

MICROSOFT FABRIC 75+ Q&A

BY - SHUBHAM WADEKAR

Copyright © Shubham Wadekar. All rights reserved.

This material is for personal use only. No part may be copied, shared, resold, or published without written permission. Unauthorized distribution is strictly prohibited and may result in action.

For permissions or licensing, contact: [shubham.p.wadekar@gmail.com]

Disclaimer: This content is for educational purposes. Accuracy is attempted but not guaranteed. No job outcome is promised.

1. What is Microsoft Fabric, and how does it relate to Azure Data Engineering?

Microsoft Fabric is a new all-in-one analytics platform provided by Microsoft. It combines multiple data services into one unified solution, including data integration, data engineering, data warehousing, data science, real-time analytics, and business intelligence. Everything runs on top of OneLake, which is a single data lake storage system for the entire platform.

From an Azure Data Engineering perspective, Microsoft Fabric changes the way we handle data pipelines and transformations. Earlier, we used services like Azure Data Factory, Azure Synapse, and Azure Data Lake separately. Now with Microsoft Fabric, we can use one platform to do all of this in a more seamless way. As a data engineer, I can build data pipelines, write data transformation code using notebooks or dataflows, and manage the entire data lifecycle inside Fabric, which helps reduce complexity and improves collaboration with other teams like data scientists or BI developers.

2. What are the main workloads available in Microsoft Fabric?

Microsoft Fabric provides several workloads, which means different types of tools and environments for different types of tasks. The main workloads include:

- **Data Engineering:** This is used for big data processing. It allows us to use Spark notebooks to clean, transform, and process large datasets.
- **Data Factory:** This is for building and scheduling data pipelines using Dataflows and activities similar to Azure Data Factory.
- **Data Science:** This helps data scientists build and train machine learning models using notebooks and AI tools.
- **Real-Time Analytics:** This is for analyzing streaming data or near real-time data from sources like IoT devices or event hubs.
- **Data Warehouse:** This provides a scalable and high-performance SQL-based warehouse engine to run analytical queries on large volumes of structured data.
- **Power BI:** This is integrated into Fabric and used for creating reports and dashboards to visualize the data.
- **Synapse Data Engineering and Synapse Data Warehousing** are now part of the Fabric workloads but improved and more integrated.

These workloads are all connected and work with the same data stored in OneLake, which makes it easier for teams to collaborate and access consistent data.

3. What are the main components included in the Microsoft Fabric unified platform?

Microsoft Fabric includes several important components that make it a complete and unified platform:

- **OneLake:** This is the central data lake storage system in Fabric. All data from different workloads is stored in OneLake so that it can be shared and reused easily across the platform.
- **Lakehouse:** This combines features of both data lakes and data warehouses. It allows us to store data in open formats and use SQL or Spark for analysis and transformation.
- **Notebooks:** These are used mainly in the Data Engineering and Data Science workloads. They allow users to write and run code using languages like Python, SQL, and Spark.
- **Pipelines and Dataflows:** These are used to move and transform data. Pipelines help schedule tasks and connect different data sources, while Dataflows allow us to clean and shape the data using low-code or no-code approaches.
- **Power BI:** This is used for data visualization. Since it's part of Fabric, we can easily build reports directly from the data in OneLake.
- **Shortcuts and Mirroring:** Shortcuts let us link external data sources into OneLake without copying the data, and mirroring allows syncing of external databases like Cosmos DB or SQL to Fabric.

These components work together in a single platform so we don't have to switch between multiple Azure services. Everything is integrated, which makes data engineering tasks smoother and more efficient.

4. How does Fabric unify data engineering, data science, real-time analytics, and BI within one environment?

Microsoft Fabric brings all the major data activities into one platform by combining different tools that usually worked separately into one unified system. In the past, we used different Azure services for each task. For example, data engineers used Azure Data Factory or Synapse pipelines, data scientists used Azure ML, and Power BI was separate. Now, in Fabric, everything is available in a single workspace that runs on the same storage layer, which is called OneLake.

For data engineering, we have notebooks and pipelines that can be used to transform large data sets using Spark or SQL. For data science, the same notebooks can be used to build and train machine learning models using Python or other languages. Real-time analytics is supported through tools like KQL and event streams, which let us process streaming data in near real-time. For business intelligence, Power BI is deeply integrated and can directly read data from Lakehouse or Warehouse to create dashboards and reports.

Since all these workloads share the same data stored in OneLake and are available in the same environment, teams can collaborate easily without moving data between services. This saves time, reduces complexity, and helps ensure that everyone is using the same, updated version of the data.

5. What is a semantic model in Fabric, and how does it relate to Lakehouse and Power BI?

A semantic model in Fabric is a layer that helps us define the business meaning of data. It includes tables, relationships, measures, and hierarchies that make the data easier to understand for reporting and analysis. It's like a middle layer between the raw data and the visual reports.

In Fabric, semantic models are automatically created when we build Power BI reports on top of Lakehouse or Warehouse. For example, if I have cleaned data stored in a Lakehouse table, I can use that table to create a Power BI report. Behind the scenes, Fabric creates a semantic model that defines how those tables connect, how calculations are done, and how data is grouped.

This model helps business users understand the data without needing to know the technical details. They can use the model to slice and filter data easily in Power BI. The key point is that the semantic model uses data from Lakehouse or Warehouse and makes it ready for visualization and self-service analytics in Power BI.

6. Can you explain how Delta Parquet format is used internally by Fabric Lakehouse and why it's beneficial?

Delta Parquet is a file format used inside Fabric Lakehouse to store data in an efficient and reliable way. It combines the features of Parquet files, which are columnar and compressed, with extra features from the Delta Lake format like versioning, schema enforcement, and transaction logs.

When data is stored in the Lakehouse, it is saved in Delta Parquet format. This means we get the fast read performance of Parquet, and also the ability to track changes in the data over time. Delta format adds a transaction log that keeps a record of all changes made to the data, like inserts, updates, or deletes. This helps with data consistency and makes it possible to do features like time travel, where we can query the data as it existed at an earlier point.

Another benefit is that this format allows multiple users or processes to read and write data at the same time without causing data corruption. For data engineers, this means we can build reliable pipelines and workflows that support large-scale data processing.

So overall, using Delta Parquet in Fabric Lakehouse gives us better performance, reliability, and flexibility when working with large data sets.

7. What is the difference between using a Warehouse's SQL endpoint vs querying Delta tables in a Lakehouse?

In Microsoft Fabric, both Warehouse and Lakehouse allow us to store and query data, but they serve slightly different purposes and have different performance and interface characteristics.

A Warehouse is optimized for structured data and is fully based on a SQL engine. When we use a Warehouse's SQL endpoint, we are connecting to a traditional SQL interface, similar to what we use in databases. It supports fast, complex analytical queries using T-SQL. Warehouses are ideal for business users and analysts who prefer using SQL and want to build reports using structured and well-defined tables. The performance is usually faster for pure SQL workloads because it is built for this purpose.

A Lakehouse, on the other hand, is designed to handle both structured and unstructured data. Data is stored as Delta Parquet files in OneLake. We can query this data using notebooks (with Spark or SQL), and also with Power BI or Lakehouse SQL analytics endpoint. Lakehouses are more flexible because they support both code-first and low-code approaches. However, the SQL performance might not always match a Warehouse, especially for very complex joins or aggregations.

So, the main difference is that Warehouses are built for high-performance SQL analytics with a traditional database feel, while Lakehouses give more flexibility for big data and support multiple engines including Spark and SQL.

8. How can you share data securely across Fabric workspaces without physically copying it?

In Microsoft Fabric, we can share data securely across workspaces using a feature called shortcuts. Shortcuts allow us to link a folder or file from one workspace's Lakehouse or OneLake location to another workspace, without making a physical copy of the data. This is useful because it avoids duplication and ensures that everyone is working with the same, up-to-date version of the data.

When we create a shortcut, it looks like the data exists in the new workspace, but it's actually pointing to the original location. Access to the data is controlled by permissions, so only users with the right access can view or query the data. This makes sharing easy and secure at the same time.

Using shortcuts is especially helpful when multiple teams need to use the same data for different purposes, like engineering, science, and reporting. Since the data stays in one place, it also helps with governance and avoids confusion caused by having multiple versions of the same data.

9. What is the role of Lakehouse tables vs files, and how do they differ in terms of querying and governance?

In a Lakehouse, we can store data either as tables or as raw files. Both have their use cases, but they serve different purposes and behave differently in terms of querying and governance.

Lakehouse tables are structured and are usually stored in Delta Parquet format. These tables are registered in the Lakehouse metadata, which means they can be queried using SQL. Tables support schema, data types, and transactions, which make them suitable for reporting, transformations, and joining with other tables. They are easier to govern because they are managed like database tables. Features like row-level security and lineage tracking can be applied more effectively to tables.

Files in a Lakehouse are raw data files stored in OneLake folders. They can be in any format like CSV, JSON, or Parquet. These files are more flexible and are often used for data ingestion, staging, or temporary processing. However, they are harder to govern because they don't have a defined schema and are not registered in the metadata. To query them, we usually need to define the schema manually or use external tools like notebooks.

So in summary, Lakehouse tables are used for structured, query-ready data with strong governance, while files are used for raw or semi-structured data that needs further processing. Both play important roles in the data lifecycle inside the Lakehouse.

10. What are the key limitations of Microsoft Fabric currently in preview or GA that a data engineer should know?

As of now, Microsoft Fabric is still evolving, and there are a few limitations that a data engineer should be aware of:

- Some features are still in preview, which means they are not fully stable for production. For example, features like mirroring or certain connectors may not support all data sources yet.
- Real-time analytics is improving, but not all streaming connectors or complex processing features are available like in Azure Stream Analytics.
- There is limited support for custom compute sizing or fine-tuning of Spark clusters. Unlike Azure Databricks, we cannot control the cluster configuration in detail.
- Data pipeline features are still catching up with Azure Data Factory. For example, not all activities and triggers are available yet, and complex control flow logic can be limited.
- Integration with Azure DevOps or GitHub for version control and CI/CD is present but not as mature or flexible as in Synapse or Azure Data Factory.
- Access control and advanced governance features are improving, but some features like row-level security across all workloads or unified policy enforcement may not be available yet.
- Exporting or archiving large volumes of data from Fabric to external systems is possible but not as straightforward or automated in some cases.

It's important for data engineers to stay updated as the platform is being developed actively, and many of these gaps are expected to be filled soon.

11. How does Fabric support multilingual development environments (e.g., PySpark, SQL, DAX)?

Fabric supports multiple languages for different workloads, which makes it easier for different types of users to work in the same environment.

- For data engineering and data science, we can use notebooks that support languages like PySpark, Scala, SQL, and also plain Python. These notebooks run on Spark behind the scenes, which means we can do big data processing using PySpark, or use native Spark SQL when needed.
- For reporting and analysis, especially in Power BI, Fabric supports DAX. DAX is used for creating calculated measures, KPIs, and custom aggregations in the semantic model. It helps business users and analysts build complex logic inside reports.
- In the Data Warehouse workload, we use T-SQL for writing queries, creating tables, and running transformations. This is helpful for database developers and BI professionals who are used to traditional SQL.
- Even inside pipelines, we can use SQL queries or map data flows with expressions, which supports more low-code or SQL-based development.

So Fabric provides flexibility for different teams to use their preferred language in the same platform. A data engineer can use PySpark for processing, an analyst can use SQL for queries, and a BI user can use DAX for reporting, all within the same workspace.

