

# Assignment 10 - Azure Data Factory

1. 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.

The screenshot shows the Microsoft Azure Data Factory interface. On the left, the 'Factory Resources' sidebar lists various pipelines and datasets. In the center, a pipeline named 'CarsDataflow' is selected. The 'Copy data' activity is highlighted. The 'General' tab of the activity configuration pane is visible, showing the activity is activated and has a timeout of 0:12:00. Below the general tab, the 'Source' tab is selected, showing the source dataset is 'CarsSqlTable'. Other tabs include 'Sink', 'Mapping', 'Settings', and 'User properties'.

This screenshot is identical to the one above, but the 'Source' tab is explicitly highlighted. It shows the 'Source dataset' is set to 'CarsSqlTable'. Other settings visible include 'Use query' (set to 'Table'), 'Query timeout (minutes)' (set to 120), and 'Partition option' (set to 'None'). A note at the bottom says 'Please preview data to validate the partition settings.'

Factory Resources

Filter resources by name

Azure SQL Database  
CarsSqlTable

Connection Schema Parameters

Linked service \* AzureSqlOutputDB Test connection Edit + New Learn more

Table dbo.Cars Refresh Preview data Enter manually

Factory Resources

Filter resources by name

CarsSqlTable

Copy data

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

Factory Resources

Filter resources by name

Azure Cosmos DB for NoSQL  
CosmosDbNoSqlContainer1

Connection Schema Parameters

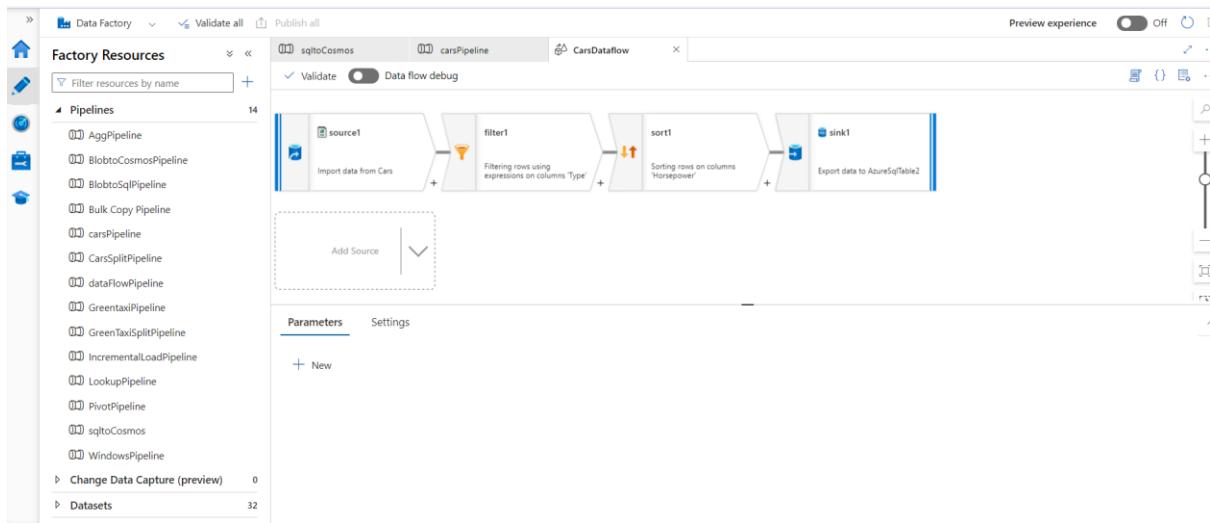
Linked service \* CosmosDbNoSql1 Test connection Edit + New Learn more

Container \* MyContainer Preview data Enter manually

The screenshot shows the 'Pipeline runs' section of the Azure Data Factory interface. On the left, a sidebar lists navigation options: Dashboards, Runs, Pipeline runs (selected), Trigger runs, Change Data Capture (preview...), Runtimes & sessions, Integration runtimes, Data flow debug, Notifications, and Alerts & metrics. The main area displays a summary of 'All pipeline runs' for the 'weather - Activity runs' pipeline. A single activity run, 'Copy data', is listed with a green checkmark indicating it succeeded. The run details show it started at 9/30/2024, 11:23:04 AM, took 1m 49s, and was run by 'AutoResolveIntegration'. There are buttons for 'Rerun', 'Cancel', 'Refresh', 'Update pipeline', and tabs for 'List' and 'Gantt'.

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

The screenshot shows the 'Pipelines' section of the Azure Data Factory interface. On the left, a sidebar lists 'Factory Resources' under 'Pipelines', including 'AggPipeline', 'BlobtoCosmosPipeline', 'BlobtoSqlPipeline', 'Bulk Copy Pipeline', 'carsPipeline' (selected), 'CarsSplitPipeline', 'dataFlowPipeline', 'GreentaxiPipeline', 'GreenTaxiSplitPipeline', 'IncrementalLoadPipeline', 'LookupPipeline', 'PivotPipeline', 'sqloCosmos', and 'WindowsPipeline'. The main area shows a 'Data flow' configuration for 'CarsDataflow'. The 'General' tab is selected, displaying fields for 'Name' (set to 'CarsDataflow'), 'Description' (empty), 'Activity state' (set to 'Activated'), 'Timeout' (set to '0:12:00'), and 'Retry' (set to '0'). There are tabs for 'Settings', 'Parameters', and 'User properties'.



**Source settings**

**Output stream name**: source1

**Description**: Import data from Cars

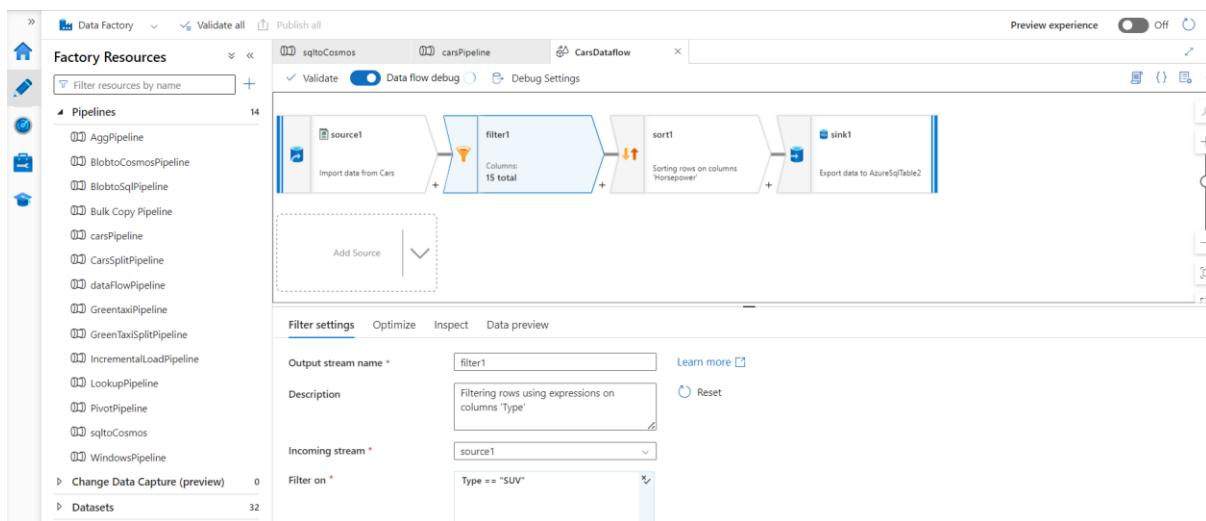
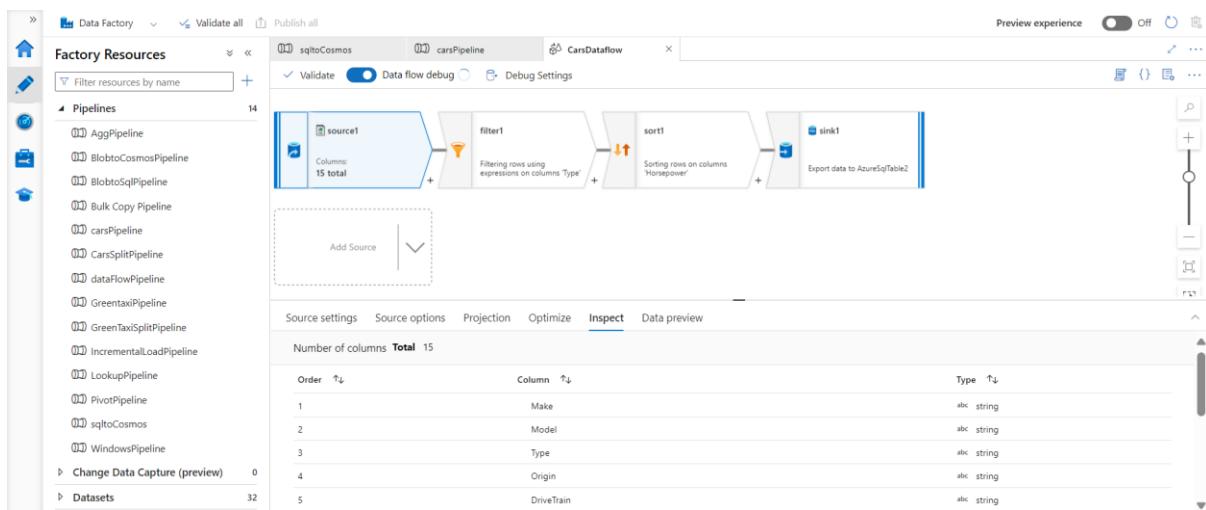
**Source type**: Dataset

