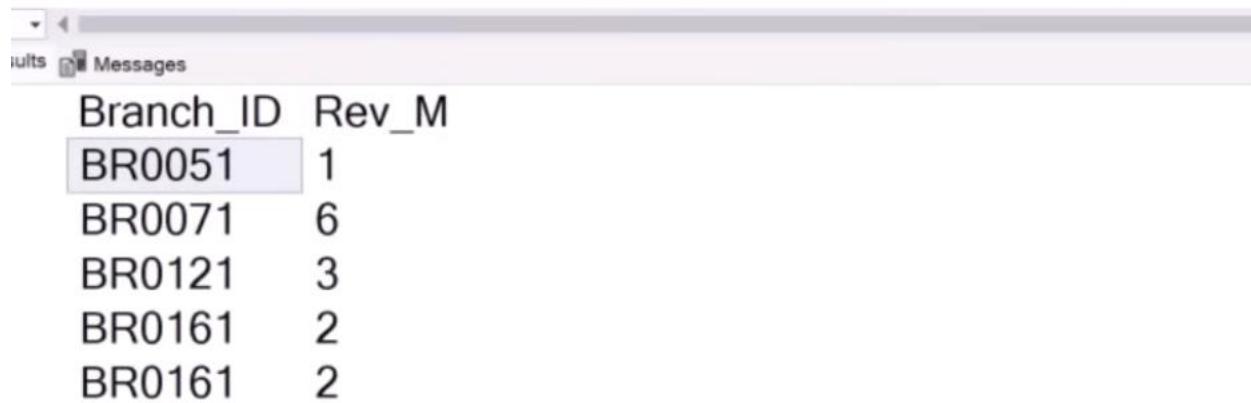
Final statement of CTE is select stmt to return sth out of CTE

```
-- Common Table Expressions
WITH top_5_rev AS (
    SELECT TOP 5 Branch_ID, SUM(Revenue) AS TotalRevenue FROM dbo.FactRevenue
    GROUP BY Branch_ID
    ORDER BY SUM(Revenue) DESC
),
branches AS (
    SELECT Branch_ID, Branch_NM FROM dbo.DimBranch
)
SELECT t5.Branch_ID, b.Branch_NM, t5.TotalRevenue FROM top_5_rev t5
LEFT JOIN branches b ON t5.Branch_ID = b.Branch_ID
```

Fabric Analytics Concepts Notes

-- LAG and LEAD

```
with revT as (
    SELECT TOP 5 Branch_ID
        , floor(Revenue/1000000) as Rev_M
    FROM dbo.FactRevenue
)
SELECT Branch_ID
    , Rev_M
    , LAG(Rev_M, 2) OVER(ORDER BY Rev_M) as lag_col
from revT
```



The screenshot shows the SQL Server Management Studio interface with the 'Results' tab selected. The results pane displays the output of the query, which is a table with two columns: 'Branch_ID' and 'Rev_M'. The data is as follows:

Branch_ID	Rev_M
BR0051	1
BR0071	6
BR0121	3
BR0161	2
BR0161	2

Fabric Analytics Concepts Notes

-- LAG and LEAD

```
with revT as (
    SELECT TOP 5 Branch_ID
        , floor(Revenue/1000000) as Rev_M
    FROM dbo.FactRevenue
)

SELECT Branch_ID
    , Rev_M
    , LAG(Rev_M, 1) OVER(ORDER BY Rev_M) as lag_col
from revT
```

Branch_ID	Rev_M	lag_col
BR0051	1	NULL
BR0161	2	1
BR0161	2	2
BR0121	3	2
BR0071	6	3