12. What options are available to back up or export data stored in a Fabric Lakehouse?

Backing up or exporting data from a Fabric Lakehouse involves a few approaches since the platform uses OneLake for storage and data is stored in Delta Parquet format.

- One option is to use the shortcut feature to link data to another workspace or Lakehouse. This doesn't create a copy, but it allows shared access to the same data, which can be helpful for redundancy.
- Another method is to use notebooks or pipelines to write the data from Lakehouse tables or files to external storage accounts like Azure Data Lake Storage Gen2 or Blob Storage. We can use Spark or Python code to copy the data.
- For exporting specific tables, we can use the export functionality in notebooks or use Dataflows Gen2 to move data to external destinations.
- There's also the option to sync or mirror data from external sources into Fabric, and some of these sources might support two-way sync, so the original system acts as the backup.
- Currently, automatic backup or versioning inside Fabric is limited, but since the Lakehouse uses Delta format, it does keep a transaction log. This means we can technically query previous versions of a table (time travel), which can help recover from some changes.

It's important to note that full backup and disaster recovery features are still being developed, so for now, exporting data to external storage is the safest backup approach.

13. What is OneLake in Microsoft Fabric, and how does it differ from traditional data lakes like Azure Data Lake Storage Gen2?

OneLake is the built-in, unified data lake storage that powers Microsoft Fabric. It is designed to be the single place where all data is stored and accessed across different workloads like data engineering, data science, real-time analytics, data warehouse, and Power BI. It is similar in concept to Azure Data Lake Storage Gen2, but it has some key differences.

In traditional Azure Data Lake Storage Gen2, we need to manage storage accounts, set up access control, define folder structures, and connect each service separately. Each service like Synapse, Data Factory, or Power BI had to integrate with the storage manually.

In OneLake, this storage layer is already built into Fabric and doesn't require extra configuration. All the workloads in Fabric use OneLake by default, and the data stored is automatically organized by workspaces. OneLake also uses Delta Parquet as its default format, so data is query-ready for both Spark and SQL engines. Another key difference is that OneLake supports shortcuts, which allow data sharing without physical movement. That feature is not available in regular ADLS Gen2.

So, while both are cloud-based data lakes, OneLake is more tightly integrated, automated, and unified with Fabric services, making it easier to use and manage for analytics projects.

14. How does OneLake provide a unified storage layer across multiple data domains and services in Microsoft Fabric?

OneLake acts as a central storage layer that connects all Fabric workloads like Lakehouse, Warehouse, Real-Time Analytics, Data Science, and Power BI. It is organized around the concept of workspaces, and each workspace automatically gets its own folder structure in OneLake.

All data saved through any Fabric service whether it's a pipeline, notebook, or Power BI is automatically stored in OneLake. Since OneLake supports open formats like Delta Parquet, this makes it easy for all workloads to access and use the same data. For example, a data engineer can write data to OneLake using a notebook, and then a Power BI user can directly build a report from the same dataset without needing to copy or transform it again.

OneLake supports multi-tenant architecture and keeps data organized by domains through its workspace structure. It enforces consistent security and access policies, which means different teams can work on the same platform but only access the data they are allowed to see.

This unified layer simplifies data management, avoids silos, and improves collaboration between data engineering, analytics, and business teams.

15. Can you explain how the shortcut feature works in OneLake and how it helps in data virtualization and reducing data duplication?

The shortcut feature in OneLake allows us to create a virtual link from one location in OneLake to another location or even to an external data source like Azure Data Lake or Amazon S3. Instead of copying the data, the shortcut just points to the original data location and lets users access it as if it existed in their own workspace.

For example, if a data engineer in one team has cleaned and stored customer data in a Lakehouse, another team can create a shortcut to that folder or table in their own workspace. Now they can query or use the data without making a copy. This saves storage space and ensures everyone works with the same, updated version of the data.

Shortcuts are important for data virtualization because they allow data to be shared across teams and domains without physical movement. They also make it easier to enforce governance and control access from one place, since the data remains in its original location.

This feature is very useful in large organizations where multiple teams need access to the same core datasets, and it helps prevent problems caused by duplicated or outdated copies of data.

16. What role do workspaces and domains play in organizing data within OneLake, and how do they relate to security and access control?

Workspaces and domains are used to organize data and control access inside Microsoft Fabric's OneLake.

A workspace is a container for related resources like Lakehouses, Warehouses, Pipelines, Notebooks, and Power BI reports. Each workspace has its own dedicated folder structure in OneLake, and all the data created or stored in that workspace is automatically organized under it. This makes it easy to manage and locate data, especially in large projects.

Domains are a higher-level concept that represent business areas like sales, finance, marketing, or operations. Domains group workspaces that belong to the same business area. This helps organizations manage their data in a more logical and business-aligned way.

In terms of security, workspaces act as the boundary for access control. Only users who have been given permission to a specific workspace can access the data and objects inside it. Permissions can be given at the workspace level or on specific items like a Lakehouse table or a report. This model provides fine-grained control, and admins can easily manage who can read, write, or modify the data.

So, workspaces and domains help with both organizing the data clearly and enforcing proper access control to make sure that only the right people can access the right data.

17. How does OneLake integrate with different engines like Synapse, Data Engineering, Power BI, and Real-Time Analytics without requiring data movement?

OneLake integrates all the different analytical engines in Microsoft Fabric by acting as the single storage location for all data, regardless of which service is using it. This is possible because OneLake uses open data formats like Delta Parquet, which can be read and written by multiple engines.

For example, when a data engineer stores data in a Lakehouse using Spark notebooks, the same data is immediately available for Power BI to build reports, for the SQL analytics engine to run queries, or for Real-Time Analytics to process live streams. There's no need to move or copy the data between services, because all of them are connected directly to OneLake.

Each engine in Fabric is designed to work natively with OneLake. This means Synapse SQL can query the data in Delta tables, Power BI can build visuals from the same source, and notebooks can use PySpark or SQL on the same tables. Everything happens within the same storage layer, and the metadata is shared as well.

This tight integration saves time, avoids duplication, reduces errors, and improves collaboration across different roles like data engineers, data scientists, analysts, and business users. It also helps maintain data consistency because everyone works with the same source of truth in OneLake.

18. What is a Lakehouse in Microsoft Fabric, and how does it unify the features of both data lakes and data warehouses?

A Lakehouse in Microsoft Fabric is a storage and compute environment that combines the best parts of a data lake and a data warehouse. It allows us to store large amounts of structured, semi-structured, or unstructured data using open formats like Delta Parquet, while also letting us run SQL queries and BI reports just like in a data warehouse.

From a data lake perspective, a Lakehouse allows us to store data in raw file formats, supports big data processing using Spark, and can handle massive volumes of data in a cost-effective way. From a data warehouse perspective, the Lakehouse supports tables, schemas, and transactional processing using the Delta format, which allows for reliable querying and data consistency.

This unification helps data engineers use a single environment for both flexible storage and structured querying. Data can be ingested, processed, and analyzed all within the Lakehouse, without needing to move it between different systems. Also, it integrates smoothly with Power BI, so reports can be created directly from Lakehouse tables, making it useful for end-to-end analytics.

19. How does the Lakehouse architecture in Fabric differ from a traditional data lake (like ADLS Gen2) and a conventional data warehouse (like Synapse Dedicated SQL Pool)?

The Lakehouse in Fabric is different from both traditional data lakes and traditional data warehouses in the following ways:

- Compared to a traditional data lake like ADLS Gen2, a Fabric Lakehouse is more structured and comes with built-in compute, metadata management, and SQL querying capabilities. In ADLS Gen2, we only get raw storage and have to connect it manually to compute engines like Synapse or Databricks. In contrast, the Fabric Lakehouse already includes a query engine, table definitions, and tight integration with notebooks, pipelines, and Power BI.
- Compared to a traditional data warehouse like Synapse Dedicated SQL Pool, the Lakehouse supports open file formats and big data processing using Spark. In Synapse SQL Pool, data must be structured and loaded into SQL tables, and it uses proprietary storage formats. The Lakehouse allows for more flexibility by storing data as Delta Parquet files, which can be used by both Spark and SQL engines.
- The Lakehouse also supports features like time travel, schema evolution, and ACID transactions through the Delta format, which are not typically available in traditional data lakes. At the same time, it avoids the fixed compute and storage costs of a dedicated SQL pool.

In summary, the Lakehouse is more flexible than a warehouse and more structured and usable than a raw data lake, giving us the best of both in one place.

20. What are the default folders created when you provision a Lakehouse in Microsoft Fabric, and what is the purpose of each?

When we create a Lakehouse in Microsoft Fabric, it automatically creates a few default folders in OneLake to help organize the data:

- **Tables:** This folder contains all the structured tables in the Lakehouse. Each table here is stored in Delta Parquet format and can be queried using SQL. These are managed tables with schema, and they are registered in the Lakehouse metadata so tools like Power BI or notebooks can easily access them.
- **Files:** This folder is used for storing unstructured or semi-structured data like CSV, JSON, Excel, images, or raw Parquet files. These files are not registered as tables by default but can be read using Spark or can be transformed into tables later. This folder is often used as a staging area before transforming data.
- **Untracked (optional):** Sometimes, if users upload data through external sources or APIs, those files might land in untracked folders. These are not automatically tracked by the system metadata but still exist in the OneLake structure.

These folders help separate raw data from structured data and make it easier to manage data ingestion, transformation, and querying inside the Lakehouse.

21. How does the Delta Lake format enhance the functionality of a Lakehouse in Fabric?

The Delta Lake format is a key technology that enhances how a Lakehouse works in Microsoft Fabric. It combines the advantages of Parquet file storage with additional features like transaction logs, version control, and schema management. These improvements make data more reliable, easier to manage, and better suited for both big data processing and analytics.

Here are the main ways Delta Lake format improves the Lakehouse:

- **ACID Transactions:** Delta Lake ensures that data operations like inserts, updates, and deletes are reliable and consistent. This means even if multiple users or jobs are accessing the same table at the same time, the data will stay correct and won't get corrupted.
- **Time Travel:** Because Delta Lake keeps a log of all changes made to the table, we can go back and query previous versions of the data. This is useful if we want to audit changes or recover from accidental deletions.
- **Schema Enforcement and Evolution:** Delta Lake allows us to define a schema for each table, and it can either enforce this schema strictly or allow it to evolve over time. This helps ensure data quality and makes it easier to track changes in the data structure.
- **Performance Optimization:** Delta format supports features like file compaction and indexing through data skipping, which makes queries faster even on large datasets.

Because of all these features, the Lakehouse becomes a reliable, high-performing environment for both batch processing and analytics, while still supporting flexible and open data formats.

22. How can various Fabric engines like Power BI, Data Engineering (Spark), and SQL Analytics interact with the same Lakehouse without copying the data?

In Microsoft Fabric, all engines interact with the same Lakehouse storage by directly reading from and writing to OneLake using open file formats like Delta Parquet. This shared access is possible because the Lakehouse acts as a central and open storage layer, and all the engines in Fabric are built to understand and work with the same format.

Here's how different engines use the same data without copying:

- **Data Engineering (Spark):** Spark notebooks and jobs can read and write to Delta tables stored in the Lakehouse using PySpark or Spark SQL. Spark understands the Delta format and can handle large-scale transformations efficiently.
- **SQL Analytics:** The SQL endpoint in the Lakehouse can directly query the Delta tables as if they were traditional SQL tables. There is no need to move or import the data. SQL queries are run directly on the stored Delta files.
- **Power BI:** Since Delta tables in the Lakehouse are registered and structured, Power BI can connect to them directly to create reports. The data model can be built using these tables without any extraction or duplication.