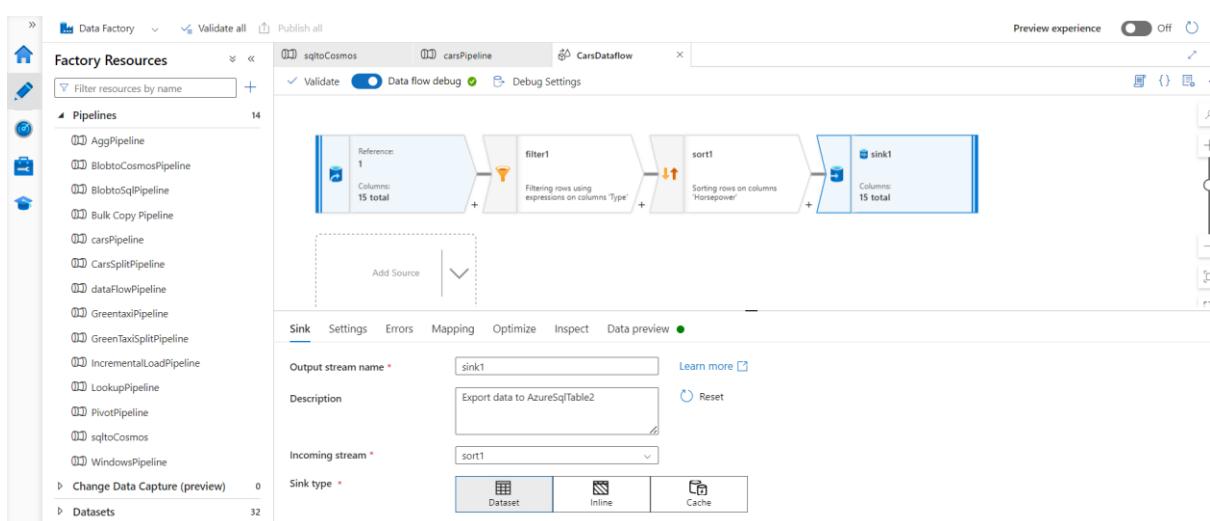
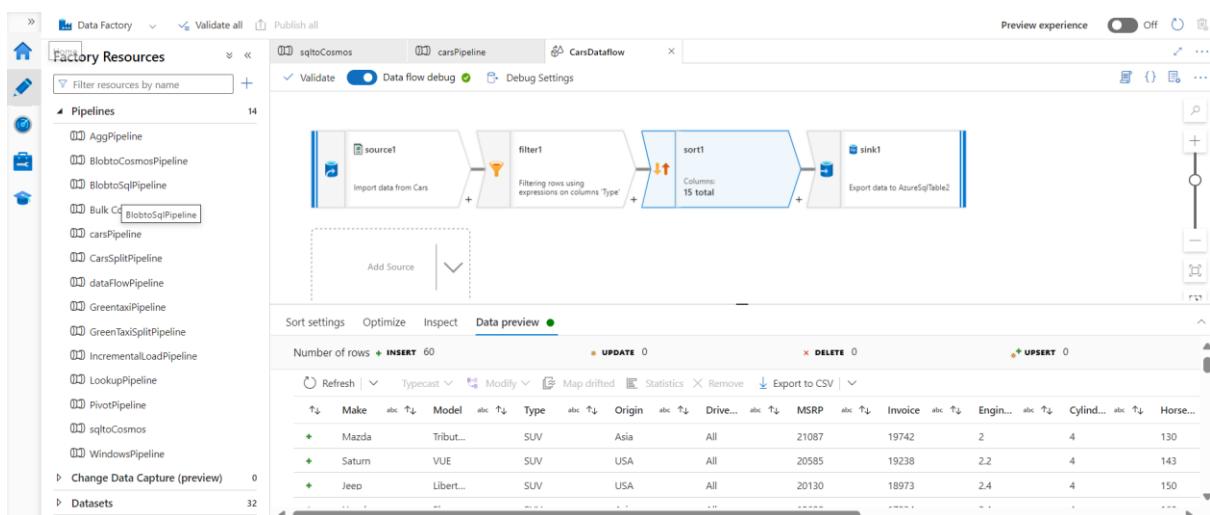
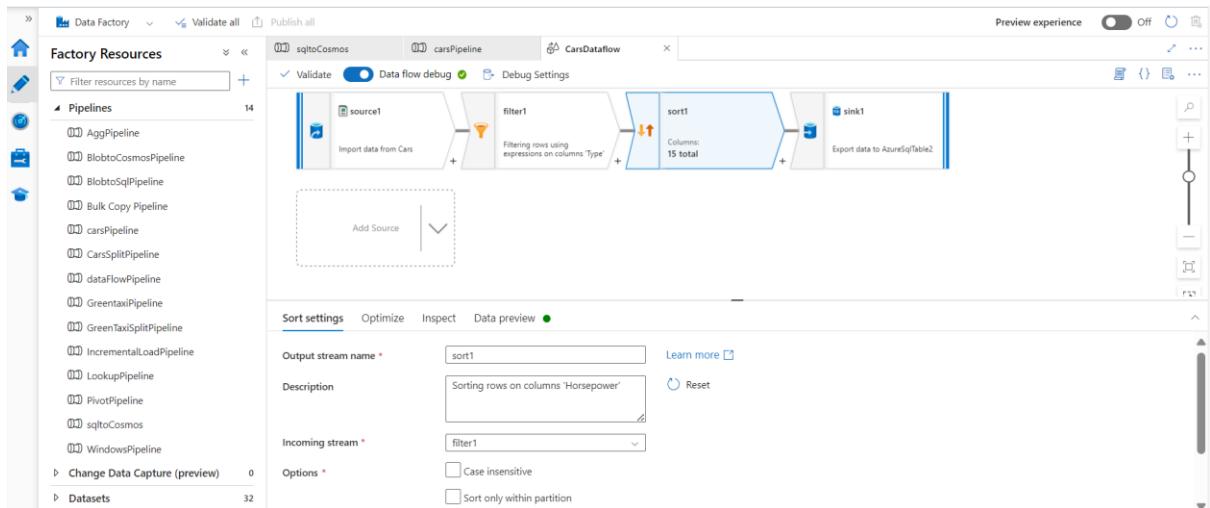
**Dataset**: Cars

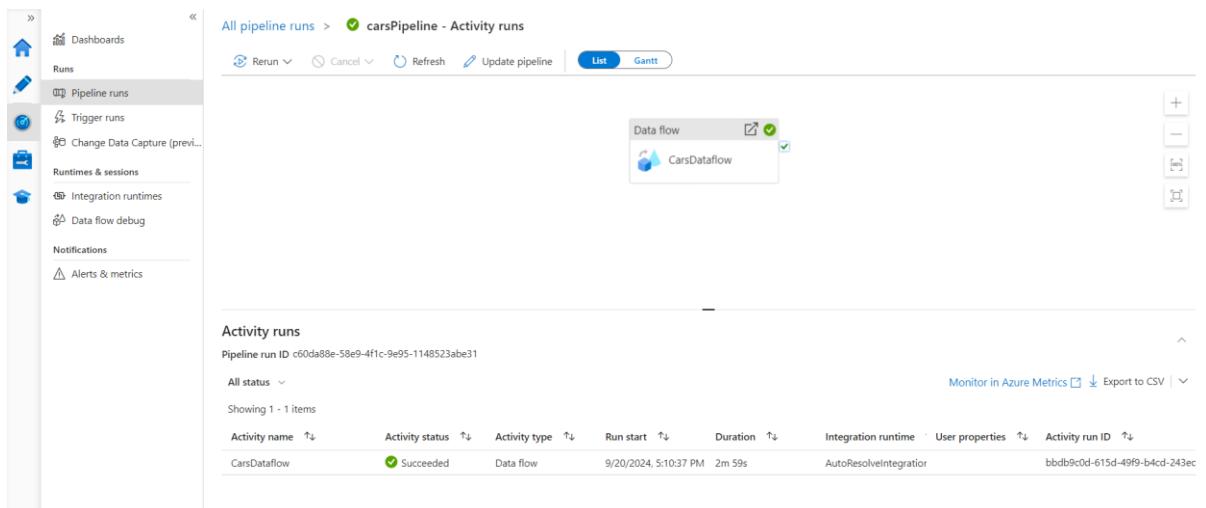
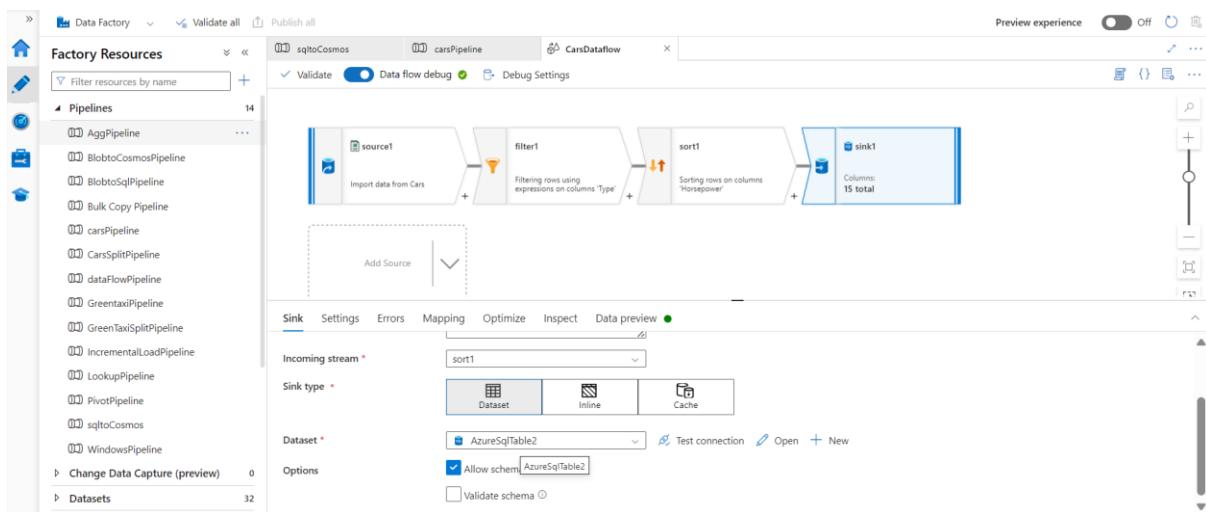
**Source options**, **Projection**, **Optimize**, **Inspect**, **Data preview**

**Projection**

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	abc string	Specify format







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

Factory Resources

- Pipelines
  - AggPipeline
  - BlobtoCosmosPipeline
  - BlobtoSqlPipeline
  - Bulk Copy Pipeline
  - carsPipeline
  - CarsSplitPipeline
  - dataFlowPipeline
  - GreentaxiPipeline
  - GreenTaxiSplitPipeline
  - IncrementalLoadPipeline
  - LookupPipeline
  - PivotPipeline
  - sqltoCosmos
  - WindowsPipeline
- Change Data Capture (preview) 0
- Datasets 32

Preview experience: Off

Data flow

Aggdataflow

Parameters Variables Settings Output

+ New

Factory Resources

- Pipelines
  - AggPipeline
  - BlobtoCosmosPipeline
  - BlobtoSqlPipeline
  - Bulk Copy Pipeline
  - carsPipeline
  - CarsSplitPipeline
  - dataFlowPipeline
  - GreentaxiPipeline
  - GreenTaxiSplitPipeline
  - IncrementalLoadPipeline
  - LookupPipeline
  - PivotPipeline
  - sqltoCosmos
  - WindowsPipeline
- Change Data Capture (preview) 0
- Datasets 32