-- Subqueries

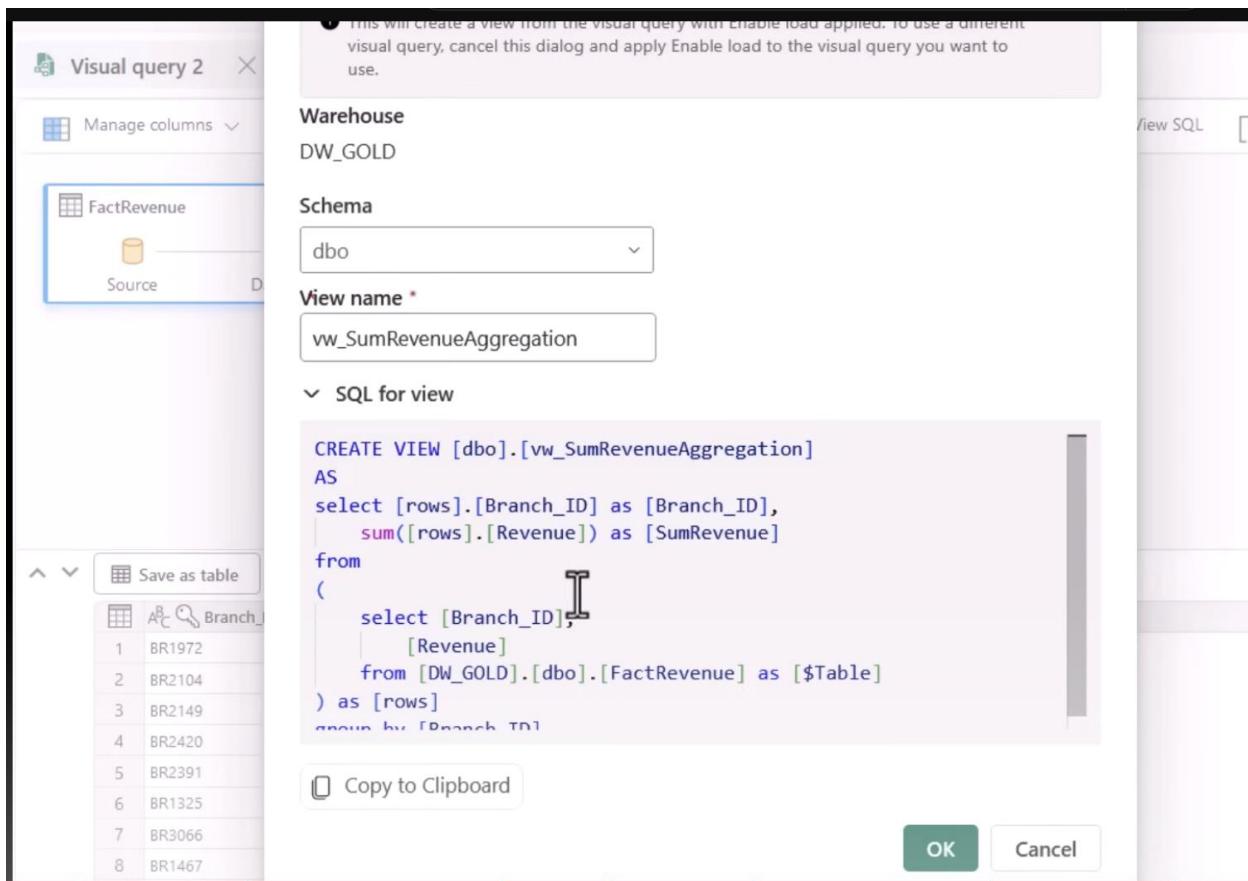
```
SELECT * from (
    SELECT * from dbo.FactRevenue
    WHERE Dealer_ID in ('DLR0001', 'DLR0017')
) sub
```

Fabric Analytics Concepts Notes

The screenshot shows the Microsoft Fabric Analytics interface, specifically the Visual query editor. At the top, there's a toolbar with various icons: Manage columns, Reduce rows, Sort, Transform, Combine, Save as view, View SQL, and settings. Below the toolbar, there are two tables listed: FactRevenue and DimDate. A context menu is open over the FactRevenue table, with the 'Transform' option selected. The menu includes options like 'Transform text column - Extract', 'Add column', and 'Save as table'. On the right side of the interface, there's a preview of the data with columns labeled 'Units_Sold' and 'Revenue'.