Because all these services are built into Fabric and connected to OneLake, they can access the same data at the same time, in-place. This reduces the need for data movement, avoids duplication, and ensures that everyone is working with the most current version of the data. It also simplifies governance because all access is controlled through the workspace and Lakehouse permissions.

23. What are Dataflows Gen2 in Microsoft Fabric, and how are they used to perform data transformation and ingestion tasks?

Dataflows Gen2 are low-code or no-code tools in Microsoft Fabric that allow users to perform data ingestion and transformation using a visual interface. They are built using Power Query, which lets users connect to various data sources, apply transformations like filtering, merging, and pivoting, and then load the cleaned data into destinations like Lakehouse or Warehouse.

Dataflows Gen2 are especially useful for data engineers or analysts who want to perform ETL tasks without writing code. We can use them to pull data from databases, files, APIs, or cloud sources, transform that data using step-by-step logic, and load the result into structured storage. Each step in the Dataflow can be reviewed, tested, and reused, which makes it easy to debug and manage.

These dataflows run on the same platform as the rest of Fabric, so they take advantage of the shared compute and storage systems. They can be scheduled to run at regular intervals, allowing automated data refresh and ingestion workflows.

24. How do Dataflows Gen2 improve upon classic Power BI Dataflows in terms of architecture, performance, and integration?

Dataflows Gen2 offer several improvements over the older Power BI Dataflows, especially in how they are built, how fast they run, and how well they connect with other parts of Microsoft Fabric:

- **Architecture:** Classic Power BI Dataflows store data in Power BI's internal storage. In contrast, Dataflows Gen2 are fully integrated with OneLake and can write directly to Lakehouse or Warehouse. This makes the data more accessible and usable by other Fabric services like notebooks, pipelines, and Power BI.
- **Performance:** Dataflows Gen2 use the compute power of the Fabric platform, which can scale better and handle larger datasets more efficiently than Power BI dataflows. They also support parallel processing and can handle more complex transformation logic.
- **Integration:** Dataflows Gen2 are part of the Fabric environment, so they can easily connect with other tools like Spark notebooks, SQL analytics, and Power BI. This makes it easier to build end-to-end data pipelines in one platform. The output of a Dataflow Gen2 can be used directly by a Power BI report, or can feed into a notebook or another Dataflow.

Overall, Dataflows Gen2 are more flexible, faster, and better integrated, making them a powerful tool for modern data transformation tasks in Fabric.

25. How do Dataflows Gen2 integrate with the Lakehouse in Fabric, and what are the advantages of storing data directly in Delta Lake format?

Dataflows Gen2 can write their output directly into a Lakehouse in Microsoft Fabric. When you select a Lakehouse as the destination, the data is automatically stored in Delta Lake format inside the Lakehouse's "Tables" folder. This integration allows seamless collaboration between low-code users (using Dataflows) and code-first users (using notebooks or SQL).

The advantages of storing data in Delta Lake format are:

- **Open Format:** Delta Lake is based on Parquet and is an open, industry-standard format. This allows the data to be read by Spark, SQL, Power BI, and other tools without conversion.
- **Query Ready:** Once stored in Delta format, the data can be queried directly using SQL in the Lakehouse or used in Spark notebooks. There's no need for additional processing or registration.
- **Reliable and Consistent:** Delta supports ACID transactions and keeps a log of all changes, which makes the data reliable and easy to track. This is useful for building trustworthy data pipelines.
- **Version Control and Time Travel:** Delta format allows querying historical versions of the data, which can be helpful for auditing or debugging.

By writing directly to Delta Lake through Dataflows Gen2, users get the benefits of a structured, high-performance data storage layer without writing any code. This improves collaboration and ensures the data is immediately available for analytics and reporting.

26. What are the key components of a Dataflow Gen2 pipeline, and how can you orchestrate data refresh and error handling within it?

A Dataflow Gen2 pipeline is built using Power Query and includes several key components that help in designing, running, and managing data ingestion and transformation tasks. The main components are:

- **Source connectors:** These let us connect to various data sources such as SQL databases, Excel files, Azure services, or REST APIs.
- **Transformation steps:** Using the Power Query editor, we can apply transformations like filter, join, merge, pivot, remove columns, and change data types. Each step is applied in sequence and shown in the query editor pane.
- **Destination settings:** After transforming the data, we select a destination usually a Fabric Lakehouse or Warehouse. When saving to a Lakehouse, the data is written in Delta format to the "Tables" folder.
- **Schedule refresh:** This allows us to automate the refresh of data at defined intervals, like daily, hourly, or weekly.
- **Error controls:** Within Power Query, we can add steps to handle errors, like replacing error values, filtering out invalid rows, or using conditional logic to clean data. In the transformation editor, error rows are highlighted so they can be handled before the load step.

For orchestration, we can configure refresh schedules and even chain Dataflows together. For example, we can have one Dataflow prepare data and another one consume that output. While Dataflows Gen2 don't have complex control flow like pipelines, we can monitor their status, check for errors in the refresh history, and re-run failed loads.

27. How do you control data lineage, incremental refresh, and schema evolution in Dataflows Gen2 when working with large datasets?

When working with large datasets, it's important to manage how data is loaded, changed, and understood over time. Dataflows Gen2 in Microsoft Fabric offer several features to help with this:

- **Data lineage:** Fabric automatically tracks lineage for Dataflows Gen2. This means we can see where the data came from and where it is used across other Fabric items like Lakehouses, Power BI reports, and other Dataflows. This helps with impact analysis, auditing, and understanding data dependencies.
- **Incremental refresh:** For large datasets, loading the entire data every time is inefficient. In Dataflows Gen2, we can enable incremental refresh by specifying a date column and configuring the range of data to be loaded (such as "load only data from the last 7 days"). This reduces load time and improves performance. It also helps reduce cost by only processing new or changed records.
- **Schema evolution:** If the source schema changes, like a new column being added or renamed, Dataflows Gen2 can detect this. In the Power Query editor, we can handle these changes by refreshing the preview, updating transformation steps, and adding logic to deal with new or missing columns. Since we're working with a visual interface, any errors caused by schema changes are usually flagged during design or refresh, so we can fix them before the next run.

By combining these features, Dataflows Gen2 provides a reliable way to manage complex, growing datasets while keeping performance high and reducing manual rework.

28. How do Notebooks in Microsoft Fabric differ from those in Azure Synapse Analytics in terms of integration, use cases, and performance?

Notebooks in Microsoft Fabric and Azure Synapse Analytics both provide a way to write and run code for data engineering, machine learning, and analysis, but there are several differences between them in terms of integration, use cases, and performance.

- **Integration:** Fabric Notebooks are deeply integrated with the entire Fabric platform, especially with Lakehouses, Power BI, Dataflows, and other workloads. This makes it easier to read from or write to a Lakehouse, visualize data, or pass results to other tools within the same workspace. In Synapse, integration is good but usually requires more configuration to connect with services like Power BI or ADLS.
- **Use Cases:** Fabric Notebooks are designed for unified use across data engineering, data science, and reporting. Since Fabric uses OneLake as the common storage, it's easy to share data across teams. In Synapse, Notebooks are mostly focused on data preparation or transformation tasks with Spark pools and are not as tightly linked to Power BI.
- **Performance:** In Fabric, Spark clusters are auto-managed and run in the background, so users don't need to manually configure or manage them. The environment is simpler to use, but offers less control over cluster size or tuning. In Synapse, users have more control over the Spark pool configuration, which may help for very large jobs but requires more effort.

Overall, Fabric Notebooks are easier to use in an end-to-end analytics environment and provide a more seamless experience for collaborative analytics and reporting tasks.

29. Which programming languages are supported in Fabric Notebooks, and how can multi-language cells be leveraged within a single notebook?

Fabric Notebooks currently support several common languages used in data engineering and data science:

- PySpark (Python with Spark)
- Spark SQL
- Scala (for advanced Spark users)
- Markdown (for notes, comments, and headings)

Each cell in a Fabric Notebook can be set to a different language using a language selector or magic commands like %python, %sql, or %md. This means we can use multiple languages in the same notebook depending on the task.

For example, we can:

- Use PySpark for data transformation
- Switch to SQL to quickly preview data with a simple query
- Add a markdown cell to document what we are doing
- Use Scala if we need to access some advanced Spark features not available in PySpark

This flexibility helps different users work in their preferred language and keeps the entire workflow in one place.

30. How can you read from and write to a Lakehouse using PySpark in Fabric Notebooks? Can you give a simple example?

In Microsoft Fabric Notebooks, we can easily read from and write to Lakehouse tables using PySpark. Since the Lakehouse is already connected to the notebook environment, we don't need to manually set up access or authentication.

Here is a simple example:

Reading from a Lakehouse table:

```
df = spark.read.format("delta").load("Tables/sales_data")  
  
df.show()
```

This reads data from the sales_data table stored in the Lakehouse's Tables folder. The format is delta because Lakehouse tables are stored in Delta Lake format.

Writing to a Lakehouse table:

```
# Transform or create a DataFrame  
  
from pyspark.sql import Row  
  
data = [Row(product="Phone", sales=1000), Row(product="Laptop", sales=2000)]  
  
df_new = spark.createDataFrame(data)  
  
# Write the DataFrame to the Lakehouse as a new table  
  
df_new.write.format("delta").mode("overwrite").save("Tables/product_sales")
```

This saves the new DataFrame as a Delta table named `product_sales` in the Lakehouse.

These paths like `"Tables/..."` are relative to the Lakehouse connected to the notebook, and they point to the actual storage in OneLake. Since the format is Delta, we get all the benefits like ACID transactions and schema enforcement.

31. What are the common use cases for Fabric Notebooks in data engineering pipelines, and how do they integrate with other Fabric components like Lakehouse and Dataflows Gen2?

Fabric Notebooks are widely used in data engineering pipelines because they offer flexibility, code-based control, and deep integration with the Fabric ecosystem. Some common use cases include:

- **Data ingestion:** Reading data from external sources like APIs, Azure storage, databases, or files, and loading it into the Lakehouse.
- **Data transformation and cleansing:** Applying business rules, filtering bad data, handling missing values, and reshaping datasets using PySpark or Spark SQL.
- **Data validation and quality checks:** Writing logic to check for duplicates, outliers, nulls, or schema mismatches before writing data downstream.
- **Building reusable logic:** Defining parameterized logic, reusable functions, or machine learning steps for inclusion in larger data workflows.
- **Scheduling as part of pipelines:** Notebooks can be triggered from Fabric pipelines to run as part of a scheduled or event-driven workflow.

Fabric Notebooks integrate smoothly with:

- **Lakehouse:** Notebooks are directly connected to a specific Lakehouse, so we can read from or write to Delta tables using simple PySpark commands without any extra setup.
- **Dataflows Gen2:** We can use Dataflows Gen2 for low-code ingestion and transformation, then use notebooks for more advanced logic or custom operations that are not possible in Power Query.
- **Power BI:** The tables created or updated by notebooks in the Lakehouse can be consumed directly in Power BI reports since the integration is built-in.

This integration allows notebooks to act as a powerful bridge between low-code data preparation and advanced code-based processing, all within the same workspace.

