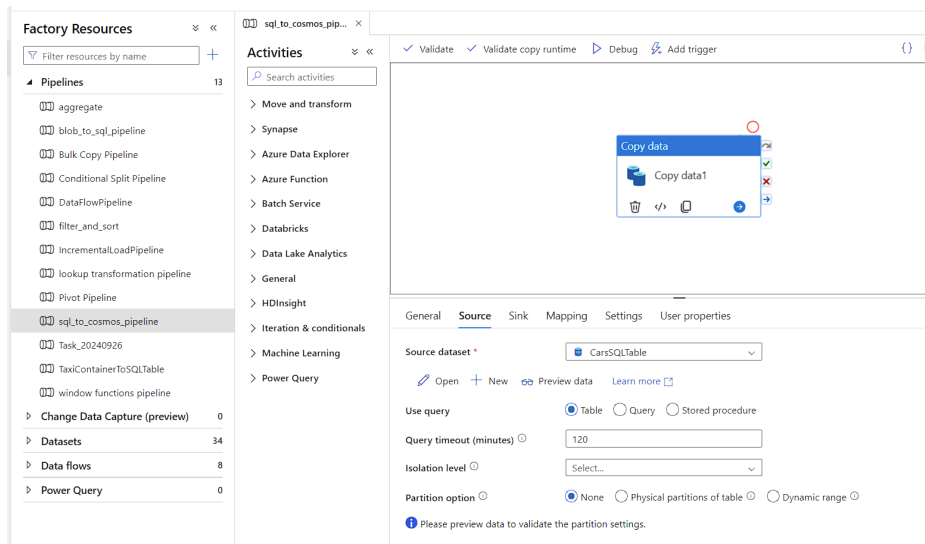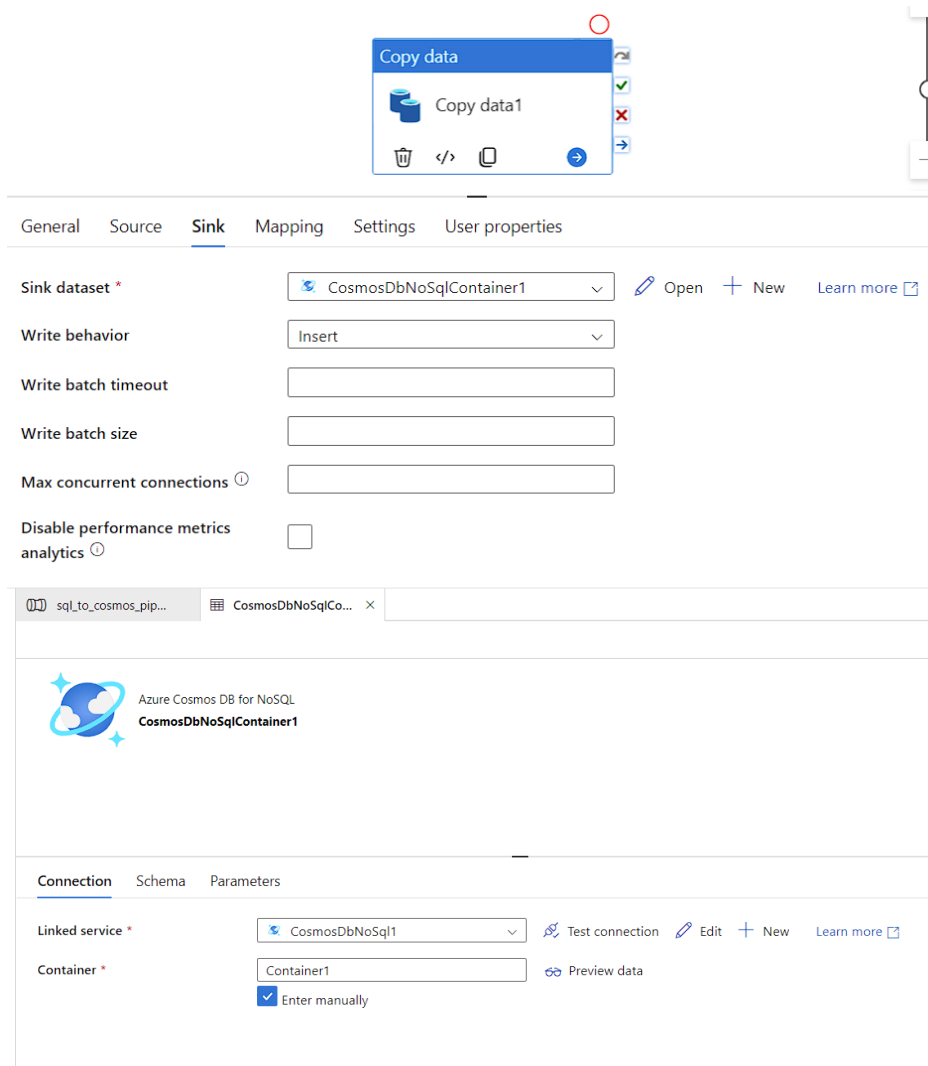# Assignment 10 - Azure Data Factory

