# Data Engineering in the Cloud

## Hybrid Transactional Analytical Processing

Xuemao Zhang
East Stroudsburg University

January 18, 2025

# Outline

- In the labs, we Query Azure Cosmos DB with Apache Spark Pool and Serverless SQL Pools
- Lab 1: Querying Azure Cosmos DB with Apache Spark Pool
- Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

# Prerequisites

- Prerequisites (10-15 minutes):
  - Create a datalake (storage account with hierarchical namespace enabled)
    - ⋆ To lower the cost, you may choose `Redundancy` as `LRS`
    - ⋆ Create a container in the storage account
  - Create an Azure Synapse workspace and
    - ⋆ create a Apache spark pool in the `Manage` tab
  - Create an Azure Cosmos DB database with a containter
    - ⋆ with `Analytical store` enabled and then create two key-value items
  - In the Synapse Studio, go to the `Manage` tab to add a new linked service for `Azure Cosmos DB for NoSQL`
  - Go to the `Data` tab of the Synapse Studio, click on + and choose "Integration dataset" to create a dataset.
    - ⋆ We do not need the Copy pipeline

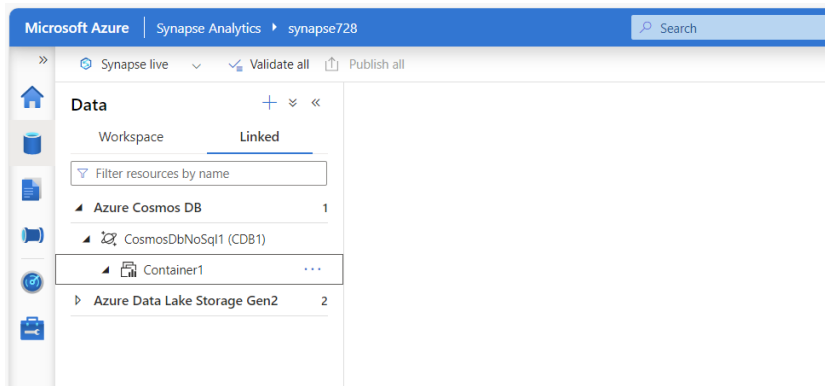# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool

- 1. Open "synapse studio", go to "Manage" tab, and add a new Cosmos DB `linked service`.
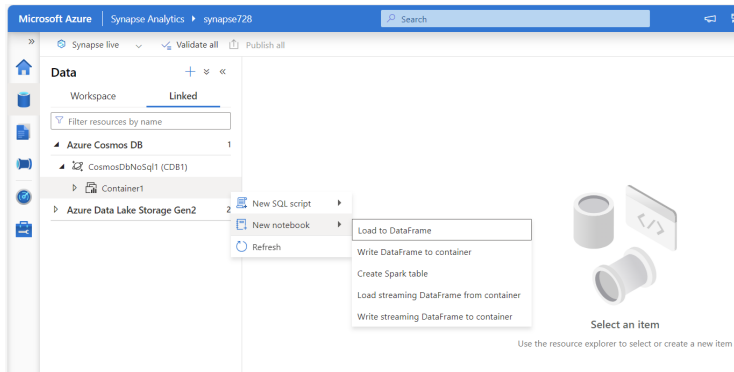
# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool
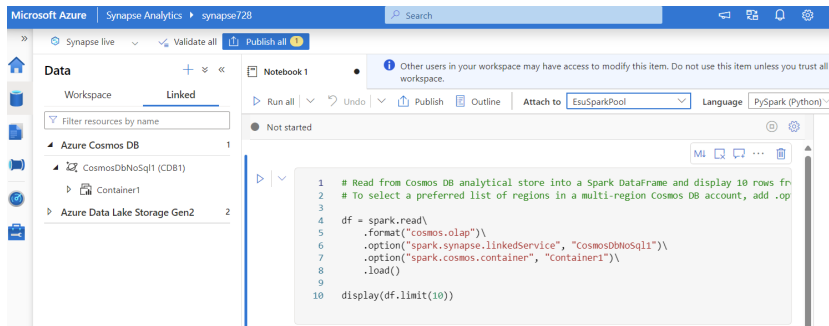
- ② Go to "Data" tab, and choose "Linked".

# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool

- - 3 Expand the Azure Cosmos DB, choose the container "Container1" created in the Cosmos DB, click on three dots, choose "New notebook", and choose "Load to Dataframe".

# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool
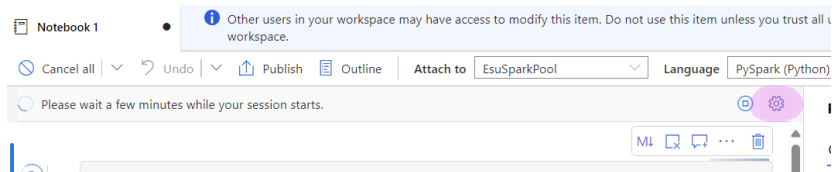
- 4 Attach the spark pool to the notebook.

# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool

- ⑤ Click on "Run All". It takes some time to run.
  - ▶ If you see similar message, then update the configuration of the spark session
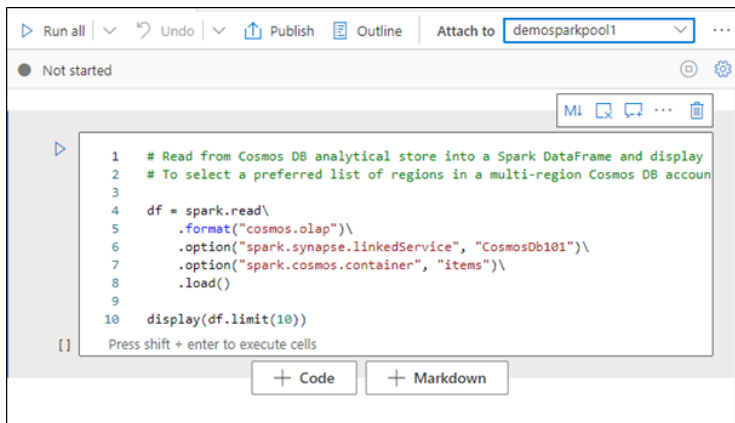


```
InvalidHttpRequestToLivy: Your Spark job requested 48 vcores.
However, the workspace has a 12 core limit. Try reducing the
numbers of vcores requested or increasing your vcore quota.
Quota can be increased using Azure Support request https://learn.microsoft.com/en-us
```

# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool

# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool

- ⑥ You must be able to see the output as given below:



```
1    # Read from Cosmos DB analytical store into a Spark DataFrame and display 10 rows from the DataFrame
2    # To select a preferred list of regions in a multi-region Cosmos DB account, add .option("spark.cosmos.preferredRegions", "<Region1>,<Region2>")
3
4    df = spark.read\
5        .format("cosmos.olap")\
6        .option("spark.synapse.linkedService", "CosmosDbNoSql1")\
7        .option("spark.cosmos.container", "Container1")\
8        .load()
9
10   display(df.limit(10))
```

✓ 3 min 28 sec - Apache Spark session started in 3 min 2 sec 344 ms. Command executed in 25 sec 972 ms by niuyuxia728 on 2:59:48 PM, 7/28/24

› **Job execution** Succeeded   **Spark** 1 executors 4 cores                        View in monitoring   Open Spark UI

View   Table   Chart        ↦ Export results ∨

| _rid | _ts | id | category | description | _etag |
|------|-----|----|----------|-------------|-------|
| 2J4oAMexroYBAAAAAAAAAA== | 1722190945 | 1 | fruits | vrious fruits | "06007145-0( |
| 2J4oAMexroYCAAAAAAAAAA== | 1722190981 | 2 | vegetables | various vegetables | "06005846-0( |

# Lab 1: Querying Azure Cosmos DB with Apache Spark Pool

- ⑦ You can remove the unwanted column using the below command:

```python
unwanted_cols = {'_rid','_ts','_etag'}
cols = list(set(df.columns) - unwanted_cols)
df2 = df.select(cols)

display(df2.limit(10))
```

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

- 1. Go to "Data" in the Synapse Studio and right click `Container1` to add an SQL script.