Validate Data flow debug Debug Settings

Data flow

source1 aggregate1 sink

Add Source

Source settings

Output stream name: source1

Description: Import data from CarsSqlTable

Source type: Dataset

Dataset: CarsSqlTable

Options: Allow schema drift

Factory Resources

- Pipelines
  - AggCars
  - AzureSqlTable1
  - AzureSqlTable2
  - Cars
  - CarSplit
  - CarSplit1
  - CarSplit2
  - CarSplit3
  - CarSplit4
  - CarsSqlTable
  - CopyTablestocsv
  - CosmosDbNoSqlContainer1
  - CosmosDbNoSqlContainer2
  - CSVfiles
  - customers
  - DelimitedText1
  - DelimitedText2

sqltoCosmos carsPipeline CarsDataflow AggPipeline Aggdataflow CarsSqlTable

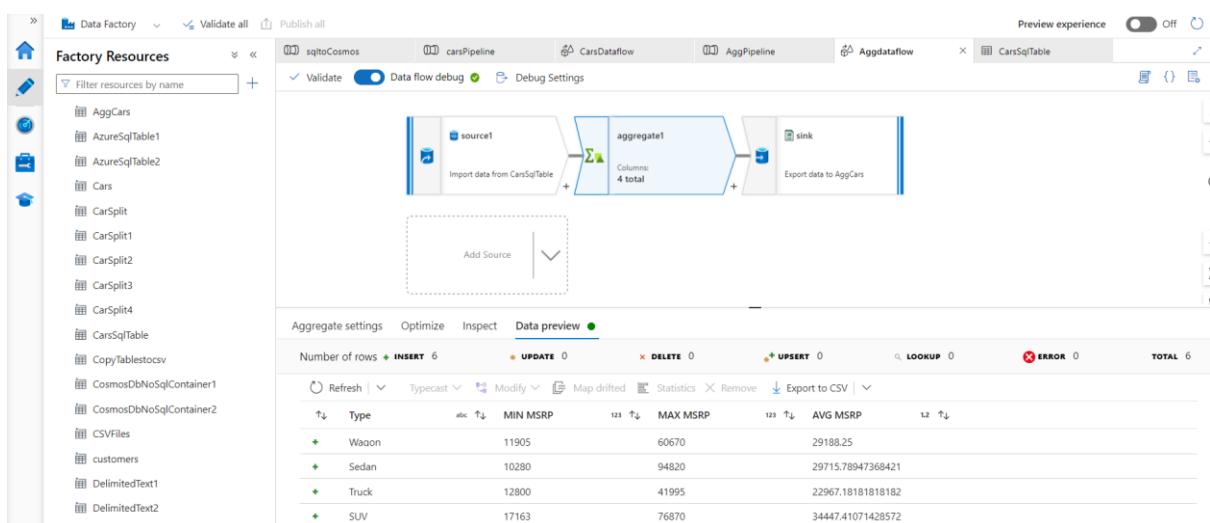
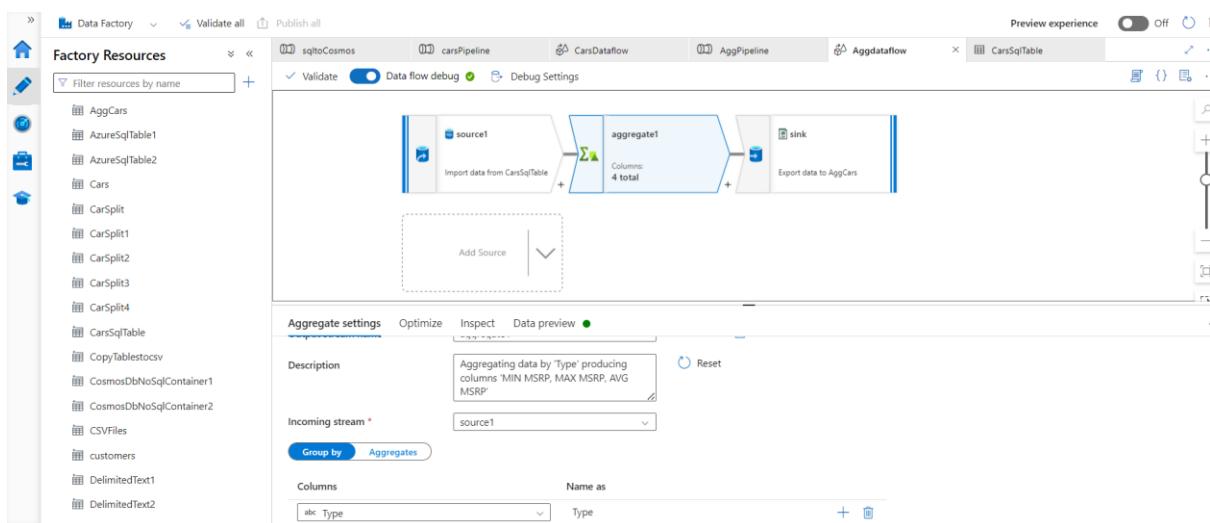
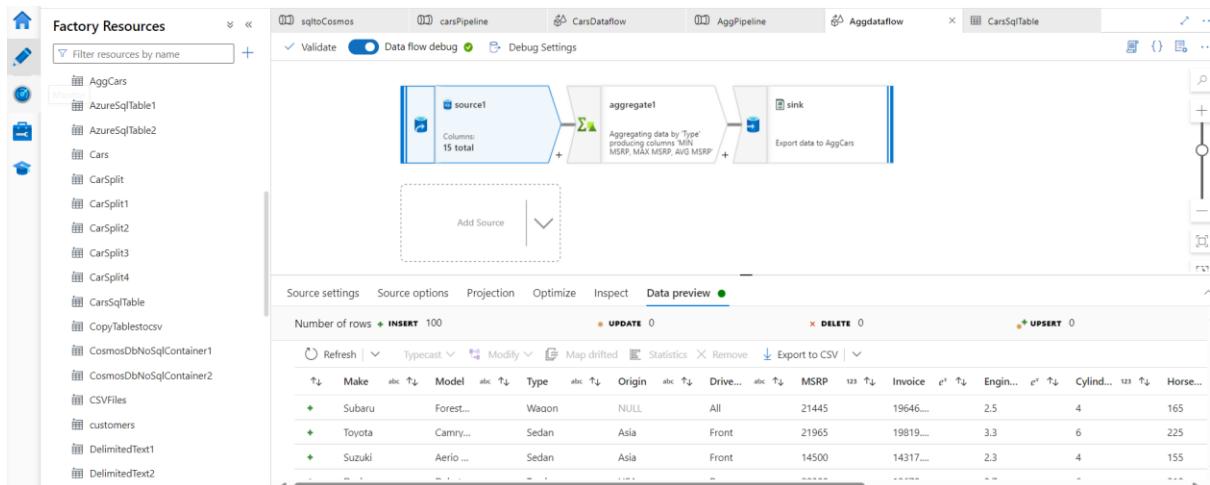
Azure SQL Database

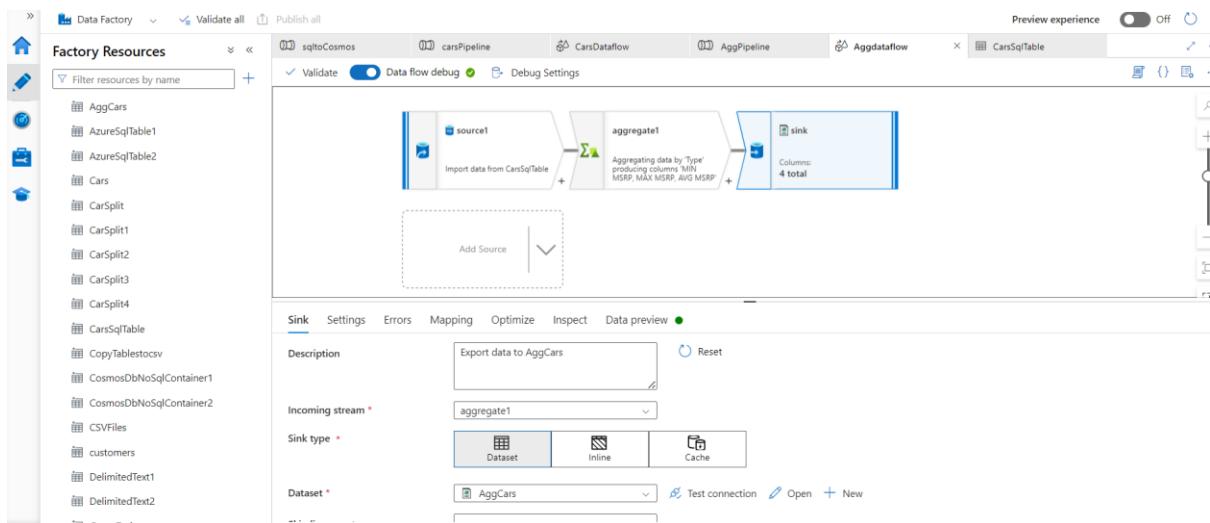
CarsSqlTable

Connection Schema Parameters

Linked service: AzureSqlOutputDB

Table: dbo.Cars





Factory Resources

Validate all Publish all

sqtoCosmos carsPipeline CarsDataflow AggPipeline Aggdataflow CarsSqlTable

Preview experience Off

Number of rows: 6

INSERT N/A UPDATE N/A DELETE N/A UPSERT N/A LOOKUP N/A ERROR N/A TOTAL 6

Type	MIN MSRP	MAX MSRP	AVG MSRP
Wagon	11905	60670	29188.25
Sedan	10280	94820	29715.78947368421
Truck	12800	41995	22967.18181818182
SUV	17163	76870	34447.41071428572

The screenshot shows the 'All pipeline runs' view for the 'AggPipeline - Activity runs' section. On the left, a sidebar lists navigation options: Dashboards, Runs, Pipeline runs (selected), Trigger runs, Change Data Capture (previous versions), Runtimes & sessions, Integration runtimes, Data flow debug, Notifications, and Alerts & metrics. The main area displays a summary of the pipeline run ID: 8d9acdfb-2da4-4ab4-ba35-80d59ed9d0a1. Below this, a table titled 'Activity runs' shows one item: 'Aggdataflow' with status 'Succeeded'. The table includes columns for Activity name, Activity status, Activity type, Run start, Duration, Integration runtime, User properties, and Activity run ID.

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.

The screenshot shows the Azure Data Factory Pipeline designer for a 'Bulk Copy Pipeline'. The pipeline consists of two main activities: a 'Lookup' activity followed by a 'ForEach' activity. The 'Lookup' activity is configured to query the 'ListTablesSql' dataset using a dynamic query. The 'ForEach' activity is configured to loop over all tables ('LoopAllTables') and perform a 'Copy data' operation for each table. The pipeline is currently in 'Validate' mode.

Factory Resources

Bulk Copy Pipeline

ListTablesSql

Azure SQL Database

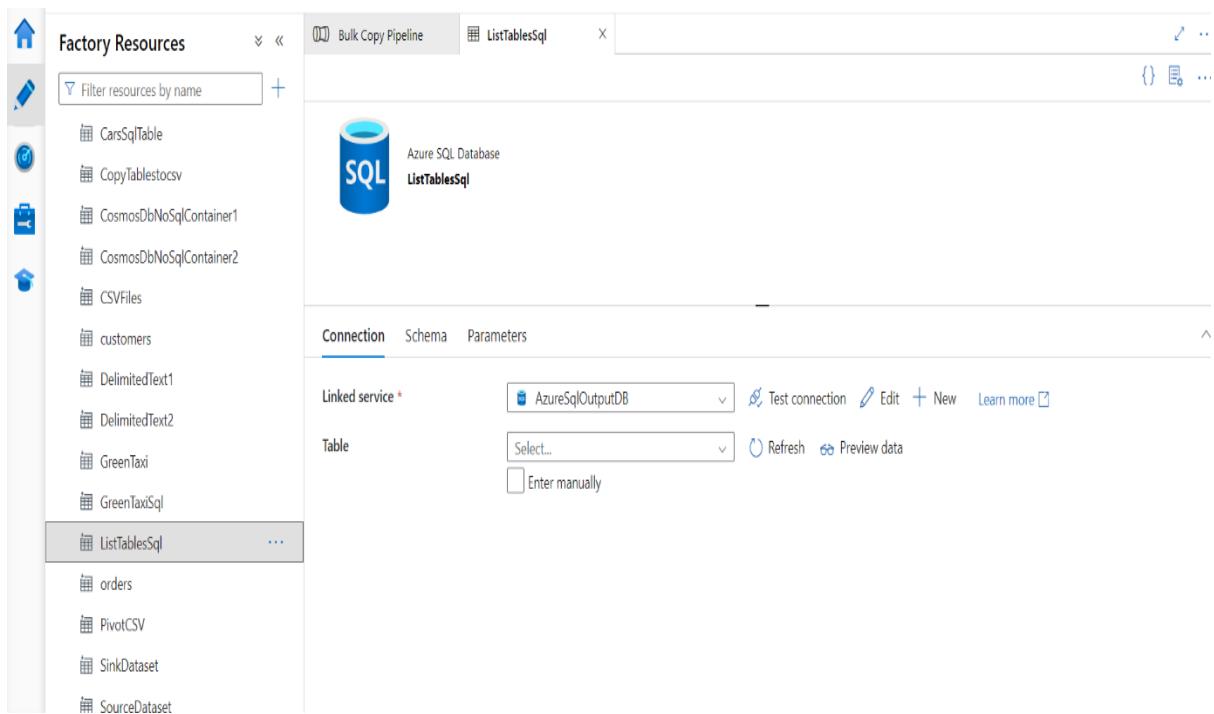
ListTablesSql

Connection Schema Parameters

Linked service \* AzureSqlOutputDB Test connection Edit New Learn more

Table Select... Refresh Preview data Enter manually

CarsSqlTable CopyTablestocsv CosmosDbNoSqlContainer1 CosmosDbNoSqlContainer2 CSVFiles customers DelimitedText1 DelimitedText2 GreenTaxi GreenTaxiSql ListTablesSql orders PivotCSV SinkDataset SourceDataset