**Q1. Design an ADF pipeline to copy data from an on-premise Azure SQL database to Azure Cosmos DB, ensuring data consistency and performance optimization. Pick correct options of partitioning for better performance.**

a. Main component of the pipeline is the copy data tool. Selected source is an SQL table named CarsSQLTable containing various details of different models of cars from multiple manufacturers.
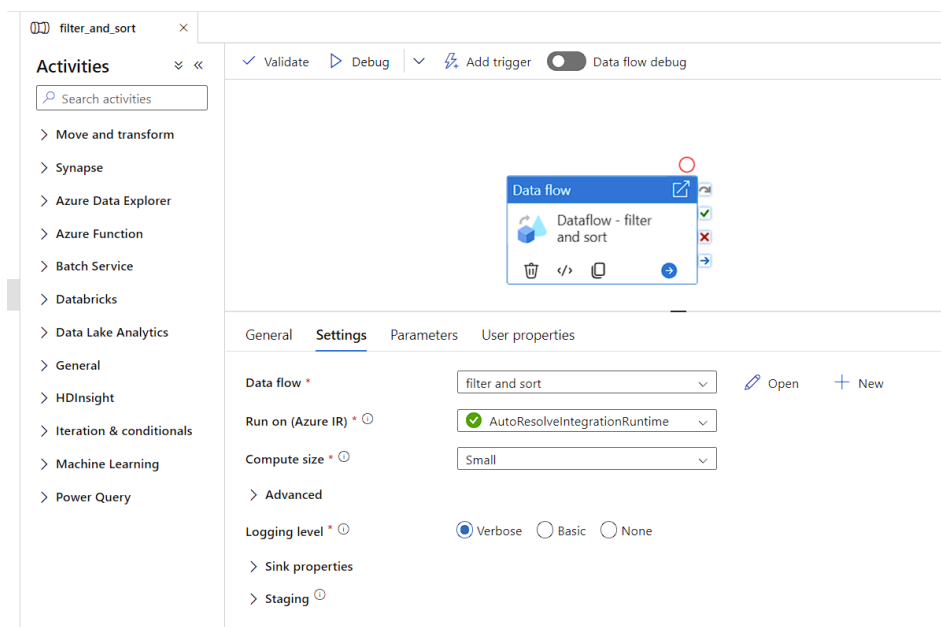


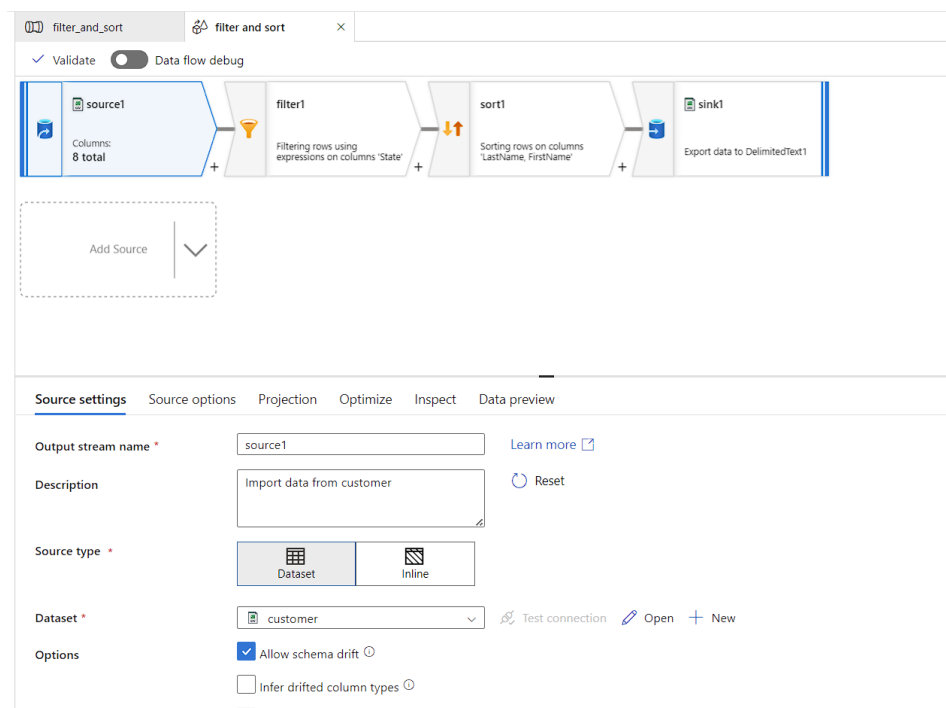b. Sink for the Copy Data tool is a Container in Azure CosmosDB.