32. How do Fabric Notebooks handle session state, Spark execution, and resource scaling, especially compared to Databricks or Synapse Spark pools?

Fabric Notebooks run on managed Spark environments, but the way sessions and compute resources are handled is simplified compared to Databricks or Synapse.

- **Session state:** When we open a Fabric Notebook, it automatically creates a Spark session in the background. The session stays active while we work and automatically shuts down after a period of inactivity. The session maintains variables, temporary tables, and libraries during the active period. However, once the session times out, any session-specific state is lost, so we should save outputs to the Lakehouse if they need to be reused.
- **Spark execution:** Fabric manages Spark execution behind the scenes. We do not need to set up or choose Spark pool sizes like we do in Synapse or Databricks. Jobs are distributed and executed on managed Spark compute with reasonable defaults for most common workloads.
- **Resource scaling:** In Databricks or Synapse, we can configure the number of nodes, memory, and compute size manually to optimize performance or handle very large datasets. In Fabric, this control is abstracted away. It's easier to use, especially for beginners or non-experts, but advanced users have less tuning flexibility. As of now, resource scaling in Fabric is automatic and focused on simplicity rather than full customization.

In summary, Fabric Notebooks are designed for ease of use and fast integration, with a managed Spark experience that works well for most workloads. For more complex or large-scale Spark jobs that need detailed cluster control, platforms like Databricks or Synapse Spark might still offer more tuning options.

33. What is Eventstream in Microsoft Fabric, and how does it help with real-time data ingestion and processing?

Eventstream in Microsoft Fabric is a component that allows you to collect, process, and route real-time data from various streaming sources into destinations like Lakehouse, Warehouse, or Real-Time Analytics. It provides a simple, visual, drag-and-drop interface to build streaming pipelines without writing complex code.

With Eventstream, we can ingest data from sources like Azure Event Hubs, IoT Hub, Kafka, or custom APIs. Once the data arrives, we can apply real-time transformations, filters, and enrichments as it flows through the stream. The processed data can then be sent to destinations such as a Lakehouse for historical storage, a KQL database for real-time querying, or Power BI for live dashboards.

This helps data engineers and analysts process and act on data as it is generated, which is important for use cases like monitoring sensors, tracking transactions, or detecting anomalies immediately instead of waiting for batch jobs.

34. How does Synapse Real-Time Analytics (KQL) integrate with Eventstream, and what role does it play in querying streaming data?

Synapse Real-Time Analytics in Microsoft Fabric uses the Kusto Query Language (KQL) to perform fast, interactive queries on streaming or near real-time data. It integrates tightly with Eventstream, allowing us to send incoming events directly into a KQL database or table.

When we configure Eventstream, we can add a KQL table as a sink (or output). This means that as the streaming data flows through the Eventstream pipeline, it can be ingested into the KQL table in real time. Then, we can write KQL queries to analyze the data immediately, detecting patterns, trends, or anomalies as they occur.

KQL is optimized for time-series and telemetry data, making it ideal for dashboards, alerts, and operational monitoring. It supports fast filtering, aggregation, and windowed queries, which are very useful when dealing with continuous event streams.

So, the role of Synapse Real-Time Analytics is to enable fast querying and exploration of live data coming from Eventstream, without needing to store or process the data in batch first.

35. How can you visualize real-time data in Power BI using Microsoft Fabric, and what components are involved in making this end-to-end pipeline work?

To visualize real-time data in Power BI using Microsoft Fabric, we need to build an end-to-end pipeline that connects the data source to Power BI through streaming and analytics components. The typical components involved are:

1. **Eventstream:** This component ingests real-time data from sources like Event Hubs or IoT Hub. It processes and routes the data as it arrives. You can apply light transformations here if needed.
2. **KQL Database (Real-Time Analytics):** The streaming data from Eventstream is written into a KQL table in the Synapse Real-Time Analytics workload. This table stores the live events and supports fast querying using KQL.
3. **Power BI:** In Power BI, we create a report or dashboard that connects directly to the KQL table. Power BI supports DirectQuery for KQL tables, which means every time the dashboard refreshes, it pulls the latest data from the streaming table in real time.

By combining these three components Eventstream for ingestion, KQL for analysis, and Power BI for visualization we can build a live dashboard that updates continuously with real-world events. This setup is ideal for scenarios like monitoring system metrics, viewing live sales data, or tracking IoT sensor updates.

36. How does Eventstream in Fabric handle different input sources (like Event Hubs, IoT Hub, Kafka), and how do you route that data to destinations like KQL DB or Lakehouse?

Eventstream in Microsoft Fabric is designed to connect easily to various real-time data sources and route that data to one or more destinations with minimal setup and no coding. It provides a visual interface where you can add input sources and define the flow of data to output destinations.

Here's how it handles input sources:

- **Event Hubs:** You connect by providing the Event Hub namespace, event hub name, and authentication settings. Once connected, Eventstream starts receiving the real-time events as they are published.

- **IoT Hub:** You can connect to Azure IoT Hub in a similar way. Eventstream can read telemetry or device messages directly from the hub.
- **Kafka:** If you're using Kafka (either on-prem or cloud), you can provide the broker URL, topic name, and credentials. Eventstream can consume messages from the Kafka topic in real time.

Once the data is coming in, you can route it to destinations using the Eventstream designer:

- **KQL Database (Real-Time Analytics):** You can add a KQL table as an output. This allows the streaming data to be ingested directly into the table. You can then query it using KQL for near-instant analysis.
- **Lakehouse:** You can also choose a Lakehouse as a destination. In this case, Eventstream writes the data into a Delta table under the "Tables" folder, which can be used for historical analysis or further processing in notebooks or Power BI.

You can send the same incoming stream to multiple destinations at once. For example, you might send one copy to a KQL table for real-time dashboards and another copy to a Lakehouse for long-term storage.

37. What are some key performance and design considerations when building real-time analytics solutions using Eventstream and KQL in Fabric?

When building real-time analytics solutions using Eventstream and KQL in Fabric, it's important to design for both speed and stability. Here are some key considerations:

- **Schema design:** Define a clear and consistent schema for your incoming events. This helps with smooth ingestion into KQL tables and avoids schema mismatches or dropped data. Flatten nested JSON if possible.
- **Partitioning and batching:** For high-throughput scenarios, make sure your source system (like Event Hub or Kafka) is partitioned properly to handle parallel streams. Eventstream handles events in real time, but source configuration affects throughput.
- **Data volume:** If you are dealing with very high event rates, consider writing only essential fields to the KQL table and sending raw or detailed data to a Lakehouse for deeper analysis later. This helps keep real-time querying fast.
- **Windowing and aggregation:** In KQL, use appropriate windowing functions (like `summarize`, `bin()`, or `take`) to group data over time intervals. This improves performance and reduces the number of rows scanned in each query.
- **Retention policies:** KQL databases allow you to set data retention periods. For real-time dashboards, you might only need the last few hours or days of data, which reduces storage and improves query performance.
- **Monitoring and alerts:** Use KQL queries to detect patterns or thresholds in the data and connect the results to Power BI or alerting systems. For example, alert when error rates spike or sensor readings exceed a limit.
- **Error handling:** Use Eventstream's transformation steps to validate incoming data and drop or redirect invalid records. This keeps your real-time pipeline healthy and prevents bad data from entering your system.

Designing with these points in mind helps ensure that your real-time analytics pipeline is fast, reliable, and scalable, while also being easy to manage and extend as data volume grows.

38. What are Domains in Microsoft Fabric, and how do they help organize data products within a data mesh architecture?

Domains in Microsoft Fabric are logical containers that help group related workspaces, datasets, and other analytics artifacts by business area or functional unit. They are designed to support the principles of data mesh, where each business domain (like sales, finance, operations, or HR) is responsible for its own data products.

In a data mesh architecture, instead of having one centralized data team handle everything, each domain can manage its own data pipelines, Lakehouses, reports, and security rules. This approach promotes ownership, scalability, and faster delivery of analytics solutions.

Domains help organize these assets clearly. For example, the "Marketing" domain might include several workspaces like "Campaign Analytics", "Customer Segmentation", and "Web Traffic Analysis", all grouped under the same domain. This structure makes it easier to manage access, monitor usage, and ensure data governance at the business level.

By using domains, organizations can better align their data strategy with how the business is structured and empower each domain team to manage and deliver their own data products responsibly.

39. How do Domains differ from Workspaces in Microsoft Fabric in terms of structure, purpose, and governance?

Domains and Workspaces serve different purposes in Microsoft Fabric:

- **Structure:** A domain is a higher-level grouping that contains one or more workspaces. A workspace is a more detailed container that holds the actual items like Lakehouses, Warehouses, Notebooks, Pipelines, and Reports.
- **Purpose:** Domains are used to organize data and analytics content by business function or organizational unit. Their goal is to support decentralization and data ownership. Workspaces, on the other hand, are used for collaboration by teams. They provide a shared environment where users can create, edit, and manage specific projects or solutions.
- **Governance:** Domains provide a layer for governance at the business level. You can monitor usage, apply naming standards, and organize workspaces within domains for easier tracking and alignment. However, detailed access control is still handled mostly at the workspace level. Workspaces are where role assignments like Admin, Member, Contributor, or Viewer are applied.

In short, domains help with higher-level business organization and visibility, while workspaces focus on team-level development, collaboration, and security.

40. How is Role-Based Access Control (RBAC) applied at the domain, workspace, and item levels in Microsoft Fabric?

Microsoft Fabric uses Role-Based Access Control (RBAC) at multiple levels to manage permissions and ensure secure access to data and tools:

- **Domain Level:** RBAC at the domain level is currently more about visibility and governance than detailed permissions. Users with domain roles (like domain admins or stewards) can oversee the domain's structure, monitor workspace alignment, and apply naming or tagging policies. However, domain roles don't automatically give access to the content inside its workspaces.
- **Workspace Level:** This is where most access control is applied. Users can be assigned roles like:
 - Admin: Full control, including managing permissions.
 - Member: Can edit, add, and manage items.
 - Contributor: Can add and modify content but with fewer permissions than members.
 - Viewer: Can only view content, like dashboards and reports.

Workspace roles define who can access and manage the Lakehouses, Pipelines, Reports, and other artifacts inside.

- **Item Level:** For more granular control, specific permissions can also be applied at the item level (like a single Lakehouse or Report). This is useful when some users need access to a specific data product but not the entire workspace. You can grant or restrict actions such as read, write, or manage at the item level.

This layered RBAC approach helps organizations enforce the right level of control from broad domain oversight down to fine-grained access on individual assets ensuring that data is secure while still enabling collaboration.

41. What are the best practices for managing security, ownership, and discoverability of data assets across domains and workspaces?

To manage data effectively in Microsoft Fabric, especially when using multiple domains and workspaces, it's important to follow best practices in three key areas: security, ownership, and discoverability.

- **Security:**
 - Always assign roles using the principle of least privilege. Give users only the permissions they need (e.g., Viewer for reports, Contributor for notebooks).
 - Use workspace-level RBAC to control who can access Lakehouses, Pipelines, or Power BI reports.
 - Apply item-level permissions for sensitive data when you want to restrict access within a shared workspace.
 - Protect data using Microsoft Purview integration for labeling, classification, and sensitivity tagging.