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

- ❷ Run the first block query see if a credential is needed.

```
IF (NOT EXISTS(SELECT * FROM sys.credentials WHERE name = 'esucosmos'))
    THROW 50000, 'As a prerequisite, create a credential with Azure Cosmos DB key in SECRET option:
    CREATE CREDENTIAL [esucosmos]
    WITH IDENTITY = ''SHARED ACCESS SIGNATURE'', SECRET = ''a4WXnMntqa0yJfMR4D84rcfNgDCDR33dfo7OPMeum8nOvNeshAEqkD0TdyycfyHUFyrDrQUAbV
GO
```

- ❸ Create a credential first

```
CREATE CREDENTIAL [esucosmos]
WITH IDENTITY = 'SHARED ACCESS SIGNATURE',
SECRET = 'a4WXnMntqa0yJfMR4D84rcfNgDCDR33dfo70
PMeum8nOvNeshAEqkD0TdyycfyHUFyrDrQUAbWUNACDb8oZYVg==';
```

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

- ④ Run the query:

```
SELECT TOP 100 *
FROM OPENROWSET( PROVIDER = 'CosmosDB',
                CONNECTION = 'Account=esucosmos;Database=CDB1',
                OBJECT = 'Container1',
                SERVER_CREDENTIAL = 'esucosmos'
) AS [Container1]
```

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

```
10
11   SELECT TOP 100 *
12   FROM OPENROWSET(PROVIDER = 'CosmosDB',
13                   CONNECTION = 'Account=esucosmos;Database=CDB1',
14                   OBJECT = 'Container1',
15                   SERVER_CREDENTIAL = 'esucosmos'
16   ) AS [Container1]
17
```

**Results**    Messages

View   [ Table ]   Chart        ↦ Export results ∨

🔍 Search

| description | category | _rid | _etag | _ts | id |
|---|---|---|---|---|---|
| vrious fruits | fruits | 2J4oAMexroYBAAA... | "06007145-000... | 1722190945 | 1 |
| various vegetables | vegetables | 2J4oAMexroYCAAA... | "06005846-000... | 1722190981 | 2 |

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

- ● ⑤ Now let's creat a database in Azure Synapse

```sql
-- Create the Profiles database if it does not exist
USE master;
GO
IF DB_ID(N'Profiles') IS NULL
BEGIN
    CREATE DATABASE Profiles;
END
GO

-- Switch to the Profiles database
USE Profiles;
GO

-- Drop the items view if it exists
DROP VIEW IF EXISTS items;
GO
```

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

```sql
-- Create the items view to access Cosmos DB
CREATE VIEW items
AS
SELECT
    *
FROM OPENROWSET(
PROVIDER = 'CosmosDB',
CONNECTION = 'Account=esucosmos;Database=CDB1',
OBJECT = 'Container1',
SERVER_CREDENTIAL = 'esucosmos'
)
WITH (
    id varchar(1000),
    category varchar(50),
    description varchar(max)
) AS profiles;
GO
```

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

- 6 Check the new database

```
USE Profiles;
GO

SELECT TOP 10 *
FROM items;
```

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

- 7 Check the new database: Refresh the Synapse Studio, then go to "Data" tab, expand the "Profiles (SQL)" and under "Views" there must be a view created.

# Lab 2: Querying Azure Cosmos DB with Serverless SQL Pools

- 8 Click on "dbo.items", choose "New SQL script", click "Select Top 100 rows" and you will be able to see the auto generated query. Run the query.

# License