| General | Source | Sink | Mapping | Settings | User properties |
|---------|--------|------|---------|----------|-----------------|

**Sink dataset** *          [CosmosDbNoSqlContainer1 ▾]    ✎ Open    + New    Learn more ⧉

**Write behavior**          [Insert ▾]

**Write batch timeout**     [                    ]

**Write batch size**        [                    ]

**Max concurrent connections** ⓘ  [                    ]

**Disable performance metrics analytics** ⓘ  ☐

| 🗐 sql_to_cosmos_pip... | 🖽 CosmosDbNoSqlCo... ✕ |
|---|---|

Azure Cosmos DB for NoSQL
**CosmosDbNoSqlContainer1**

| Connection | Schema | Parameters |
|------------|--------|------------|

**Linked service** *    [CosmosDbNoSql1 ▾]    ⚡ Test connection    ✎ Edit    + New    Learn more ⧉

**Container** *         [Container1          ]    👓 Preview data
                        ☑ Enter manually

Multiple partitioning options are availlable while choosing source - here, Noone is choosen due to the smaller size of availabe data which will not provide significant performance benefits. For bigger datasets, using partitions will help provide performance benefits by parallelizing data extraction

**Q2. Create Pipeline using Azure Data Flow in Azure Data Factory to apply Filter and Sort transformations on datasets.**

a. Main component of the pipeline is a dataflow task.

b. Dataflow first involves taking the data from a FlatFile source named Customers.



c. Filter involves taking only rows from the data where State equals 'TX'

**Filter settings**   Optimize   Inspect   Data preview

| | | |
|---|---|---|
| Output stream name * | filter1 | Learn more ↗ |
| Description | Filtering rows using expressions on columns 'State' | ↻ Reset |
| Incoming stream * | source1 | |
| Filter on * | State=="TX" | |

d. Sort involves sorting the output from Filter step based on descending order of LastName and FirstName columns



**Sort settings**   Optimize   Inspect   Data preview

| | | |
|---|---|---|
| Output stream name * | sort1 | Learn more ↗ |
| Description | Sorting rows on columns 'LastName, FirstName' | ↻ Reset |
| Incoming stream * | filter1 | |
| Options * | ☐ Case insensitive | |
| | ☐ Sort only within partition | |

Sort conditions *

| filter1's column | Order | Nulls first | | |
|---|---|---|---|---|
| abc LastName | Descending | ☑ | + | 🗑 |
| abc FirstName | Descending | ☑ | + | 🗑 |

e. Final destination is a CSV file.



**Sink**   Settings   Errors   Mapping   Optimize   Inspect   Data preview

| | | |
|---|---|---|
| Output stream name * | sink1 | Learn more ↗ |
| Description | Export data to DelimitedText1 | ↻ Reset |
| Incoming stream * | sort1 | |
| Sink type * | ▦ Dataset   ▨ Inline   ▤ Cache | |
| Dataset * | DelimitedText1   Test connection   ✎ Open   + New | |
| Skip line count | | |
| Options | ☑ Allow schema drift ⓘ | |
| | ☐ Validate schema ⓘ | |

**Q3. Design an ADF pipeline to implement aggregate operations, such as sum, average, max, min and count, within an Azure Data Flow.**

a. Main component of PIpeline is a dataflow task.



b. Data source here is a CSV file named Cars. Change datatype of each column to the appropriate type from the projection section.