Units_Sold	Revenue
1	1196976
1	6083151
1	3450696
1	2221530
1	2221530
1	2221530
1	2221530
1	5440414
1	4988693

Fabric Analytics Concepts Notes



Question

The following SQL script creates the result shown below. What is {FUNCTION}?

```
SELECT
    Branch_ID,
    , Col1
    , {FUNCTION}
From Sales
```

A) LAG(Col1, 2) OVER(ORDER BY Col1) as Col2

B) LEAD(Col1, 2) OVER(ORDER BY Col1) as Col2

C) LAG(Col1, -2) OVER(ORDER BY Col1) as Col2

D) LEAD(Col1, -2) OVER(ORDER BY Col1) as Col2

Branch_ID	Col1	Col2
BR0051	1	NULL
BR0161	2	NULL
BR0161	2	1
BR0121	3	2
BR0071	6	2

It is lag func because col 2 is 2 rows behind col1. NULL at top → Lag

Lead looks ahead in time rather than back in time

Offset is always +ve

Fabric Analytics Concepts Notes

You have the following query analyzing Sales data for various products:

```
SELECT s.ProductName, d.Year, SUM(s.SalesAmount)
FROM Sales s
LEFT JOIN DimDate d ON s.DateKey = d.DateKey
```

Your goal is to analyze the Sales data by ProductName and Year, but only for products that have a yearly Sales Amount of more than \$50,000.

How would you complete the query?

You are trying to inspect a join between two tables to spot referential integrity violations.

Which of the following T-SQL join types would be easiest to identify keys on both sides of the join that do not have a match on the other side of the join?

Cross join – tells all different combinations

Referential integrity is a database concept that ensures relationships between tables remain consistent and valid. It guarantees that any foreign key value in a child table must match a primary key value in the parent table (or be NULL).

Key Rules of Referential Integrity

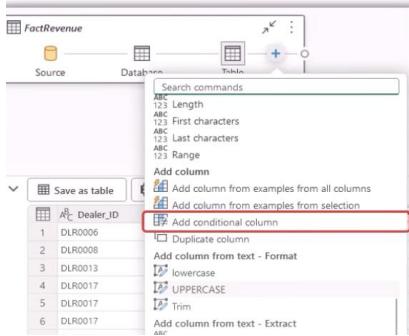
1. **Primary Key (Parent Table):**
 - Must be unique and non-null (e.g., customer_id in a Customers table).
2. **Foreign Key (Child Table):**
 - Must reference an existing primary key or be NULL (e.g., customer_id in an Orders table).
3. **Actions Enforced by DBMS:**
 - **On Delete:** What happens to child rows if a parent row is deleted?
 - CASCADE: Delete child rows automatically.
 - SET NULL: Set foreign key to NULL.
 - RESTRICT: Block deletion if child rows exist.
 - **On Update:** Automatically update foreign keys if the parent key changes.

- A) GROUP BY p.ProductKey, d.DateKey
HAVING SUM(s.SalesAmount) > 50000
- B) **GROUP BY s.ProductName, d.Year**
HAVING SUM(s.SalesAmount) > 50000
- C) WHERE s.SalesAmount > 50000
GROUP BY p.ProductKey, d.DateKey
- D) WHERE s.SalesAmount > 50000
GROUP BY p.ProductName, d.Year

- A) LEFT JOIN
- B) RIGHT JOIN
- C) INNER JOIN
- D) FULL OUTER JOIN**
- E) CROSS JOIN

Fabric Analytics Concepts Notes

You are using the T-SQL visual query editor and your goal is to add a new column to your dataset called 'Salary Bins'. Which functionality could you use to add this new column?



A) Add conditional column

B) Add columns from examples from all columns

C) Add column from examples from selection

D) Duplicate column

E) Add column from text

Reference

https://www.youtube.com/watch?v=Bjk93hi21QM&ab_channel=LearnMicrosoftFabricwithWill