Data Factory Validate all Publish all

Factory Resources Filter resources by name

Bulk Copy Pipeline ListTablesSql

Validate Debug Add trigger

ForEach

Lookup List tables

Activities

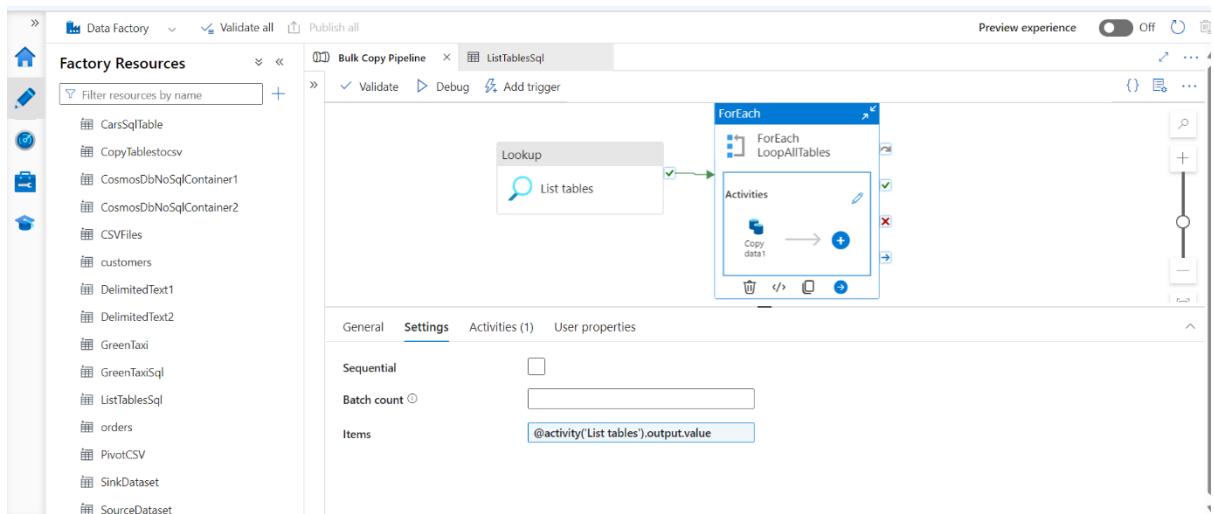
Copy data1

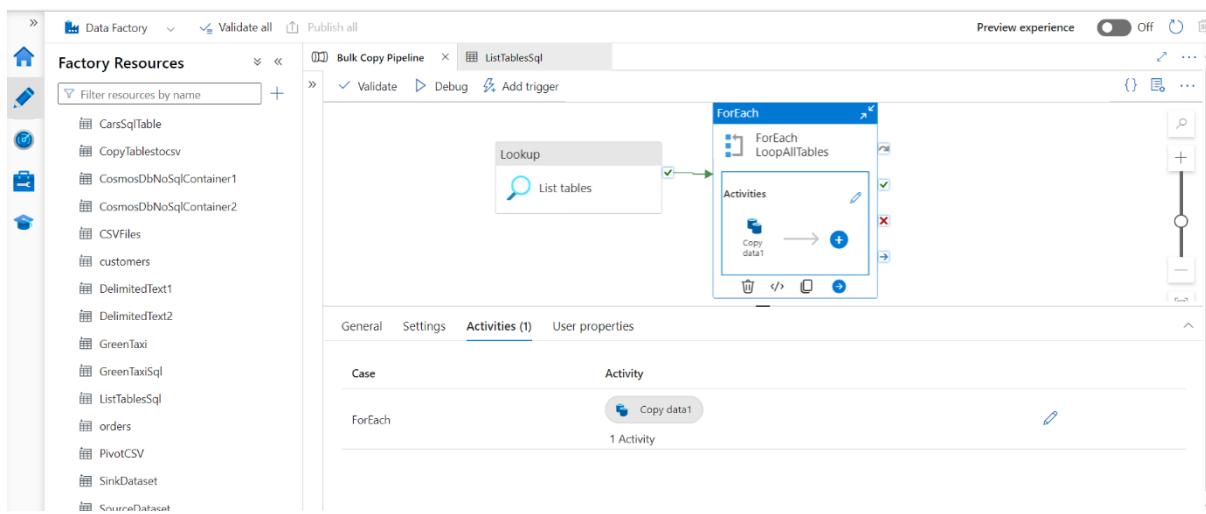
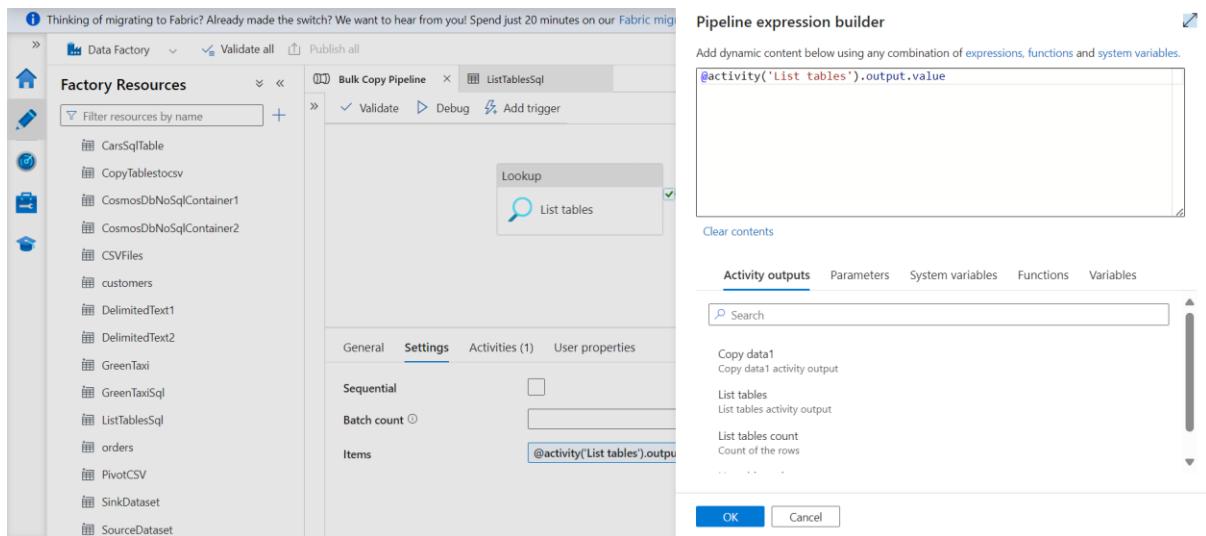
General Settings Activities (1) User properties

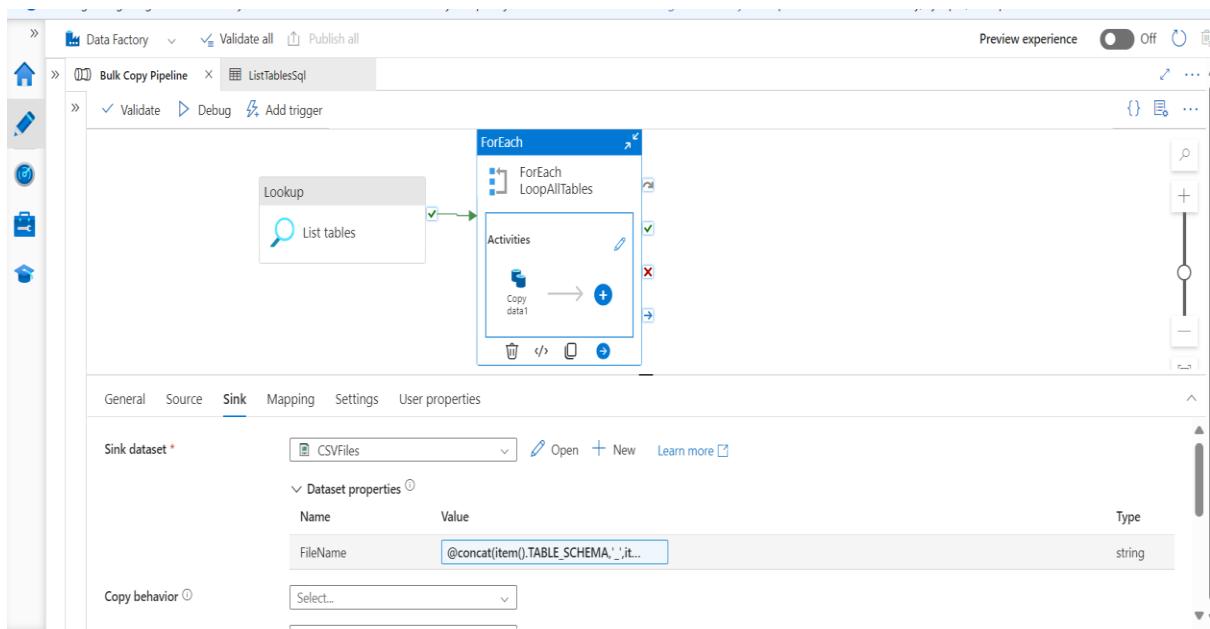
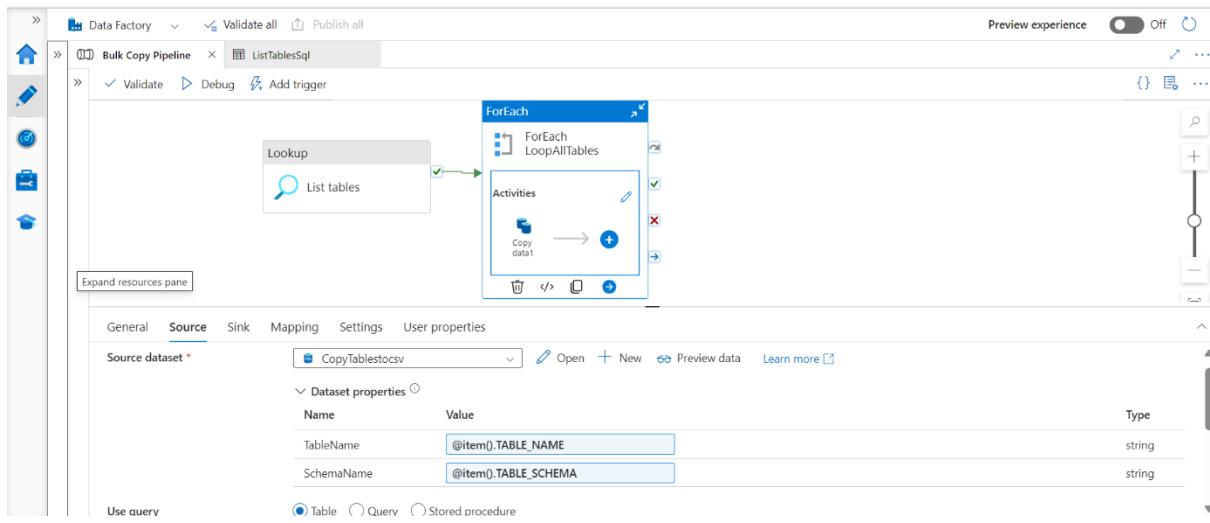
Sequential