## Source settings | Source options | Projection | Optimize | Inspect | Data preview

| | |
|---|---|
| Output stream name * | source1 |
| Description | Import data from Cars |
| Source type * | Dataset / Inline |
| Dataset * | Cars |
| Options | ☑ Allow schema drift ⓘ |
| | ☐ Infer drifted column types ⓘ |
| | ☐ Validate schema ⓘ |
| Skip line count | |
| Sampling * ⓘ | ◯ Enable  ● Disable |

Learn more ↗

↻ Reset

Test connection    Open    + New

Define default format    Detect data type    Import projection    Reset schema

| Column name | Type | Format |
|---|---|---|
| Make | abc string | Specify format |
| Model | abc string | Specify format |
| Type | abc string | Specify format |
| Origin | abc string | Specify format |
| DriveTrain | abc string | Specify format |
| MSRP | t.2f float | Specify format |
| Invoice | abc string | Specify format |
| EngineSize | abc string | Specify format |
| Cylinders | 123 integer | Specify format |
| Horsepower | t.2f float | Specify format |
| MPG_City | t.2f float | Specify format |
| MPG_Highway | t.2f float | Specify format |
| Weight | t.2f float | Specify format |
| Wheelbase | t.2f float | Specify format |
| Length | t.2f float | Specify format |

c. Aggregating data by 'Make' producing columns 'avg_MSRP, avg_MPG_City, avg_MPG_Highway, stddev_Weight'

**Aggregate settings**   Optimize   Inspect   Data preview

| | | |
|---|---|---|
| Output stream name * | aggregate1 | Learn more 🔗 |
| Description | Aggregating data by 'Make' producing columns 'avg_MSRP, avg_MPG_City, avg_MPG_Highway, stddev_Weight' | ↻ Reset |
| Incoming stream * | source1 | |

[ Group by ] [ Aggregates ]

**Columns**                          **Name as**

| abc  Make | Make | ➕ 🗑️ |
|---|---|---|

---

**Aggregate settings**   Optimize   Inspect   Data preview

| | | |
|---|---|---|
| Output stream name * | aggregate1 | Learn more 🔗 |
| Description | Aggregating data by 'Make' producing columns 'avg_MSRP, avg_MPG_City, avg_MPG_Highway, stddev_Weight' | ↻ Reset |
| Incoming stream * | source1 | |

[ Group by ] [ Aggregates ]

Grouped by: Make

➕ Add    ⧉ Clone    🗑️ Delete    ⧉ Open expression builder

| | Column | Expression | |
|---|---|---|---|
| ☐ | avg_MSRP | avg(MSRP) 1.2 | ➕ 🗑️ |
| ☐ | avg_MPG_City | avg(MPG_City) 1.2 | ➕ 🗑️ |
| ☐ | avg_MPG_Highway | avg(MPG_Highway) 1.2 | ➕ 🗑️ |
| ☐ | stddev_Weight | stddev(Weight) 1.2 | ➕ 🗑️ |

## d. Final output is stored to a CSV file

| ▭ filter_and_sort | ⧉ filter and sort | ▭ aggregate | ⧉ aggregate1 ✕ | |

✓ Validate   ⬤ Data flow debug

| 📄 source1 | | aggregate1 | | 📄 sink1 |
|---|---|---|---|---|
| Import data from Cars | Σ⧉ | Aggregating data by 'Make' producing columns 'avg_MSRP, avg_MPG_City, avg_MPG_Highway,' | ➡ | Columns: 5 total |

**Sink**   Settings   Errors   Mapping   Optimize   Inspect   Data preview

| | | |
|---|---|---|
| Output stream name * | sink1 | Learn more 🔗 |
| Description | Export data to DelimitedText2 | ↻ Reset |
| Incoming stream * | aggregate1 | |

| Sink type * | | | |
|---|---|---|---|
| | ▦ Dataset | ▨ Inline | ⧉ Cache |

| Dataset * | 📄 DelimitedText2 | 🔌 Test connection   ✏️ Open   ➕ New |
|---|---|---|

| Skip line count | | |
|---|---|---|

| Options | ☑ Allow schema drift ⓘ |
|---|---|
| | ☐ Validate schema ⓘ |

**4. Create best approach to bulk copy data from multiple homogenous sources into Azure SQL Database using ADF pipelines. Show usage of Lookup, For Each Loop and Expressions in Azure Data Factory.**

a. Lookup component uses a query to get list of all tables in the specified database. The source dataset has only a linked service and no table connection as it is used just to have a connection the the database.