- **Ownership:**
 - Assign clear ownership for each workspace and item. Make sure each data product (like a Lakehouse or Report) has a responsible team or user.
 - Use naming conventions to indicate ownership or domain (e.g., “finance_sales_lakehouse” or “hr_payroll_workspace”).
 - Periodically review user access and ownership roles to avoid having outdated or unused access permissions.
- **Discoverability:**
 - Tag datasets and workspaces with business-friendly names and descriptions so others can understand what the data is about.
 - Use domains to group related workspaces, making it easier for users to browse and find what they need.
 - Publish important datasets as certified or promoted so analysts know which sources are trusted and high quality.
 - Enable metadata tracking and lineage features to help users see where data comes from and how it’s used across the system.

By combining strong security, clear ownership, and good discoverability practices, organizations can build a trustworthy and well-governed analytics environment in Fabric.

42. How do Domains support decentralized data ownership, and how can they be used to align Fabric with enterprise data governance strategies?

Domains in Microsoft Fabric are designed to support decentralized data ownership by organizing data products and workspaces around business units or functional areas, such as marketing, finance, or operations.

Here’s how they support this model:

- **Clear boundaries for ownership:** Each domain represents a business area, so data teams within that domain can take full ownership of their data pipelines, Lakehouses, and reports. This allows domain teams to act as data product owners, responsible for the quality, access, and usage of their own data.
- **Separation of responsibilities:** Because domains group related workspaces, each team can manage their data independently without waiting for a central IT team. This improves agility and allows faster delivery of data products.
- **Enterprise-wide visibility:** Even though ownership is decentralized, domains provide a structured way to view how the organization’s data is distributed. This helps governance teams monitor data usage, set policies, and ensure consistency across the platform.
- **Alignment with governance:** Domains can be used alongside tools like Microsoft Purview for classification, policy enforcement, and auditing. Governance teams can set rules at the domain level (such as naming policies or sensitivity labels), while domain teams handle the implementation within their own areas.
- **Support for data mesh principles:** Domains map well to data mesh architecture, where each domain is responsible for producing and maintaining their own data products. This allows Fabric to scale across large enterprises without relying on a central bottleneck.

Using domains in this way helps bridge the gap between business and IT, encourages accountability at the domain level, and supports a federated data governance model where policies are centralized but operations are distributed.

43. What is Data Activator in Microsoft Fabric, and how does it support event-driven data workflows?

Data Activator is a no-code tool in Microsoft Fabric that allows users to create event-driven workflows by monitoring data for specific conditions and triggering actions when those conditions are met. It enables proactive decision-making by automatically reacting to patterns, thresholds, or anomalies in your data without writing any custom logic or building a full application.

Instead of constantly checking dashboards or reports, Data Activator continuously watches data sources like Power BI reports or real-time streams and responds when a predefined rule is matched. For example, if sales drop below a certain threshold or a sensor reports an unsafe reading, Data Activator can immediately send an email, alert a team, or trigger another action.

It's useful in scenarios such as monitoring KPIs, catching operational issues early, or automating business responses based on data events. It works like an always-on watchtower for your data.

44. How does Data Activator detect patterns or thresholds in data to trigger alerts or automated actions?

Data Activator works by using rules that you define to monitor incoming data. These rules are called **triggers**, and they can detect events like values crossing a threshold, changes in a trend, or specific data combinations.

Here's how it works step by step:

1. **Bind to a data source:** You connect Data Activator to a dataset from Power BI or a streaming source like Eventstream or a KQL table.
2. **Define a trigger condition:** You create a rule based on the data fields. For example, "if temperature > 100°C" or "if sales this week are 20% lower than last week."
3. **Activate an object:** The object is usually a record or item (like a product, device, or region) that you want to monitor.
4. **Choose an action:** When the trigger is met, Data Activator can send an alert via email, Microsoft Teams, Power Automate, or custom webhooks.

It continuously checks the live data source for new values and automatically fires the alert or action when the condition is matched. This allows organizations to catch problems early or respond to opportunities faster without manual monitoring.

45. What types of data sources can Data Activator monitor, and how is it integrated with real-time and batch data in Fabric?

Data Activator in Microsoft Fabric supports both real-time and batch data sources, allowing it to react to a wide variety of analytics situations:

- **Power BI Datasets:** You can bind Data Activator to a visual or data field in a Power BI report. It monitors the refreshed data in the dataset and acts on it after every update. This is useful for KPIs and business metrics that are updated regularly.
- **KQL Tables (Real-Time Analytics):** It can monitor live data flowing through Kusto-based databases. This is ideal for telemetry, IoT, logs, and high-frequency data events.
- **Eventstream:** Data Activator can listen to events coming in from Eventstream pipelines. These can originate from Event Hubs, IoT Hub, or Kafka, and the data can be filtered and monitored as it arrives.
- **Lakehouse and Warehouse** (indirectly): While Data Activator doesn't connect directly to Lakehouse files, you can expose your data through Power BI visuals or KQL streams that are backed by Lakehouse or Warehouse data.

This integration means you can use Data Activator for both batch use cases (such as nightly updated Power BI datasets) and real-time scenarios (like detecting spikes in sensor data). It connects well with other Fabric components and allows you to turn insights into immediate actions automatically.

46. What kinds of actions can be triggered by Data Activator, and how are they typically used in business scenarios (e.g., Power Automate, email, Teams alerts)?

Data Activator can trigger several types of actions when a data condition or event is met. These actions are designed to automate responses or alert users so they can take immediate steps. The most common actions include:

- **Email Notifications:** You can automatically send an email to a user or group when a trigger condition is met. For example, if a product's inventory level drops below a certain threshold, an email can be sent to the supply chain team.
- **Microsoft Teams Alerts:** Data Activator can send messages directly to a Teams channel or user. This is helpful for real-time collaboration for instance, notifying the customer service team when there's a sudden spike in unresolved tickets.
- **Power Automate Workflows:** You can connect Data Activator to Power Automate to launch multi-step business processes. For example, if a delivery is delayed beyond a certain time, a Power Automate flow can update the CRM system, notify the customer, and create a support ticket all automatically.
- **Custom Webhooks:** For more technical or custom integrations, Data Activator can call a webhook. This lets you trigger external systems or APIs when certain events occur in your data.

These actions help businesses respond faster to real-world changes by closing the gap between insights and action. They reduce the need for constant manual monitoring and support a more automated, event-driven workflow.

47. How does Data Activator differ from traditional alerting mechanisms in Power BI or Azure Monitor, and when would you use it over those tools?

Data Activator is designed for more flexible, proactive, and intelligent data monitoring compared to traditional alerting systems like Power BI alerts or Azure Monitor.

Here's how it differs:

- **Data binding flexibility:** In Power BI, alerts are limited to visuals like KPI cards and only trigger on a single value change. Data Activator can monitor entire datasets, fields, or patterns across multiple objects (like per region or per product), not just a single number.
- **Event-driven actions:** Data Activator supports richer actions like Power Automate flows, Teams integration, and custom webhooks. Power BI alerts mostly send emails or trigger simple flows. Azure Monitor is focused more on infrastructure and system metrics, not business data events.
- **Object-level monitoring:** With Data Activator, you can activate multiple objects from one data stream (like each product or sensor) and monitor them independently. Power BI alerts don't support this level of dynamic, object-specific logic.
- **Real-time data support:** Data Activator is tightly integrated with Eventstream and KQL for real-time data. Power BI alerts are based on scheduled refreshes, and Azure Monitor is used mostly for infrastructure-level logs, metrics, or telemetry not business metrics.

You would use Data Activator when:

- You need to trigger actions across many items (like all devices or all regions) dynamically.
- You're working with real-time data and need instant reaction.
- You need to go beyond simple alerts and connect to business workflows or custom systems.
- You want to monitor data flowing through Power BI, Eventstream, or KQL with advanced conditions.

In short, Data Activator gives more flexibility and power in reacting to business data events, while Power BI alerts and Azure Monitor are better for simpler thresholds or infrastructure-focused monitoring.

48. How does Microsoft Fabric track data lineage across components like Data Pipelines, Lakehouses, Notebooks, and Power BI reports?

Microsoft Fabric automatically captures and visualizes **data lineage** across its different components, making it easier to understand how data flows from source to consumption. This helps engineers and analysts trace data origin, transformation steps, and downstream usage.

Here's how it works:

- When you build a Data Pipeline that ingests data from a source (like SQL or Blob Storage) into a Lakehouse, Fabric records that relationship.
- If a Notebook or Dataflow Gen2 transforms data from that Lakehouse and writes the output into another table or dataset, the system captures that dependency as well.
- When a Power BI report or semantic model connects to the Lakehouse or a Warehouse, the lineage is extended to the visualization layer.

All these connections are automatically tracked and shown in a lineage view, which is available in the workspace UI. The lineage diagram shows:

- Which data items are upstream (sources)
- Which items are downstream (consumers)
- What the flow path looks like from ingestion to report

This lineage helps users answer questions like “Where does this report get its data from?” or “If I change this table, what dashboards will be impacted?”

49. How does Microsoft Purview integrate with Microsoft Fabric to support data discovery, classification, and governance?

Microsoft Purview integrates with Microsoft Fabric to provide **enterprise-level data governance**, including features like data cataloging, classification, policy enforcement, and auditing.

Here's how Purview supports Fabric:

- **Data catalog and discovery:** Purview scans Fabric items such as Lakehouses, Warehouses, Power BI datasets, and Pipelines. It brings them into a central catalog where users can search for data assets across the entire organization.
- **Data classification:** Purview can automatically classify sensitive data in Fabric (like email addresses, credit card numbers, or health information). These labels help identify compliance risks and ensure privacy policies are applied.
- **Sensitivity labels and policies:** Using Microsoft Information Protection (MIP), Purview allows you to apply sensitivity labels directly to Fabric items, which can help control sharing, access, and encryption.
- **Lineage and impact analysis:** Purview enhances Fabric's native lineage by allowing cross-platform tracking. For example, if a Power BI dataset is consuming a SQL database outside Fabric, Purview can show the full path from external system to report.
- **Role-based access and auditing:** Governance teams can use Purview to define access control policies and track usage activity across domains and workspaces.

With this integration, organizations get centralized visibility and control over their Fabric data assets, while still allowing individual domain teams to manage and produce data independently.

50. What item-level metadata is captured in Fabric, and how can engineers or analysts use it to trace data origin and impact?

Microsoft Fabric captures rich item-level metadata to support data management, governance, and lineage tracking. This metadata includes both technical and business details about each data object.

Examples of metadata captured:

- **Name and type:** The item's name (like a Lakehouse table or pipeline) and its category (table, notebook, dataset, etc.)
- **Owner:** The user or group responsible for the item.
- **Created and modified dates:** Helps track activity and identify stale items.
- **Schema details:** For tables or datasets, this includes column names, data types, and partitioning.
- **Source and destination paths:** Shows where the data came from (e.g., a pipeline source or raw folder) and where it is being written to.
- **Lineage links:** Metadata that shows upstream and downstream relationships to other items like Dataflows, Notebooks, or Power BI reports.
- **Tags and sensitivity labels:** Business metadata added manually or through Purview for classification and policy control.

Engineers and analysts can use this metadata to:

- **Trace data origin:** Understand where a dataset comes from and how it was transformed.
- **Impact analysis:** See what downstream systems will be affected if a table or column is changed.
- **Debug issues:** Quickly locate the source of broken reports or data quality problems.
- **Collaborate:** Use tags, owners, and descriptions to improve communication across teams.

This metadata is accessible through the workspace UI, lineage views, and, when integrated, through Microsoft Purview's data catalog. It supports better decision-making, safer changes, and stronger trust in the data.