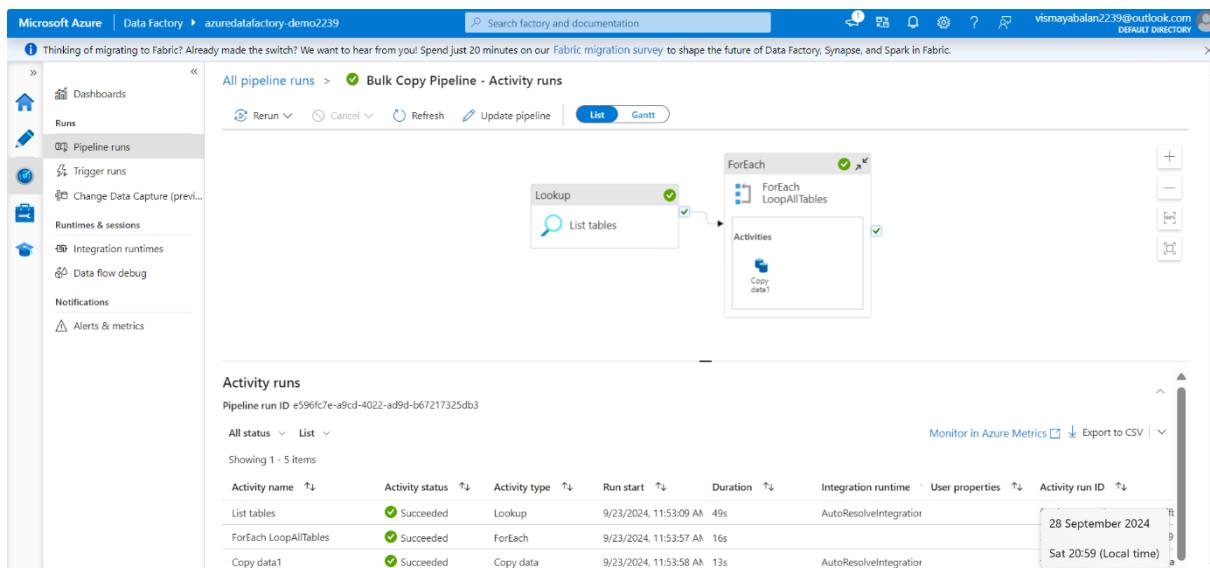
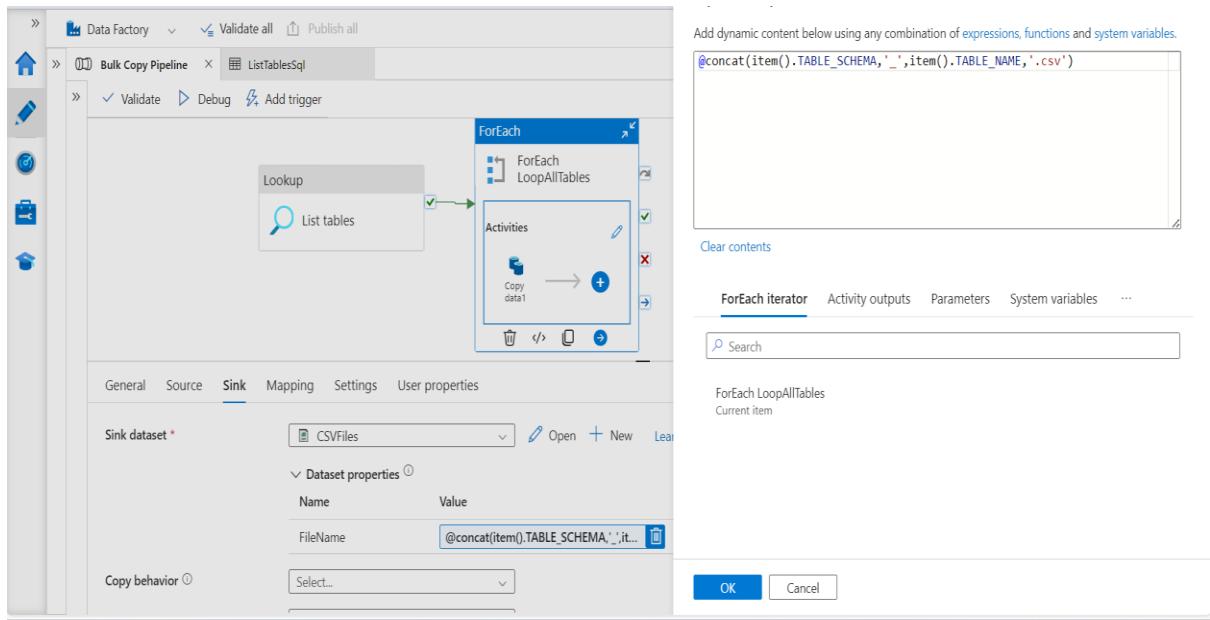
Batch count

Items @activity('List tables').output.value

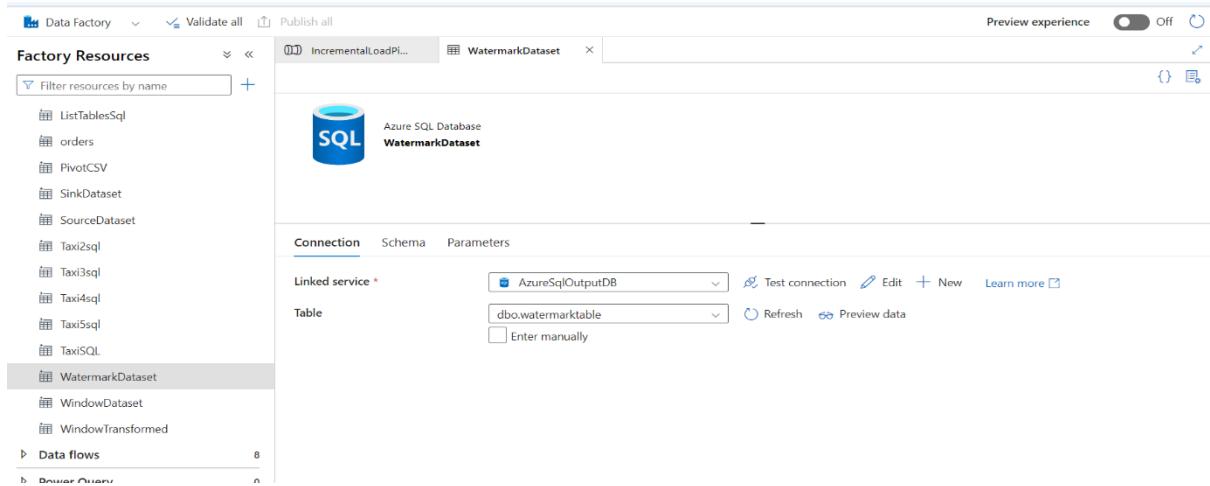
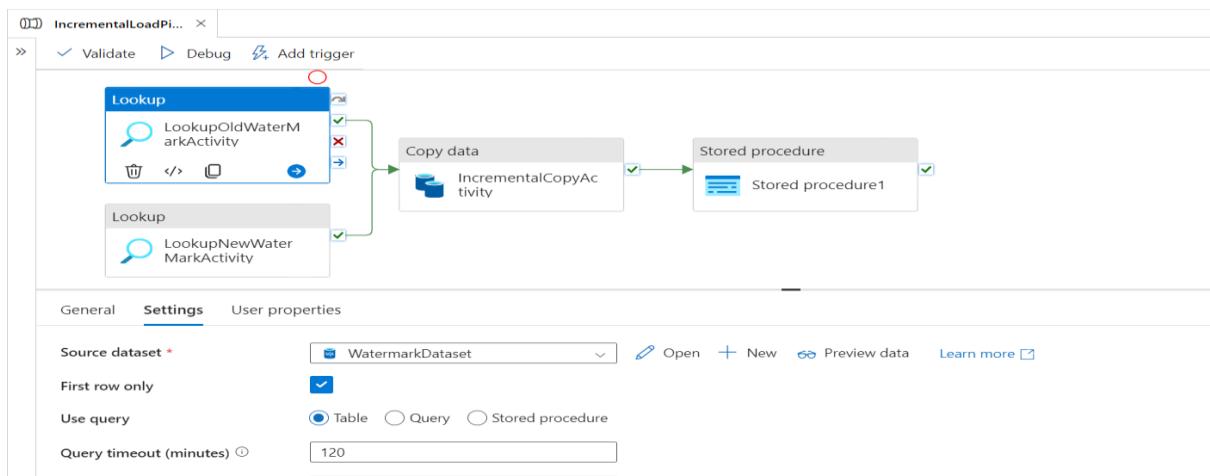
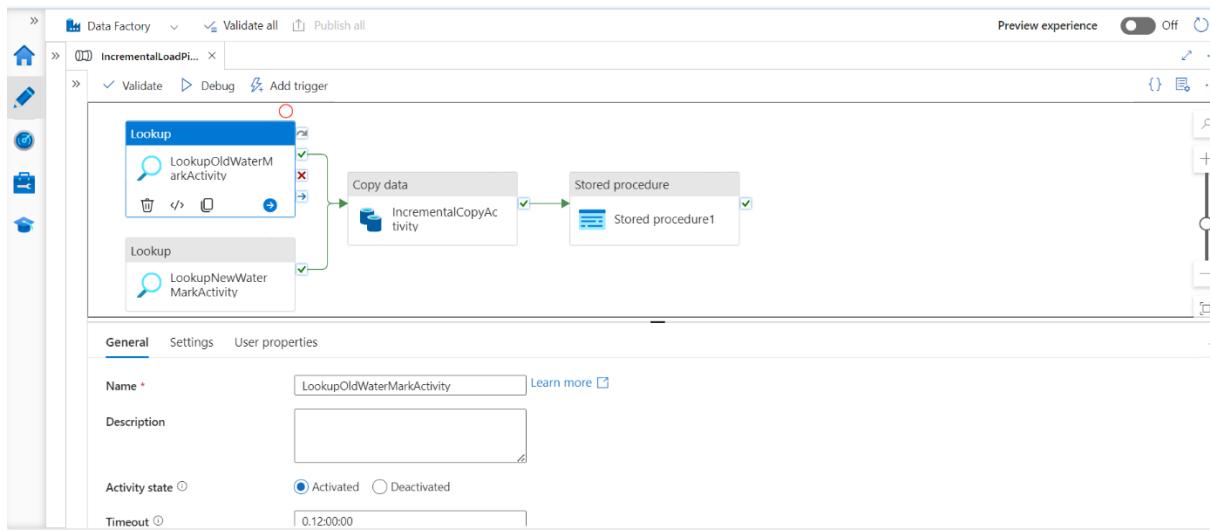








5. Implement incremental load Pipeline in Azure Data Factory for handling datasets, ensuring efficient insert/upsert/updates to the target storage without re-inserting the entire dataset?