b. ForEach loop gives the name of the tables to the 'Export Table' activity inside it using the expression `@activity('List tables').output.value`

c. In the Export Table activity, source dataset has two properties - TableName and SchemaName. These were used to access each table in database.



d. Sink dataset is a blob storage, where we store as records of each table in CSV files. Dataset has a property named TableName with value set to `@concat(item().table_schema, '_', item().table_name, '.csv')`
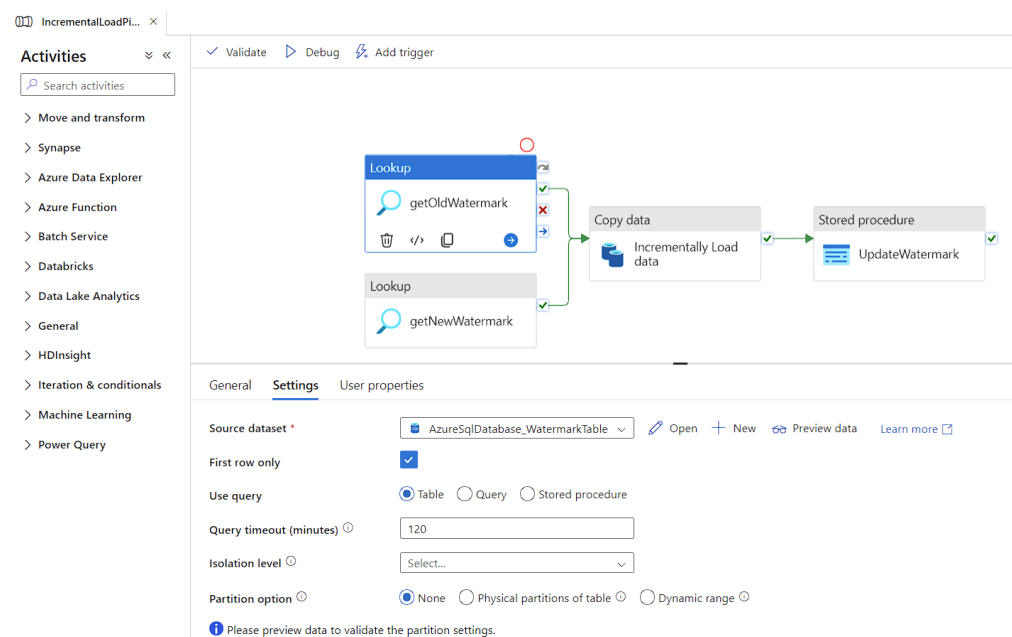
## 5. Implement incremental load Pipeline in Azure Data Factory for handling datasets, ensuring efficient insert/upsert/updates to the target storage without re-inserting the entire dataset?

## a. Main pipeline

b. The lookup activity named 'getOldWatermark' is used to retreive latest lookup value stored in the watermark table- which signifies when incremental loading last took place.



c. The lookup activity 'getNewWatermark' uses a query to retreive the time when Source Dataset was last modified.

d. Copy Data tool uses an SQL query inorder to retreive al rows where watermark value is greater than old watermark, but less than or equals newer watermark

## Pipeline expression builder

Add dynamic content below using any combination of expressions, functions and system variables.

```
SELECT * FROM data_source_table
WHERE
    LastModifytime > '@{activity('getOldWatermark').output.
        firstRow.WatermarkValue}'
    AND
    LastModifytime <= '@{activity('getNewWatermark').output.
        firstRow.NewWatermarkvalue}'
```

Clear contents

e. THe sink dataset is the destination table to which we are incrementally loading data. Upsert option is used to continuously input new data without getting overridden. Primary key for the table is employee_id.

| General | Source | **Sink** | Mapping | Settings | User properties |

| Sink dataset * | | AzureSqlDatabase_DestinationTable ∨ | ✎ Open | + New | Learn more 🔗 |