51. How can custom tags and classifications be applied to items in Fabric, and how does this support enterprise data governance policies?

In Microsoft Fabric, you can apply **custom tags** and **classifications** to data items like Lakehouses, Warehouses, Pipelines, Dataflows, and Power BI datasets. These metadata labels help organizations organize, control, and track their data assets more effectively, especially in large environments.

Here's how it works:

- **Custom tags:** You can assign business-specific tags (like “confidential,” “financial,” “customer-data,” or “gold-certified”) to individual items. Tags help categorize data for easier discovery and filtering, especially when working across multiple domains and workspaces.
- **Classifications:** These are usually automated or defined through Microsoft Purview. They help identify sensitive information such as personal identifiers, financial records, or health data using built-in data classifiers. For example, Fabric can label a column containing emails as “Personally Identifiable Information (PII).”

These metadata labels support enterprise data governance in several ways:

- **Policy enforcement:** Classifications can be tied to sensitivity labels or access control rules. For example, data marked as “confidential” might be restricted from export or limited to specific users.
- **Compliance:** Automatically classifying sensitive data helps meet data privacy regulations like GDPR, HIPAA, or ISO standards.
- **Data lifecycle management:** Tags can help track which datasets are authoritative, archived, in development, or deprecated helping teams manage their data more effectively.
- **Improved discoverability:** Users can search and filter items based on tags or classifications, making it easier to find trusted, relevant data for analysis or reporting.

By tagging and classifying items consistently, governance teams gain more control and visibility over how data is handled across the organization.

52. What are the key benefits of catalog integration in Microsoft Fabric, and how does it help in enabling data mesh and self-service analytics?

Catalog integration in Microsoft Fabric especially through its connection with Microsoft Purview provides a central way to organize, discover, and govern data assets across the platform. This is a key enabler for both data mesh and self-service analytics.

Key benefits include:

- **Unified data discovery:** Users can search across Lakehouses, Warehouses, Pipelines, and Power BI datasets using business-friendly terms. This makes it easy for analysts and data consumers to find and use the data they need, even if it's stored in another domain or workspace.
- **Trust and quality indicators:** Catalog entries can show certifications (like “gold” datasets), tags, descriptions, and owners, helping users identify which data is trustworthy and well-maintained.
- **End-to-end lineage:** Integrated catalog views show the full flow of data from ingestion through transformation to reporting. This helps analysts understand context and engineers trace impact when making changes.

- **Cross-domain visibility:** In a data mesh model, each domain owns its data products but the catalog makes them visible and usable across the organization. This allows one domain to consume data products built by another, enabling collaboration without central bottlenecks.
- **Governance and compliance:** Catalog integration enables central teams to apply classification, sensitivity labels, and access policies across all domains, even though the data is managed locally.

Together, these features help balance decentralization (through domain ownership) with central governance and discoverability. It empowers users to explore and analyze data confidently, while still keeping the organization in control of compliance, quality, and security.

53. What is an F SKU capacity in Microsoft Fabric, and how does it determine the performance and concurrency of workloads?

An F SKU capacity in Microsoft Fabric is a capacity unit that defines the amount of dedicated compute power available to run various workloads in a Fabric tenant. The "F" stands for Fabric, and these SKUs are similar in concept to Power BI Premium capacities but extended to support the full set of Fabric services.

Each F SKU (like F2, F4, F8, etc.) provides a certain number of Capacity Units (CUs). These units are used to execute and manage jobs across Fabric components, including Lakehouses, Notebooks, Pipelines, Warehouses, and Real-Time Analytics.

Here's how it impacts performance:

- **Higher F SKUs** mean more compute is available. This results in faster execution times, more parallel jobs, and better handling of large data volumes.
- **Concurrency** is controlled by the capacity. A higher SKU allows more users or workloads to run simultaneously without delays or throttling.
- **Workload queuing** can happen if the capacity is fully used. Jobs may wait for resources unless the SKU is increased or usage is optimized.

In short, F SKU defines how much power and parallelism your environment has. Choosing the right SKU is important for ensuring good performance and user experience.

54. How is compute capacity shared and managed across different Fabric workloads such as Data Engineering, Data Factory, Power BI, and Real-Time Analytics?

In Microsoft Fabric, compute capacity from your F SKU is pooled and shared across all workloads within a capacity, making it a multi-engine platform with a unified resource pool.

Here's how it works:

- **Single pool:** Whether you're running a Power BI report, executing a Spark notebook, processing a Dataflow Gen2, or querying a Real-Time KQL table, all these actions consume from the same Fabric capacity pool.
- **Workload-specific limits:** While all workloads share the same pool, Microsoft applies **fair-share scheduling and workload rules** to ensure balance. For example, Spark (used in Data Engineering) and Power BI won't block each other completely, but they may compete for resources if heavily loaded.
- **Dynamic allocation:** Fabric intelligently allocates resources depending on the workload demand. If you're running more pipelines in the morning and Power BI reports in the afternoon, Fabric adjusts allocation to match that usage pattern.
- **Monitoring tools:** Admins can monitor capacity usage through the Fabric admin portal, where you can track how much compute is being used by each workload and user.

This shared model allows for better resource efficiency, especially in organizations where not all services are used heavily at the same time.

55. How does the licensing model of Microsoft Fabric (capacity-based) compare with the pay-per-use pricing of Azure Data Factory and Synapse Analytics?

Microsoft Fabric uses a capacity-based pricing model, while Azure Data Factory (ADF) and Synapse Analytics typically follow a pay-per-use model. Each model has advantages depending on the use case.

Fabric's Capacity-Based Model (F SKU):

- You pay for a fixed capacity (F SKU), which includes all services Power BI, Data Factory, Lakehouse, Real-Time Analytics, etc.
- All workloads share this capacity, and you can run as many jobs or reports as your capacity allows, without paying extra per operation.
- This model is predictable and cost-effective for organizations with high usage or many active users, since the compute cost doesn't grow with each additional activity.

ADF and Synapse Pay-Per-Use Model:

- In Azure Data Factory, you pay per pipeline run, per Data Movement (DIU-hour), or per data transformation (Mapping Data Flows).
- In Synapse, you pay for each query (in serverless SQL) or for the time a dedicated SQL or Spark pool is running.
- This model is more cost-efficient for occasional or light usage, as you only pay for what you consume.

Comparison Summary:

Feature	Microsoft Fabric (F SKU)	ADF / Synapse (Pay-Per-Use)
Billing Type	Fixed monthly per capacity	Per activity or compute usage
Cost Predictability	High	Variable based on workload
Best for	Heavy, mixed workloads	Low-volume, occasional workloads
Multi-workload integration	Unified and shared	Separated services
Governance and management	Unified under one tenant	Managed per service

Organizations often choose Fabric when they want centralized governance, all-in-one capabilities, and predictable billing, especially if they use Power BI heavily. Pay-per-use still works well when workloads are small, temporary, or isolated.

56. How do you use Data Factory in Microsoft Fabric to ingest data from sources like SQL Server, ADLS Gen2, REST APIs, and cloud storage services?

In Microsoft Fabric, we use the Data Factory workload to build data pipelines that ingest and move data from various sources. The process is similar to Azure Data Factory but is fully integrated into the Fabric environment.

To ingest data from sources like SQL Server, ADLS Gen2, REST APIs, or cloud storage, I follow these steps:

1. First, I create a new pipeline in the Data Factory section of Fabric.
2. Then, I add a Copy activity to the pipeline. This activity allows me to define both the source and the destination (sink).
3. For each source, I create a connection using a linked service. For example:
 - For SQL Server, I provide the server name, database, and authentication details.
 - For ADLS Gen2 or Blob storage, I select the storage account and container.
 - For REST APIs, I configure the base URL, headers, and authentication (if needed).
4. Once the source is connected, I define the dataset or file I want to pull. I can also apply filters or select specific columns.
5. Then, I configure the sink, such as a Lakehouse table, Warehouse table, or another storage location.
6. I can schedule the pipeline to run on a trigger or run it manually.

This setup lets me move structured and unstructured data into Fabric for further processing and reporting. The integration is smooth and allows me to ingest from a wide variety of sources using the same pipeline structure.

57. What is the role of the Copy activity in Fabric Data Pipelines, and how does it handle schema mapping, performance optimization, and sink configuration?

The Copy activity is the main component in a Fabric data pipeline that helps move data from a source to a destination. Its role is to extract the data, transform it (if needed), and load it into the target system.

For schema mapping:

- The Copy activity supports both auto-mapping and manual mapping. If the source and destination have the same column names and types, it can auto-map them.
- If there are differences in column names or data types, I can define custom mappings using the UI or JSON editor.

For performance:

- Fabric uses scalable and distributed compute behind the scenes to copy data efficiently, especially when dealing with large volumes.
- I can enable parallelism and partitioning options if the source supports it, such as when reading from a large SQL table or multiple files.
- It also allows me to set retry policies, fault tolerance, and logging to make the pipeline more reliable and efficient.

For sink configuration:

- I can choose different types of destinations, such as Lakehouse tables, SQL Warehouse tables, or storage folders.
- It supports overwrite, append, or upsert modes depending on the destination.
- I can also choose to write data in formats like CSV, Parquet, or Delta.

The Copy activity is very flexible and is used as the backbone for building ETL or ELT processes in Fabric pipelines.

58. How can you use parameterization in Fabric Data Pipelines to build dynamic and reusable workflows across environments and datasets?

Parameterization in Fabric Data Pipelines allows me to make my workflows dynamic, reusable, and easier to manage across multiple datasets or environments.

Here's how I use it:

1. I define pipeline parameters at the top level of the pipeline. These parameters could be things like source file name, database name, table name, or even environment (dev, test, prod).
2. Inside the pipeline, I use these parameters in activities like Copy or Lookup by referencing them with expressions. For example, I can pass a file name into the Copy activity's source path so the pipeline works with different files depending on input.
3. I can also create dynamic datasets and linked services using parameters. This helps me connect to different data sources without hardcoding details.
4. When triggering the pipeline, I pass in the parameter values either manually, through a schedule trigger, or from another pipeline.

This method is especially helpful when I want to deploy the same pipeline across different environments. For example, I might use a different storage path or SQL server name in dev versus prod. With parameterization, I don't need to build separate pipelines. I just change the values of the parameters.

Overall, parameterization reduces duplication, improves maintainability, and makes pipelines more scalable and flexible.

59. How do scheduling and trigger mechanisms work in Fabric Data Factory, and how do they differ from the traditional Azure Data Factory triggers?

In Fabric Data Factory, scheduling and triggers are used to automate when a pipeline should run. These mechanisms are built into the Fabric experience and are similar in purpose to Azure Data Factory but have some differences because Fabric is still evolving as a unified platform.

In Fabric, we use schedule triggers to run pipelines at specific times or on recurring intervals. For example, I can schedule a pipeline to run every day at 6 AM or every 15 minutes. To do this, I go to the pipeline trigger section and define a schedule with frequency (hourly, daily, weekly), time zone, and start time.

Fabric currently supports time-based triggers only, unlike traditional Azure Data Factory, which supports other trigger types such as:

- **Event-based triggers** (like when a file is dropped into Blob Storage)
- **Dependency triggers** (based on the completion of other pipelines)

This means that in Fabric, I cannot yet trigger a pipeline based on file arrival or a custom event, but I can still work around this by scheduling it frequently or using control logic in the pipeline to check for new data.