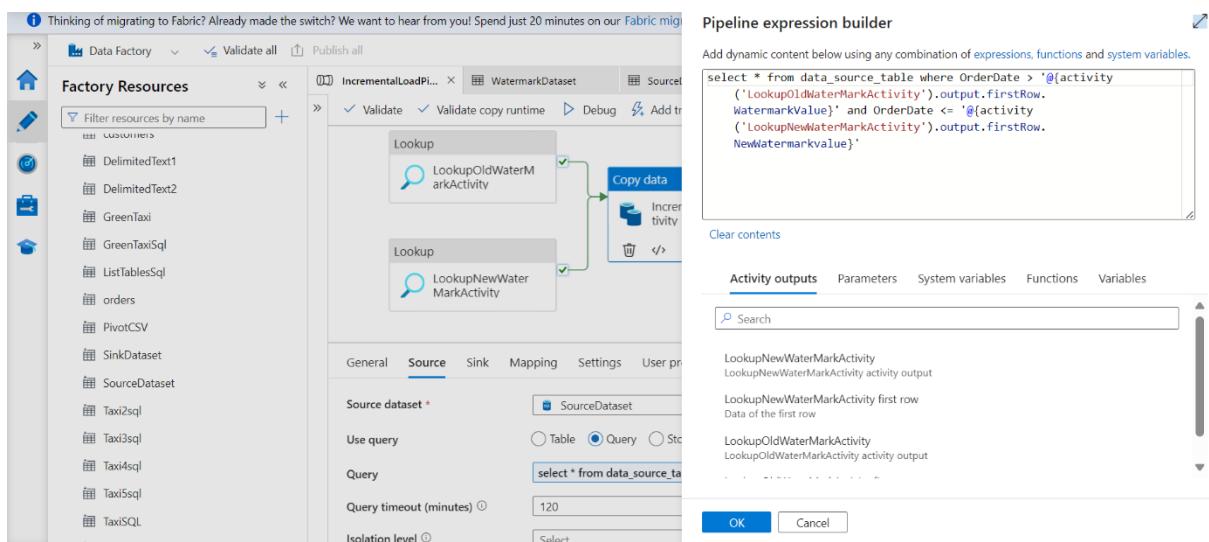
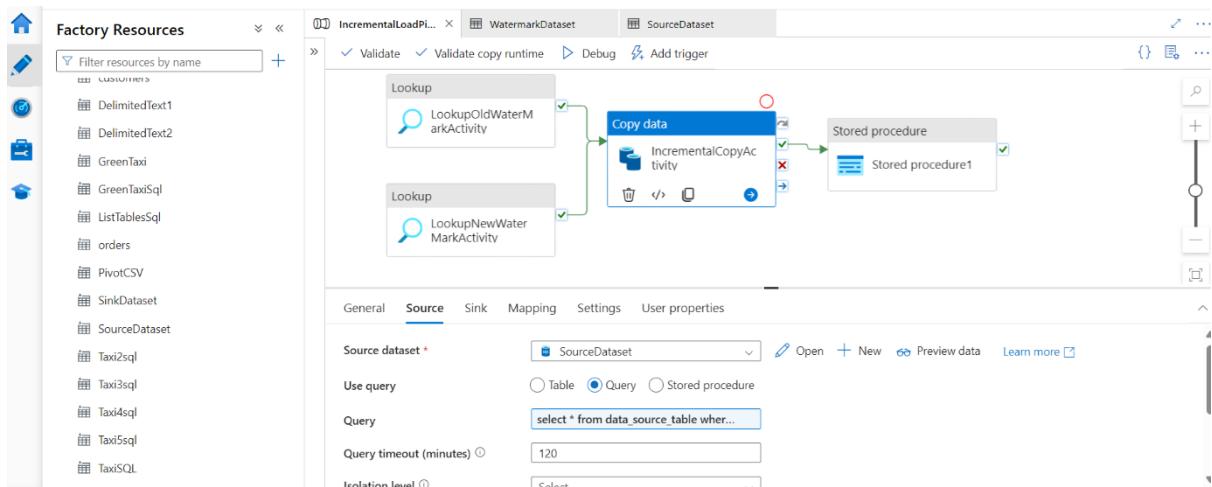


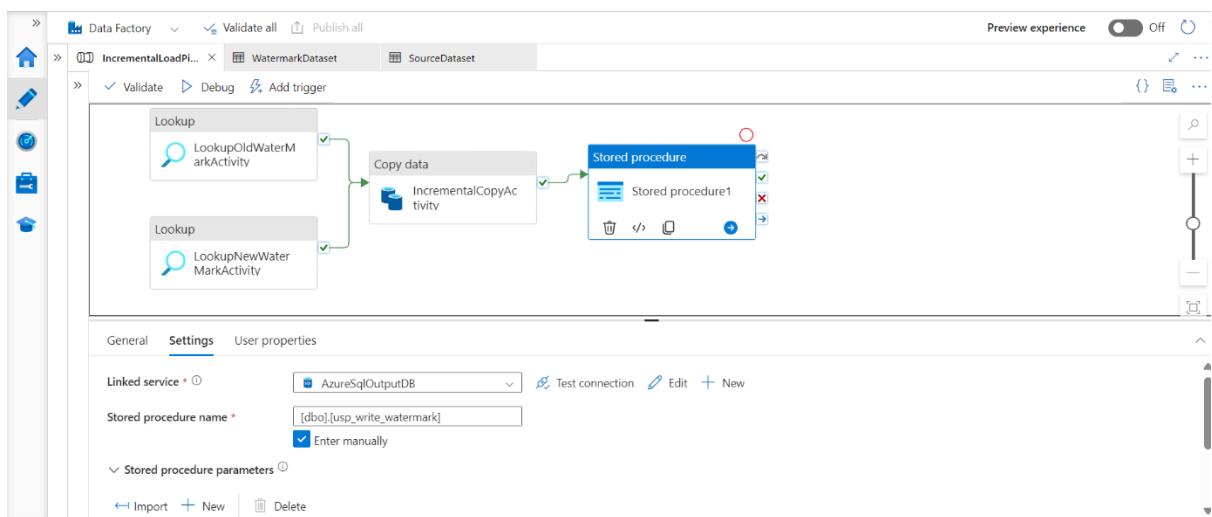
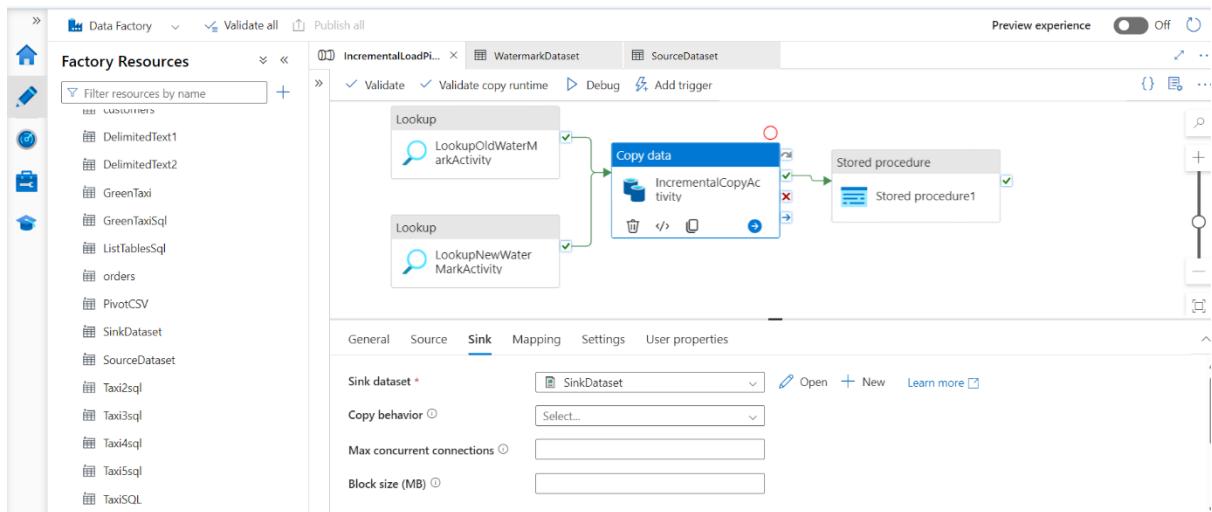
The screenshot shows the Azure Data Factory interface. In the top navigation bar, there are buttons for 'Validate all' and 'Publish all'. On the right, there's a 'Preview experience' toggle set to 'Off'. The main area displays a dataset named 'WatermarkDataset' connected to an 'Azure SQL Database'. The 'Connection' tab is selected, showing options to 'Import schema' or 'Clear'. Below this, the 'Schema' tab lists columns: 'TableName' (varchar) and 'WatermarkValue' (datetime). A toolbar at the bottom right includes icons for copy, edit, and delete.

