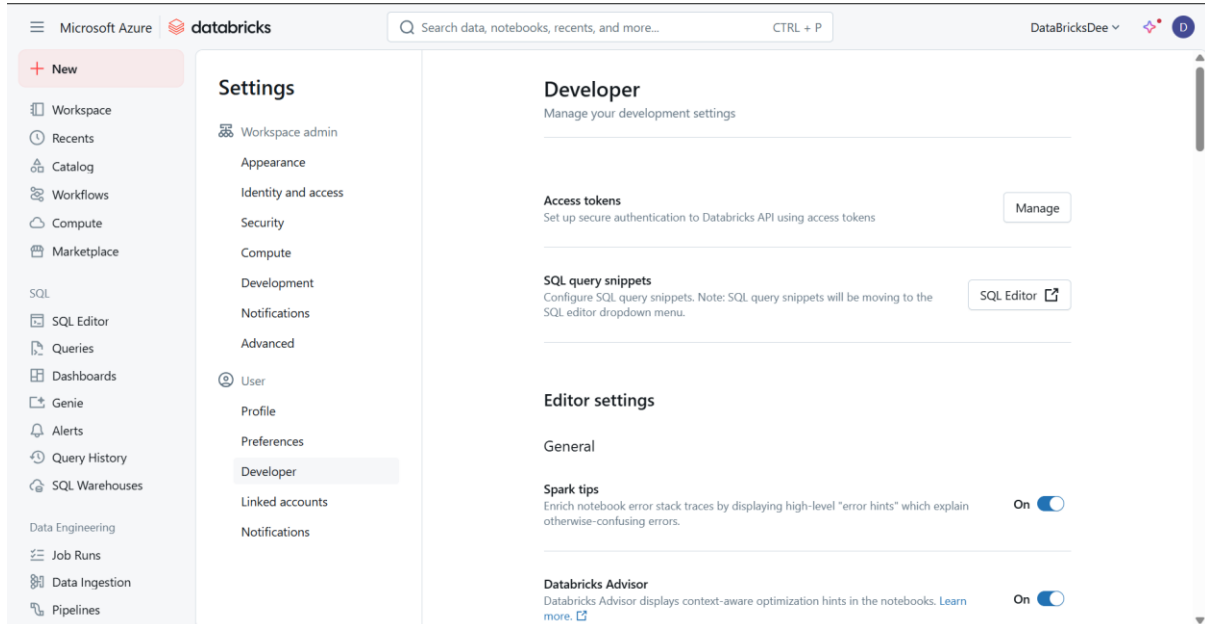# Connect notebook using synapse workspace

First, we have to generate Access token in Data bricks

For that we have to go to profile -> settings -> user -> Developer -> Access token -> Manage



Manage - > Generate new token.

Copy token: dapi6780fe307c20bfdab0ab53ec1e223779



# Generate New Token

Your token has been created successfully.

dapi6780fe307c20bfdab0ab53ec1e223779

⚠ Make sure to copy the token now. You won't be able to see it again.

Done

User settings › Developer ›

## Access tokens

Personal access tokens can be used for secure authentication to the Databricks API instead of passwords.

Generate new token

| Comment | Creation | Expiration | |
|---------|----------|------------|---|
| Connecting to synapse | 2025-04-14 11:05:22 EDT | 2025-04-15 11:05:22 EDT | 🗑 |

Now we have to create a secret key for access token in key vault

Go to key vault -> Objects -> Secrets -> Generate/ Import

**Now we have to create linked services for both Key vault and data bricks.**

Go to synapse -> Manage -> Linked services -> New



Once we click on new, search for Azure key vault as we have to first create linked services for Azure key vault.

Once you select this option, we have provided all the details, about subscription, what key vault we have to use and then test connection.



Once connection is successful, create the linked service.

Once it is created publish the changes for it to be saved

**Now create Data bricks linked service**

Click on new and search for Databricks.



Provide all required details and test connection

After linked services are created, create a pipeline, to run the notebook.

Go to Integrate -> new Pipeline -> drag and drop notebook from data bricks



Now go to Azure data bricks section and choose linked service.

As we already created a linked service, we can directly use it.



Now go to settings and choose which note book we want to run.

Validate and run pipeline.

If cluster is on it will take less time, but if cluster is off first, cluster will start and once it is up and running then it will run the notebook.



Pipeline ran successfully.

## Explore medallion architecture and write the best way of organizing data

**What is Medallion Architecture?**

The Medallion Architecture is a **modular approach** to building data lakes using **Delta Lake** (often used with Azure Synapse or Databricks), promoting **data quality, governance, and reuse**.

It splits the data lifecycle into 3 layers:

**1. Bronze Layer (Raw Data)**

**Purpose:** Store raw, unfiltered data from source systems.

- Format: CSV, JSON, Avro, or raw Delta
- Minimal transformation (e.g., basic parsing)
- Used for auditing, reprocessing, or backup

**2. Silver Layer (Cleansed & Transformed)**

**Purpose:** Cleaned and structured data for analytical use.

- Joins, filters, type conversions
- Deduplication, null handling
- Typically stored in Delta Lake format for reliability

**3. Gold Layer (Business/Curated Data)**

**Purpose:** Final layer for analytics, ML, BI, and reporting.

- Aggregated KPIs
- Business-ready metrics
- Denormalized, fast access structure

The **Medallion Architecture** provides a **layered approach** to structuring data in a data lake or Delta Lake environment. It promotes **data quality, scalability, and clarity** by segmenting data into **bronze, silver, and gold layers**. Each layer serves a distinct purpose and follows best practices in terms of storage organization, access control, and data format.

## Explain RDD vs Data frames

**RDD (Resilient Distributed Dataset):**

It is low-level distributed data abstraction, which supports Java, Scala, Python languages.

**Characteristics:**
- Immutable, distributed collection of objects
- Can perform transformations (e.g., map, filter) and actions (count, collect)
- Offers **fine-grained control** over data

- Best for **unstructured data**, complex logic, or **low-level transformations**

**Downsides:**
- **No optimization**: No automatic query optimization (no Catalyst or Tungsten)
- **More code** to perform operations
- Less efficient than Data Frames
- Performance is low

**Data Frame:**

It is high-level abstraction over structured data, and this supports Java, Scala, Python and R languages.

**Characteristics:**
- Distributed **table-like** structure with rows and named columns (like a SQL table)
- Built on top of RDDs
- Supports **SQL-like operations** using DSL or SQL queries
- Backed by **Catalyst optimizer** (logical + physical optimization)
- Best for **structured/semi-structured data**, analytics, and aggregations

**Benefits:**
- **Optimized performance**
- **Less code** needed
- Supports various data sources: Parquet, JSON, JDBC, Hive, etc.
- Can interoperate with SQL and BI tools

# Explain narrow and wide transformations

**Narrow Transformations**

Narrow transformations are those where each partition of the parent RDD is used by at most one partition of the child RDD. In simple terms, data from a single input partition only maps to one output partition. These transformations do not require data to be shuffled across the network, making them faster and more efficient.

**Examples**: map(), filter(), flatMap(), union(), sample()

Since there's no data movement between nodes, narrow transformations allow Spark to pipeline operations and execute them more efficiently.

**Wide Transformations**

Wide transformations are those where multiple partitions of the parent RDD may contribute to a single partition of the child RDD. This requires data shuffling, where Spark redistributes data across the network. Shuffles are expensive operations as they involve disk I/O, data serialization, and network latency.

**Examples**: groupByKey(), reduceByKey(), join(), distinct()

These transformations are necessary for tasks that require data aggregation or combining datasets but can significantly impact performance if not handled properly.