Another difference is that in Fabric, triggers are more workspace-aware and tightly connected with the user interface. The experience is more integrated across the platform, making it easier to manage for business users as well.

So while Fabric has more limited trigger types compared to Azure Data Factory for now, the time-based scheduling is simple, reliable, and works well for most batch data workflows.

60. What strategies do you use in Fabric Pipelines to handle errors, retries, logging, and alerting during data ingestion and transformation processes?

In Fabric Data Factory, it's very important to handle errors and monitor pipeline health to ensure data workflows are reliable and maintainable. I follow several strategies to manage errors, retries, logging, and alerts.

1. Error handling with activity outputs:

I check the output of each activity using control structures like "If Condition" and "Set Variable". For example, if a Copy activity fails, I route the workflow to a different branch to log the failure or take a different action.

2. Built-in retry settings:

Each activity in the pipeline (like Copy or Dataflow) has a retry setting. I usually configure retries (like 3 attempts) with a delay between them. This helps handle temporary issues like network hiccups or source availability problems.

3. Custom logging:

I use an additional activity (like Web or Copy to a log table) to capture details like pipeline name, activity status, start and end time, error messages, and row counts. This helps with tracking historical runs and debugging failures.

4. Failure alerts:

While Fabric doesn't yet support full alerting like Azure Monitor, I use Power Automate integrated with Fabric to send emails or Teams notifications when a pipeline fails or a certain condition is met. This lets me stay informed without checking manually.

5. Testing and control activities:

I add "Wait" and "Switch" activities to control the flow and simulate conditions before full deployment. Also, I break complex workflows into smaller modular pipelines and call them using the "Execute Pipeline" activity. This makes error isolation and recovery easier.

6. Idempotent design:

I design my pipelines so that re-running them does not corrupt data. For example, using watermarking or writing to staging tables and then merging helps ensure consistent results even after retries.

By combining these practices, I ensure that Fabric pipelines are robust, can recover from failures, and provide visibility into what's happening behind the scenes.

61. How do you use PySpark or SQL in Microsoft Fabric Notebooks to perform data transformation tasks on large datasets?

In Microsoft Fabric Notebooks, I can use PySpark or Spark SQL to transform large datasets efficiently. The Spark engine in Fabric runs behind the scenes and allows me to process data in a distributed way, which means it can handle big volumes of data with good performance.

To perform transformations using PySpark, I usually follow these steps:

1. I start by reading data from a Lakehouse table or a folder using `spark.read` functions.
2. I apply transformations using PySpark DataFrame methods like `filter`, `select`, `withColumn`, `groupBy`, and `join`.
3. After transforming the data, I write the results back to the Lakehouse in Delta format.

I can also use `%sql` magic commands in notebook cells to write SQL queries directly. This is helpful when working with users who are more comfortable with SQL syntax.

The benefit of using Fabric Notebooks is that I can run these transformations interactively, test them step by step, and view the results immediately in the output below each cell.

62. How do you read from and write to a Lakehouse using Spark in Fabric Notebooks? Can you share a basic example using PySpark?

In Fabric Notebooks, reading and writing data to a Lakehouse is very straightforward using PySpark. The Lakehouse is already connected to the Spark engine, so I don't need to configure any external storage manually.

Here's a simple example:

```
# Reading a Delta table from the Lakehouse
```

```
df = spark.read.format("delta").load("Tables/sales_data")
```

```
# Applying a transformation: filtering rows with sales above 1000
```

```
filtered_df = df.filter(df["sales_amount"] > 1000)
```

```
# Writing the transformed data back to another Lakehouse table
```

```
filtered_df.write.mode("overwrite").format("delta").save("Tables/high_sales_data")
```

In this example:

- The `load("Tables/sales_data")` path refers to a managed Lakehouse table.
- I use PySpark's `filter` function to transform the data.
- The `save("Tables/high_sales_data")` command creates a new table or overwrites an existing one.

I can also write in append mode if I want to keep adding data instead of replacing it.

This workflow allows me to process data efficiently and store the results in a structured and queryable format within the same Fabric workspace.

63. What are the key capabilities of Fabric Notebooks for interactive data exploration and profiling, and how do they compare with notebooks in Synapse or Databricks?

Fabric Notebooks are built for interactive data exploration, especially when working with Lakehouse data using Spark. Some of the key features that I use regularly are:

- **Multi-language support:** I can write cells in PySpark, SQL, or Markdown using magic commands like `%pyspark`, `%sql`, and `%md`. This is useful when switching between coding, querying, and documentation in one place.
- **Instant feedback:** The output from each cell appears immediately below it, allowing me to test code, visualize data, and debug in real time.
- **Built-in Lakehouse integration:** I don't need to connect storage manually. I can browse Lakehouse tables directly from the notebook's left pane and drag-and-drop them into a code cell.
- **Schema and data preview:** I can preview table schemas and sample data quickly, which helps when profiling a dataset before transformation.
- **Notebook scheduling and reuse:** I can trigger notebooks from pipelines or schedule them using the Fabric interface. This makes it easy to operationalize code I've tested.

Compared to Synapse Notebooks:

- Fabric Notebooks feel more tightly integrated with the Lakehouse and have a more modern UI.
- Synapse Notebooks support similar Spark operations but require more setup to connect to storage or link to external sources.

Compared to Databricks:

- Databricks offers more advanced features like MLflow integration and job clustering, but for standard data engineering and exploration tasks, Fabric Notebooks provide a simple and integrated experience, especially within the Microsoft ecosystem.

Overall, Fabric Notebooks are a strong choice for interactive data analysis and transformation, especially when working closely with Lakehouse, Power BI, and other Fabric-native tools.

64. How does Spark execution work under the hood in Fabric Notebooks, and how is compute managed and scaled during workload execution?

In Fabric Notebooks, when I run PySpark or SQL code, it uses a managed Spark engine that is fully integrated into Microsoft Fabric. I don't need to set up any Spark cluster manually because Fabric handles it for me automatically.

Here's how Spark execution works in the background:

- When I run a cell with Spark code, Fabric provisions the compute behind the scenes using a pool of resources that are part of the assigned Fabric capacity (F SKU).
- Spark jobs are broken into stages and tasks, and Fabric distributes these tasks across the available resources. This allows the job to run in parallel, which is very useful when working with large datasets.
- The driver program manages the overall job and sends work to executors, which do the actual data processing. In Fabric, these executors are dynamically managed inside the Spark environment.
- The system automatically handles session management so each notebook has a session that stays active while I'm working. If I leave the notebook idle for a while, the session can time out to save resources.
- Fabric adjusts compute usage based on capacity. If other workloads are also using the same capacity, Spark may slow down or queue jobs if resources are limited. That's why sizing the capacity correctly is important for consistent performance.

So in short, Spark in Fabric Notebooks runs like a managed service, scaling based on the capacity and workload size, and simplifying the developer experience by removing the need to manually configure clusters or compute settings.

65. What techniques do you use in Fabric Notebooks to handle missing data, outliers, and schema drift when preparing data for analytics or machine learning?

When I prepare data for analytics or machine learning in Fabric Notebooks using PySpark, I use different techniques to clean and manage the data depending on the issue I face.

1. Handling missing data:

- I use `dropna()` to remove rows with null values if those rows are not useful or if they're a small portion of the data.

```
df_clean = df.dropna()
```

- If dropping rows is not an option, I use `fillna()` to replace missing values with a default value or mean.

```
df_filled = df.fillna({'age': 30, 'income': 0})
```

2. Handling outliers:

- I use summary statistics like mean, standard deviation, and percentiles to detect outliers.

```
df.describe().show()
```

- I use filtering logic to remove or flag rows with values outside expected ranges.

```
df_filtered = df.filter(df["income"] < 100000)
```

- For machine learning tasks, I sometimes apply transformations like normalization or capping to reduce the impact of outliers.

3. Handling schema drift:

- Schema drift happens when the structure of incoming data changes over time (like new columns appearing or data types changing).
- To detect it early, I often read a sample file first and inspect the schema with `df.printSchema()`.
- I use the `mergeSchema` option when reading data with evolving columns in Delta format:

```
df = spark.read.option("mergeSchema", "true").format("delta").load("Tables/data")
```

- I also write code that checks if expected columns exist before processing, so the notebook does not break when new data arrives.

By using these techniques, I make sure the data is clean, consistent, and ready for analysis or modeling. Fabric Notebooks help a lot here because I can interactively test and fix issues step by step and see the output instantly.

66. How would you design an end-to-end data engineering workflow in Microsoft Fabric, starting from data ingestion and landing it in a Lakehouse, followed by transformation and reporting in Power BI?

To design a complete data engineering workflow in Microsoft Fabric, I follow a series of well-defined steps using different Fabric components that work together smoothly. Here's how I approach it:

1. Data ingestion using Data Factory (Pipelines):

- I create a pipeline that uses the Copy activity to ingest data from sources like SQL Server, ADLS Gen2, REST APIs, or cloud storage.
- I configure the source dataset and define a Lakehouse table or folder as the sink.
- I use scheduling to run the pipeline daily or hourly, depending on business needs.
- The raw data is written into the "Files/raw" folder or directly into "Tables" in Delta format.

2. Data transformation using Fabric Notebooks or Dataflows Gen2:

- For complex transformations, I use a PySpark notebook to read the raw data, apply filters, joins, and calculations, and then write the cleaned data into curated Lakehouse tables.
- For simpler, visual transformations, I use Dataflows Gen2, which also allows me to store the output directly in Delta format inside the Lakehouse.

3. Data modeling and reporting using Power BI:

- I create a semantic model (dataset) by connecting Power BI to the curated Lakehouse tables.
- I define relationships, calculated columns, and measures using DAX.
- I build dashboards and reports directly inside Fabric using Power BI Reports, which automatically use the same storage layer (OneLake).

4. Automation and orchestration:

- I use triggers in Data Factory to automate the entire pipeline.
- If needed, I also use Notebooks inside pipelines or chain multiple pipelines together.

This workflow allows me to move raw data to refined, report-ready data in a seamless way, fully within Fabric, without needing to copy data across tools.

67. What strategies do you use to handle schema drift and evolving data structures across different stages of the pipeline in Fabric?

Schema drift is when the structure of incoming data changes unexpectedly, such as new columns being added or data types changing. I use several strategies in Fabric to manage this:

1. Use Delta Lake format with mergeSchema option:

When reading or writing data using PySpark, I enable mergeSchema to allow new columns to be added automatically without breaking the pipeline.

```
df = spark.read.option("mergeSchema", "true").format("delta").load("Tables/sales_data")
```

2. Schema inference in Dataflows Gen2 and Copy activity:

In pipelines and Dataflows, I allow schema inference if I want the system to adjust automatically, but I review schema changes during testing to avoid unexpected results.

3. Column presence and type checks:

Before processing, I check if required columns exist using Python or SQL logic. If not, I log the issue and skip processing that batch.

4. Column versioning or staging layers:

I store incoming data in a raw layer and only move it to curated tables after validation. This gives me time to inspect schema changes and make updates to the transformation logic.

5. Monitoring and alerting:

I log schema metadata during ingestion and compare it to expected structures. If a mismatch is found, I trigger an alert (via Power Automate or email) so the data team can respond quickly.

These strategies help make my pipelines flexible and more resilient to changing source data.

68. How do you implement data validation and quality checks in Fabric pipelines to ensure accurate and trustworthy reporting?