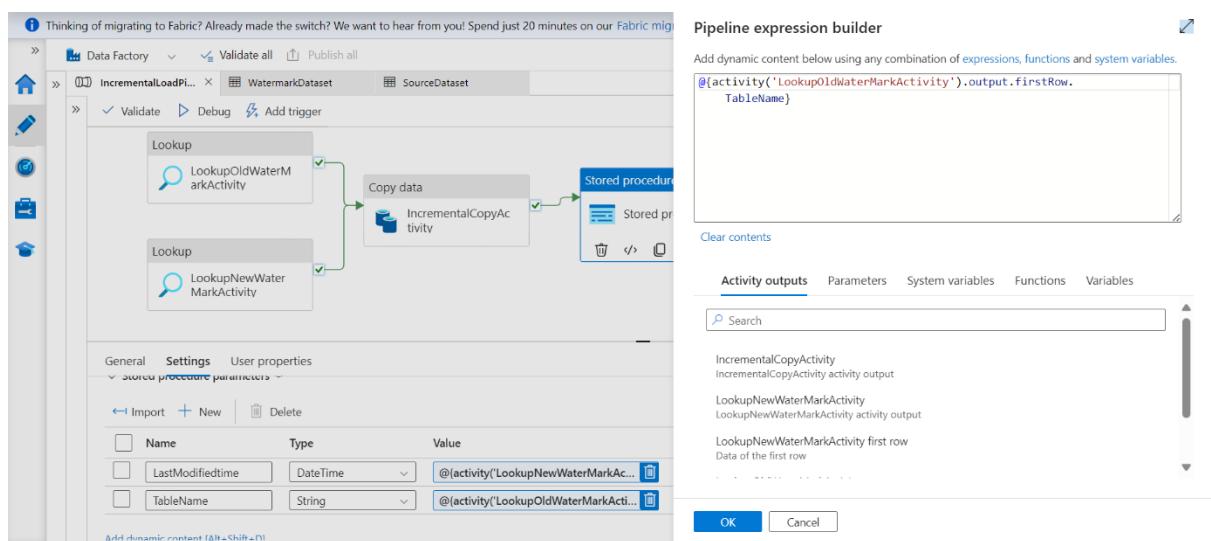
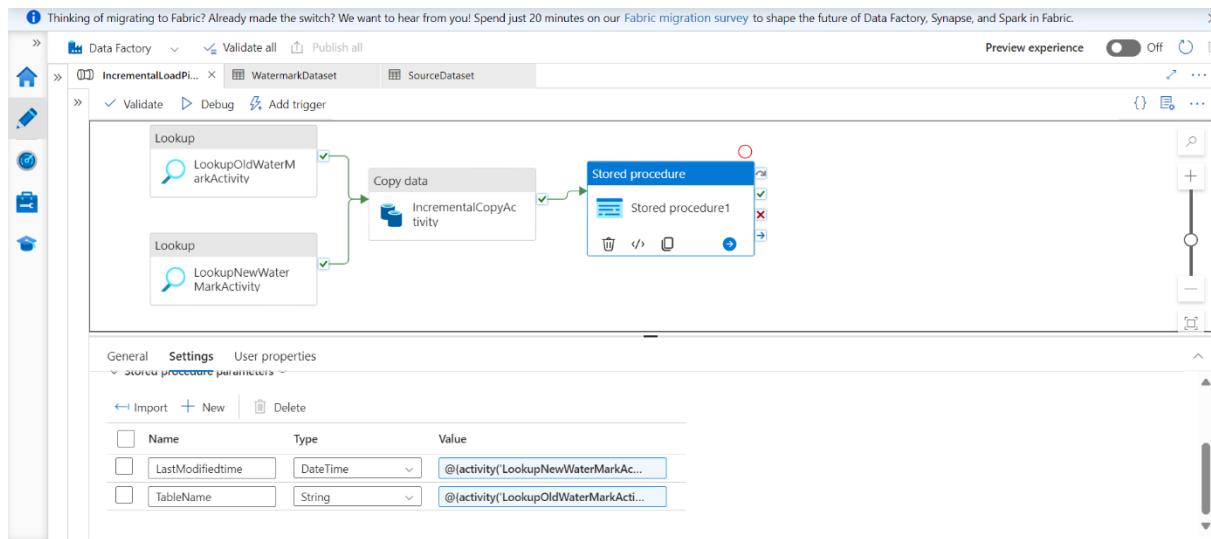
The screenshot shows the Azure Data Factory pipeline editor. At the top, there are buttons for 'Validate', 'Debug', and 'Add trigger'. The pipeline consists of three main components: a 'Lookup' activity named 'LookupOldWaterMarkActivity', another 'Lookup' activity named 'LookupNewWaterMarkActivity', and a 'Copy data' activity named 'IncrementalCopyActivity'. The 'IncrementalCopyActivity' is connected to a 'Stored procedure' named 'Stored procedure1'. Below the pipeline, the 'Settings' tab is selected for the 'Source dataset' 'SourceDataset'. It shows the 'Source dataset' dropdown set to 'SourceDataset', a 'First row only' checkbox checked, and a 'Query' section containing the following SQL query:

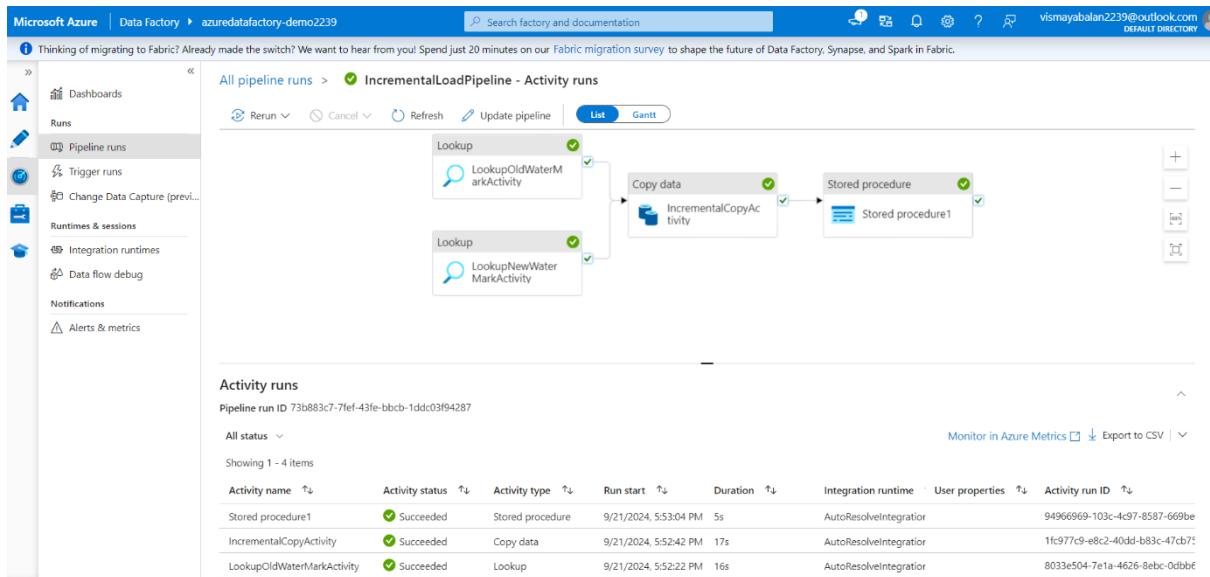
```
select MAX(OrderDate) as NewWatermarkvalue from data_source_table
```

The screenshot shows the 'Factory Resources' page in the Azure Data Factory interface. On the left, a sidebar lists various resources: DelimitedText1, DelimitedText2, GreenTaxi, GreenTaxiSql, ListTablesSql, orders, PivotCSV, SinkDataset, SourceDataset, and several TaxিSQL entries. The 'SourceDataset' item is currently selected. The main pane displays the 'SourceDataset' configuration for an 'Azure SQL Database'. It shows the 'Connection' tab selected, with the 'Linked service' set to 'AzureSqlOutputDB' and the 'Table' set to 'dbo.data\_source\_table'. There are also 'Test connection', 'Edit', 'New', and 'Learn more' buttons.

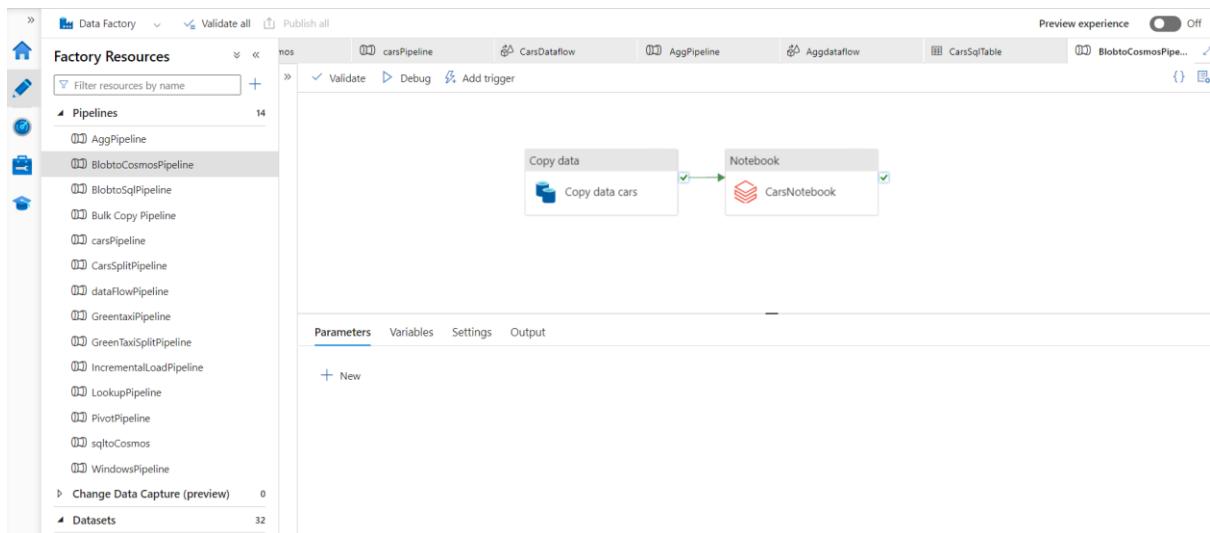








## 6. What are the key steps to connect Azure Databricks to Cosmos DB for real-time analytics and data transformation using spark and Databricks.



Factory Resources

- Pipelines (14)
  - AggPipeline
  - BlobtoCosmosPipeline
  - BlobtoSqlPipeline
  - Bulk Copy Pipeline
  - carsPipeline
  - CarsSplitPipeline
  - dataFlowPipeline
  - GreentaxiPipeline
  - GreenTaxiSplitPipeline
  - IncrementalLoadPipeline
  - LookupPipeline
  - PivotPipeline
  - sqltoCosmos
  - WindowsPipeline
- Change Data Capture (preview) (0)
- Datasets (32)

Copy data

General

Name *	Copy data cars
Description	
Activity state	Activated
Timeout	0:12:00:00
Retry	0

Factory Resources

- Pipelines (14)
  - AggPipeline
  - BlobtoCosmosPipeline
  - BlobtoSqlPipeline
  - Bulk Copy Pipeline
  - carsPipeline
  - CarsSplitPipeline
  - dataFlowPipeline
  - GreentaxiPipeline
  - GreenTaxiSplitPipeline
  - IncrementalLoadPipeline
  - LookupPipeline
  - PivotPipeline
  - sqltoCosmos
  - WindowsPipeline
- Change Data Capture (preview) (0)
- Datasets (32)

Copy data

Source

Source dataset *	Cars
File path type	File path in dataset
Filter by last modified	Start time (UTC) [ ] End time (UTC) [ ]
Recursively	<input checked="" type="checkbox"/>
Enable partitions discovery	<input type="checkbox"/>
Max concurrent connections	[ ]