Write behavior   ◯ Insert   ⦿ Upsert   ◯ Stored procedure

Use TempDB ⓘ   ☑

Key columns ⓘ

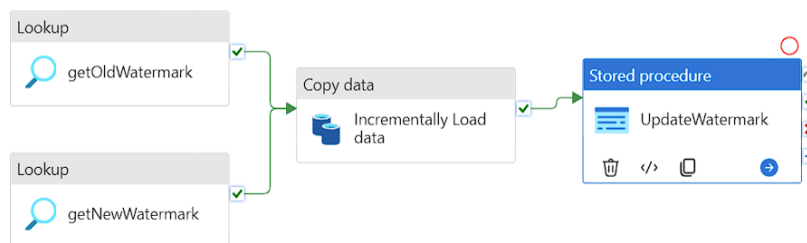+ New   |   🗑 Delete   ↻ Refresh

☐ Column name

☐ 123 employee_id ∨

Bulk insert table lock ⓘ   ◯ Yes   ⦿ No

Table option   ⦿ Use existing   ◯ Auto create table ⓘ

f. A stored procedure is used to write new watermark value and table name to the watermark table. 'LastModifiedtime' from getNewWatermark activity and 'TableName' are passed are parameters to the stored procedure.



| General | **Settings** | User properties |

Linked service * ⓘ   EmployeesSQLSource ∨   ⚡ Test connection   ✎ Edit   + New

Stored procedure name *   [dbo].[usp_write_watermark]
☑ Enter manually

∨ Stored procedure parameters ⓘ

⟵ Import   + New   |   🗑 Delete

| | Name | Type | Value |
|---|---|---|---|
| ☐ | LastModifiedtime | DateTime ∨ | @{activity('getNewWatermark').outp... |
| ☐ | TableName | String ∨ | @{activity('getOldWatermark').outpu... |

**6. What are the key steps to connect Azure Databricks to Cosmos DB for real-time analytics and data transformation using spark and Databricks.**

a. Transfer data from source to CosmonDB storage. Source here is a CSV file in blob storage, that needs to be moved to Cosmos using the Copy Data Tool.

Task_20240926 ✕

Activities ⌄ «

🔍 Search activities

> Move and transform
> Synapse
> Azure Data Explorer
> Azure Function
> Batch Service
> Databricks
> Data Lake Analytics
> General
> HDInsight
> Iteration & conditionals
> Machine Learning
> Power Query

✓ Validate   ✓ Validate copy runtime   ▷ Debug   ⚡ Add trigger

Copy data
Blob to CosmosDB

Notebook
CarsNotebook

General   **Source**   Sink   Mapping   Settings   User properties

Source dataset *          📄 Cars          ✏️ Open   + New   👓 Preview data   Learn more ⬏

File path type            ● File path in dataset   ○ Prefix   ○ Wildcard file path   ○ List of files ⓘ

                          Start time (UTC)                      End time (UTC)
Filter by last modified ⓘ  [                    ]               [                    ]

Recursively ⓘ             ☑

Enable partitions discovery ⓘ   ☐

Max concurrent connections ⓘ   [                    ]

Skip line count            [                    ]

Additional columns ⓘ       + New

---

Task_20240926 ✕

Activities ⌄ «

🔍 Search activities

> Move and transform
> Synapse
> Azure Data Explorer
> Azure Function
> Batch Service
> Databricks
> Data Lake Analytics
> General
> HDInsight
> Iteration & conditionals
> Machine Learning
> Power Query

✓ Validate   ✓ Validate copy runtime   ▷ Debug   ⚡ Add trigger

Copy data
Blob to CosmosDB

Notebook
CarsNotebook

General   Source   **Sink**   Mapping   Settings   User properties

Sink dataset *            🗄 CosmosDbNoSqlContainer2          ✏️ Open   + New   Learn more ⬏

Write behavior            [ Insert                            ⌄]

Write batch timeout       [                                   ]

Write batch size          [                                   ]

Max concurrent connections ⓘ   [                            ]

Disable performance metrics
analytics ⓘ               ☐

b. Create a linked service for connecting to databricks notebook. Go to Settings section to connect to a notebook.

c. In the databricks notebook, first install driver to connect the cluster to CosmosDB. Define endpoint URL, primary key, database name and container name. Read from cosmos as given below.

```python
spark.conf.set("spark.cosmos.accountEndpoint", cosmos_endpoint)
spark.conf.set("spark.cosmos.accountKey", cosmos_master_key)
spark.conf.set("spark.cosmos.database", database_name)
spark.conf.set("spark.cosmos.container", container_name)

df = spark.read.format("cosmos.oltp") \
    .option("spark.cosmos.accountEndpoint", cosmos_endpoint) \
    .option("spark.cosmos.accountKey", cosmos_master_key) \
    .option("spark.cosmos.database", database_name) \
    .option("spark.cosmos.container", container_name) \
    .load()

df.show()
```