To ensure that the data I load and transform is correct and reliable, I add several layers of validation and quality checks at different stages:

1. Validation during ingestion (in pipelines):

- I use Lookup or Script activities to check if source tables exist and contain rows.
- I log metadata like row counts, file size, and load timestamps after each Copy activity.
- If row counts are zero or much lower than expected, I stop the pipeline or send an alert.

2. Validation in transformation steps (in Notebooks):

- I check for nulls in important columns using filters or `df.describe()` to find data quality issues.
- I apply business rules, like validating that dates are within expected ranges or that numeric fields are not negative.
- If data fails these rules, I log it in a separate error table for review.

3. Post-load validation in Lakehouse or Power BI:

- I run summary queries to compare key metrics over time and check for sudden spikes or drops.
- I create validation dashboards in Power BI that track data freshness, row counts, and rule violations.

4. Automated alerts for failures:

- I connect pipelines to Power Automate or use custom notebook logic to send alerts when validations fail or thresholds are breached.

5. Logging and audit trail:

- I write detailed logs to Lakehouse tables or files, including batch IDs, record counts, and validation results. This helps with auditing and debugging if something goes wrong.

By applying these checks throughout the pipeline, I ensure that the final data used in reports is accurate, complete, and trustworthy. This builds confidence among stakeholders and helps avoid surprises during reporting.

69. What are your approaches for setting up monitoring, logging, and observability in Fabric pipelines to track data flow and detect anomalies?

To set up monitoring and logging in Fabric pipelines, I try to make sure that every stage of the pipeline gives me enough information to understand what happened, how much data was processed, and whether anything went wrong.

Here are the steps I follow:

- 1. Use activity outputs for tracking**

Each activity in the pipeline, like Copy or Notebook, has output metadata. I capture details such as row counts, duration, and success or failure status. I often store this information in a Lakehouse table or log file for historical tracking.

- 2. Add custom logging logic**

I use Set Variable or Script activities to log custom messages at different stages of the pipeline. For example, after each successful step, I log a message saying "Step X completed with Y records".

- 3. Create a metadata or audit table**

I usually maintain a table in the Lakehouse to store audit logs. Each pipeline run inserts a row with pipeline name, date, source, sink, row counts, and status. This helps with reviewing past runs and understanding trends or anomalies.

- 4. Detect anomalies using basic rules**

I include checks like comparing current row counts with previous runs. If there's a sudden drop or spike in data volume, I flag it in the log and mark the status as a warning.

- 5. Visual monitoring in Fabric UI**

I also use the Fabric pipeline monitoring tab to review run history, success and failure rates, and duration of activities. This helps during live debugging and post-run analysis.

By combining these methods, I build pipelines that are not only functional but also easy to monitor, debug, and maintain.

70. How do you implement automated error handling and alerting (e.g., using Data Activator or Power Automate) for pipeline failures or data issues in Fabric workflows?

For automated error handling and alerting in Fabric, I try to design pipelines that detect problems early and notify the right people automatically, so manual checks are not needed every time.

Here's how I implement it:

1. Use failure paths in the pipeline

I use the built-in failure branches in the pipeline design. For example, if a Copy activity fails, I connect it to a Web or Script activity that logs the error or calls a Power Automate flow for alerting.

2. Integrate with Power Automate

I create a Power Automate flow that sends an email, Teams alert, or creates a ticket when triggered. In the pipeline, I use a Web activity to call the Power Automate webhook and pass error details like activity name, error message, and pipeline name.

3. Use Data Activator for threshold-based alerts

When data is successfully loaded into the Lakehouse or Warehouse, I can set up a Power BI dashboard on top of it. Then I use Data Activator to monitor metrics like record counts or validation flags. If a metric crosses a certain threshold, Data Activator can automatically trigger actions like sending a message or starting a workflow.

4. Retry mechanisms for temporary failures

I configure retries on activities like Copy, Lookup, or Notebook in case the failure is temporary (for example, due to network issues). This avoids sending alerts for issues that resolve on their own.

5. Mark pipeline run status clearly

I make sure the pipeline ends with a clear success or failure status by using condition checks. If the pipeline ends in failure, it can be picked up easily by monitoring tools or custom reports.

By using these strategies, I make sure that any data failure or issue is quickly identified and that the team is informed without manual intervention. This helps keep the workflows reliable and production-ready.

71. What are the key differences between a Lakehouse and a Data Warehouse in Microsoft Fabric in terms of storage format, query engine, and use cases?

In Microsoft Fabric, both Lakehouse and Data Warehouse are used for analytics, but they are designed differently and serve different types of users and workloads.

Here's how they differ:

1. Storage format

- A Lakehouse stores data in an open file format called Delta Lake, which is built on top of Parquet. The data is stored in folders and files inside OneLake, which is a distributed storage layer.
- A Data Warehouse stores data in relational tables managed by the SQL engine. The data is abstracted from file-level storage and optimized for SQL performance.

2. Query engine

- The Lakehouse supports Spark and SQL engines. Spark is used for big data processing using notebooks or Dataflows Gen2, while SQL can be used through the Lakehouse SQL endpoint for querying.
- The Data Warehouse is built for high-performance SQL querying only. It uses a distributed SQL engine with features like automatic indexing and result caching.

3. Use cases

- The Lakehouse is ideal for data engineers and data scientists who work with semi-structured data, large files, and need flexibility for data exploration, transformation, and machine learning.
- The Data Warehouse is better for business analysts and reporting teams who rely on structured, cleaned data for dashboards, ad hoc queries, and BI reporting with strong performance.

In summary, Lakehouse gives flexibility and openness, while the Data Warehouse gives speed and simplicity for structured analytics.

72. When would you choose to use Fabric's SQL-based Data Warehouse over the Lakehouse model, and what types of workloads does it serve best?

I would choose to use the SQL-based Data Warehouse in Fabric when I need high-performance SQL queries on structured data with minimal setup and no need for Spark-based transformations.

It works best in these situations:

1. **BI-focused workloads**

If the main goal is to build Power BI reports that query clean and modeled data, the Warehouse gives fast response times and supports large numbers of users with low latency.

2. **Data marts or reporting layer**

When I want to create a data mart with star schemas or normalized tables that are easy to query with SQL and DAX, the Warehouse is a great fit.

3. **Minimal data transformation**

If the data has already been cleaned or transformed before arriving in the Warehouse, then I don't need the flexibility of Spark.

4. **Concurrency and scale**

When many users or dashboards are expected to query data at the same time, the Warehouse can handle it better than Lakehouse SQL endpoints due to better query optimization and caching.

5. **Business user-friendly access**

Analysts can use the Warehouse without learning Spark or dealing with file structures. It behaves like a traditional SQL database, which is familiar and easier to use.

In short, I choose the Data Warehouse when I need structured, fast, and scalable SQL-based reporting on clean data.

73. How does Fabric's Data Warehouse handle compute scaling, concurrency, and performance compared to Synapse Dedicated SQL Pools?

Fabric's Data Warehouse offers several improvements over Synapse Dedicated SQL Pools, especially in terms of simplicity, performance tuning, and management.

Here's how it compares:

1. Compute scaling

- In Fabric, compute is tied to the overall capacity (F SKU), so I don't need to manage or resize the Warehouse manually. It auto-scales based on the available capacity and workload.
- In Synapse Dedicated SQL Pools, I had to manually scale up or down the compute (DWUs), which required planning and could interrupt queries.

2. Concurrency

- Fabric automatically handles multiple users running queries at the same time. It uses techniques like query queuing, workload management, and result caching.
- Synapse Dedicated SQL Pools support concurrency but can be limited by workload groups and user sessions, which needed more manual tuning.

3. Performance

- Fabric's Warehouse benefits from automatic indexing, data caching, and optimizations under the hood. It delivers faster results without much manual configuration.
- Synapse performance depends heavily on how well tables are distributed and indexed. Poor design in Synapse could lead to slower queries.

4. Management and cost

- Fabric is simpler to manage because everything is managed through capacity. I don't have to pause, resume, or pay for idle time like in Synapse.
- Synapse Dedicated SQL Pools run constantly unless paused, so managing cost and performance required more effort.

So, in Fabric, I get a more modern and hands-off experience with better built-in performance tuning, while Synapse gave me more control but also more complexity to manage.

74. How can you integrate a Lakehouse and a Warehouse in Fabric to create a hybrid architecture, and what are the benefits of doing so?

In Microsoft Fabric, it is possible to build a hybrid architecture by combining a Lakehouse and a Warehouse. I use this pattern when I want the flexibility of file-based processing from a Lakehouse along with the performance and structured access of a Warehouse.

Here is how I integrate them:

1. Ingest and store raw or semi-structured data in the Lakehouse

- I use Data Factory pipelines, Dataflows Gen2, or Spark notebooks to bring in data from various sources and land it in the Lakehouse as Delta tables or raw files.
- The Lakehouse is good for data engineering and transformation using PySpark or SQL.

2. Clean and transform the data in the Lakehouse

- I perform filtering, joining, column renaming, and business logic using notebooks or SQL queries on the Lakehouse Delta tables.

3. Load the refined data into the Warehouse

- Once the data is ready, I load it into the Warehouse using a pipeline or a notebook. This gives the reporting team access to fast, structured SQL tables.
- The Warehouse becomes the source for Power BI dashboards, business reports, and ad-hoc queries.

The benefits of this hybrid setup include:

- **Best of both worlds:** The Lakehouse is flexible for big data and transformations, while the Warehouse is optimized for reporting and performance.
- **Separation of concerns:** Data engineers work in the Lakehouse, while business users and analysts work in the Warehouse.
- **Better performance for BI:** Power BI queries run faster on a Warehouse because of result caching, automatic indexing, and optimized execution plans.
- **Governance and security:** I can apply different access controls at each layer, ensuring only refined and trusted data goes into the reporting layer.

This setup is very useful in enterprise scenarios where data volume is large, but reporting needs to be fast and simple for end users.

75. What are some best practices for loading data from a Lakehouse (Delta format) into a SQL Warehouse in Fabric for BI or operational reporting?

When I move data from a Lakehouse to a Warehouse in Fabric, I follow certain best practices to ensure the process is efficient, reliable, and easy to maintain.

Here are the steps I usually take:

1. Use the Copy activity in a pipeline

- I use the Copy activity to load data from Delta tables in the Lakehouse into tables in the Warehouse.
- I set the Lakehouse Delta table as the source and the Warehouse table as the sink.
- This allows me to reuse metadata and avoid writing manual code.

2. Incremental loading

- To avoid copying the full dataset every time, I use watermarking or last-modified timestamps to load only new or changed data.
- I add filters in the source query to extract only recent data.

3. Validate data before loading

- I include checks like row count comparison or key column validation before loading data into the Warehouse.
- I log these checks for tracking and troubleshooting.

4. Use staging tables when needed

- If data transformations are needed during the load, I first load the data into a staging table in the Warehouse, then apply updates or merges into the final table.
- This makes error recovery easier and improves data consistency.

5. Monitor and alert

- I add monitoring steps in the pipeline to track load success, duration, and row counts.
- If there's an error or unusual drop in volume, I send alerts using Power Automate or log the issue in a Lakehouse log table.

6. Optimize the destination Warehouse tables

- I create appropriate indexes or clustering keys in the Warehouse to improve performance.
- I also keep the table schema clean and use proper data types to reduce storage cost and improve query speed.

By following these best practices, I ensure that my data movement from Lakehouse to Warehouse supports reliable reporting, runs efficiently, and keeps business users happy with accurate and up-to-date data.