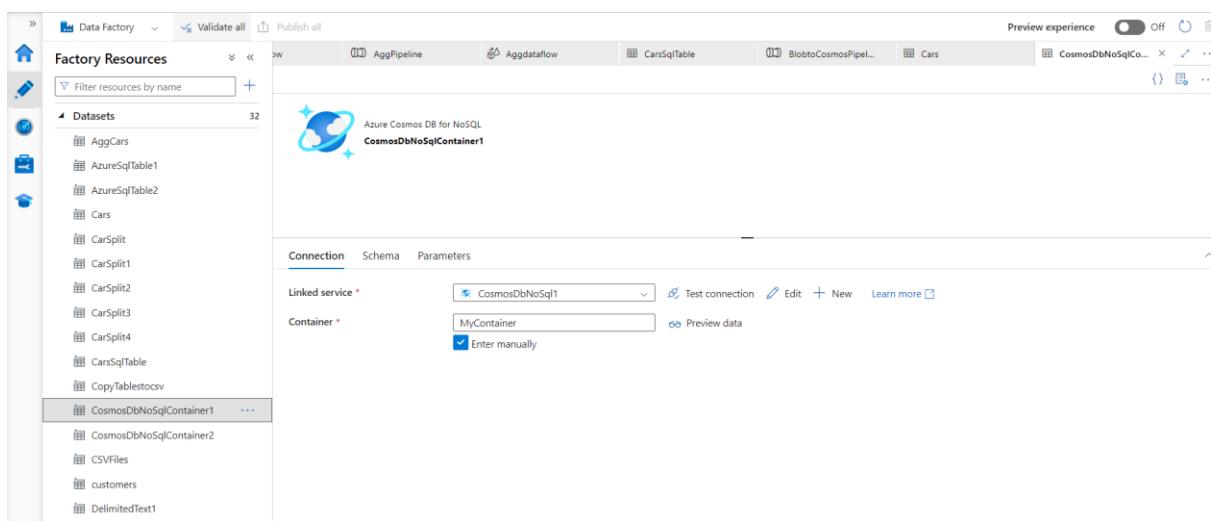
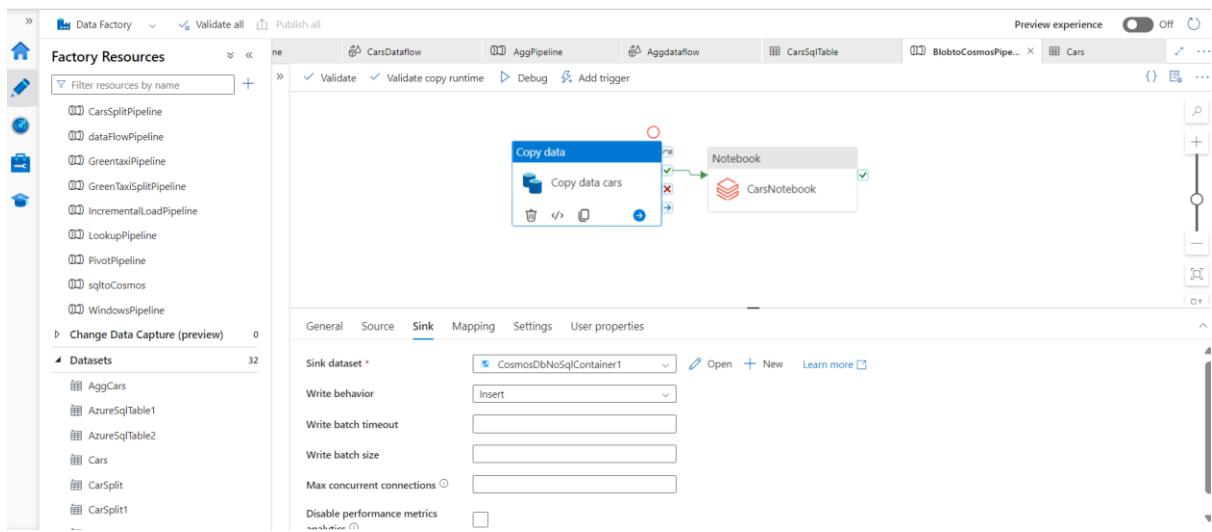
Factory Resources

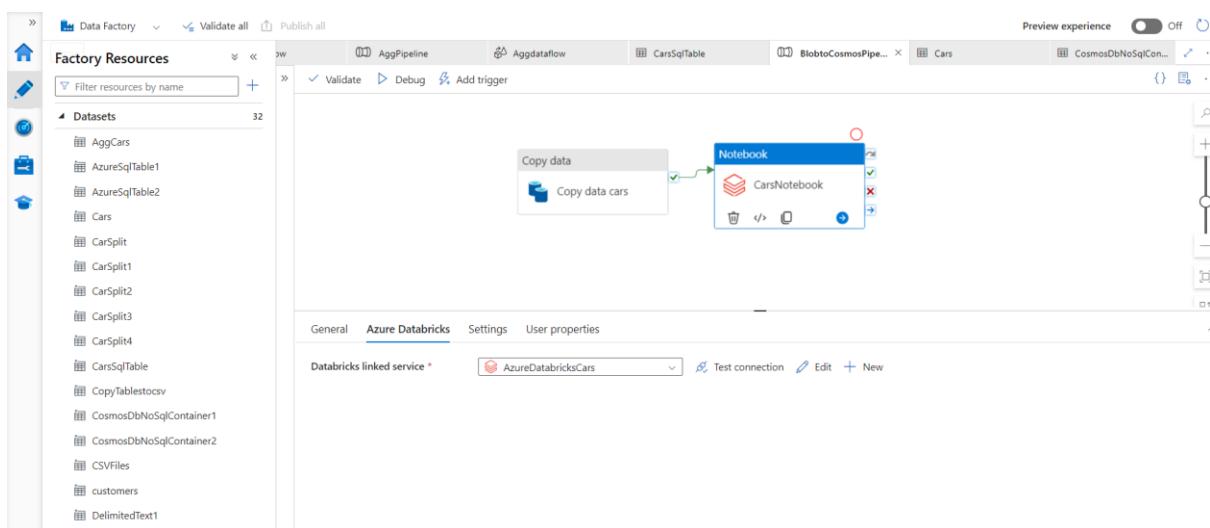
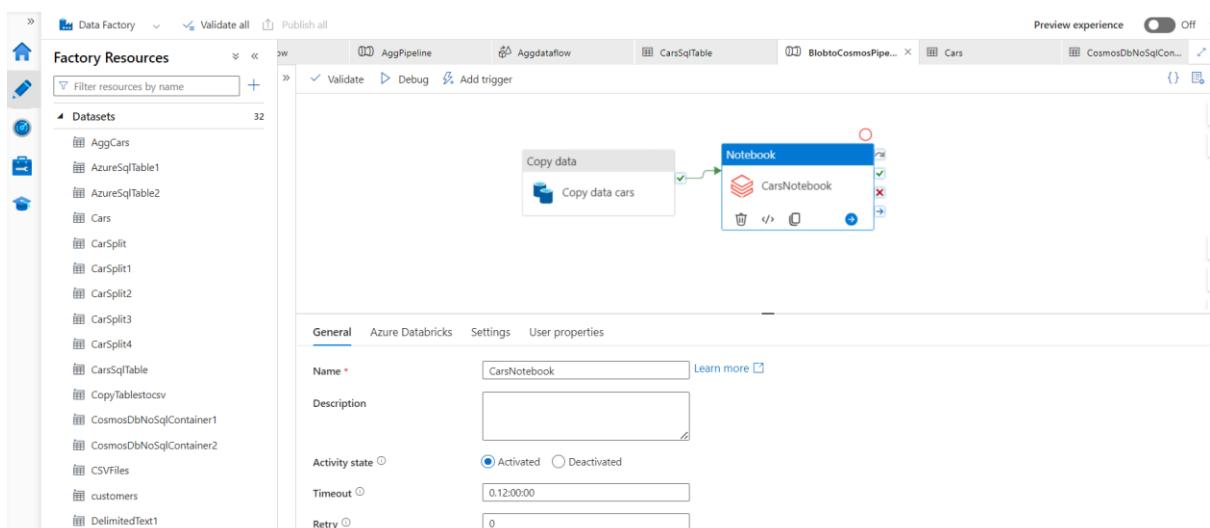
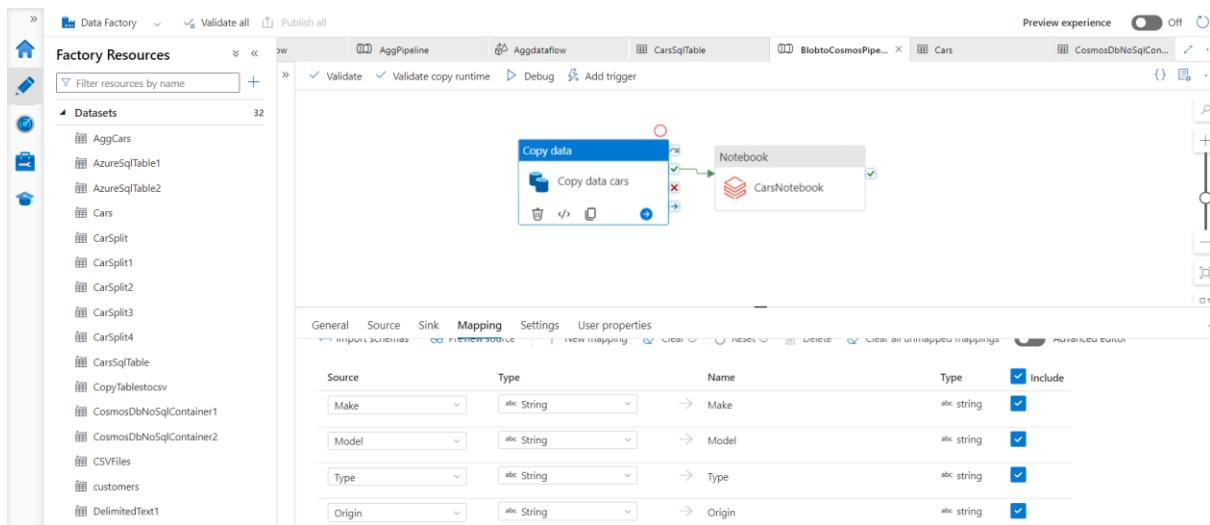
- Pipelines (14)
  - AggPipeline
  - BlobtoCosmosPipeline
  - BlobtoSqlPipeline
  - Bulk Copy Pipeline
  - carsPipeline
  - CarsSplitPipeline
  - dataFlowPipeline
  - GreentaxiPipeline
  - GreenTaxiSplitPipeline
  - IncrementalLoadPipeline
  - LookupPipeline
  - PivotPipeline
  - sqltoCosmos
  - WindowsPipeline
- Change Data Capture (preview) (0)
- Datasets (32)
  - Cars
  - CarSplit
  - CarSplit1

Csv

Connection

Linked service *	InputBlob
File path *	input / Directory / cars.csv
Compression type	Select...
Column delimiter	Comma (,)
Row delimiter	Default (\r\n or \n\r)
Encoding	Default(UTF-8)
Quote character	Double quote (")
Escape character	Backslash (\)





Factory Resources

Datasets

- AggCars
- AzureSqlTable1
- AzureSqlTable2
- Cars
- CarSplit
- CarSplit1
- CarSplit2
- CarSplit3
- CarSplit4
- CarsSqlTable
- CopyTablecsv
- CosmosDbNoSqlContainer1
- CosmosDbNoSqlContainer2
- CSVFiles
- customers
- DelimitedText1

Copy data

Notebook

General Azure Databricks Settings User properties

Notebook path: /Users/vismayabalani2239@outlook.com/... [Browse](#) [Open](#)

Base parameters

Append libraries

```
graph LR; A[Copy data: Copy data cars] --> B[Notebook: CarsNotebook]
```

All pipeline runs > BlobtoCosmosPipeline - Activity runs

Rerun Cancel Refresh Update pipeline List Gantt

Pipeline run ID: 1c3e4e63-f0a5-44a9-a14c-796a8912fa38

Activity runs

Showing 1 - 2 items

Activity name	Activity status	Activity type	Run start	Duration	Integration runtime	User properties	Activity run ID
Copy data cars	Succeeded	Copy data	9/26/2024, 5:51:57 PM	13s	AutoResolveIntegrator		247c774d-11ed-473c-be62-b1ecf
CarsNotebook	Succeeded	Notebook	9/26/2024, 5:52:11 PM	34s	AutoResolveIntegrator		def6714c-ed99-41c9-aa81-5a83bc

Monitor in Azure Metrics Export to CSV

```
graph LR; A[Copy data: Copy data cars] --> B[Notebook: CarsNotebook]
```