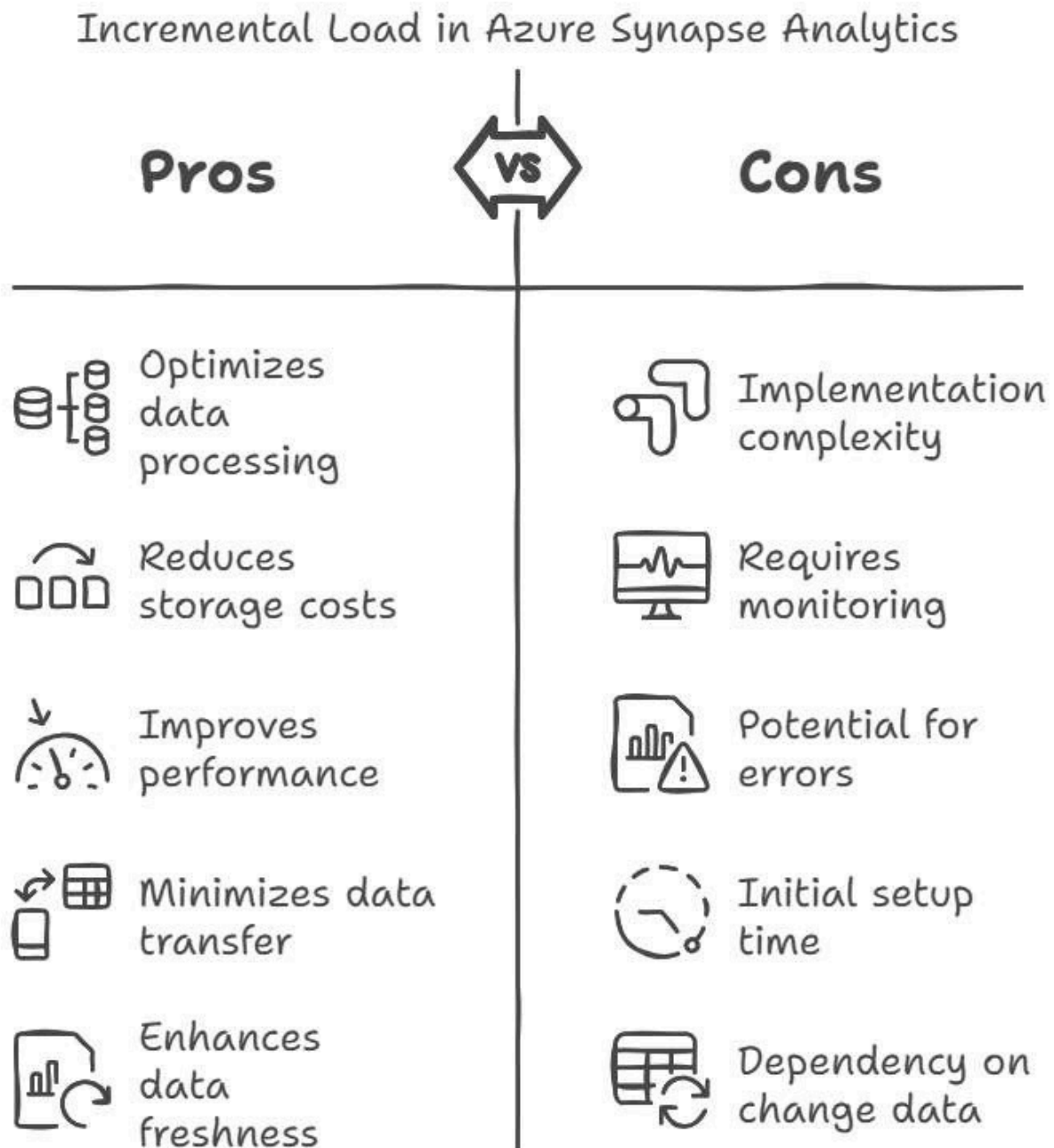


# INCREMENTAL LOADING CONCEPT



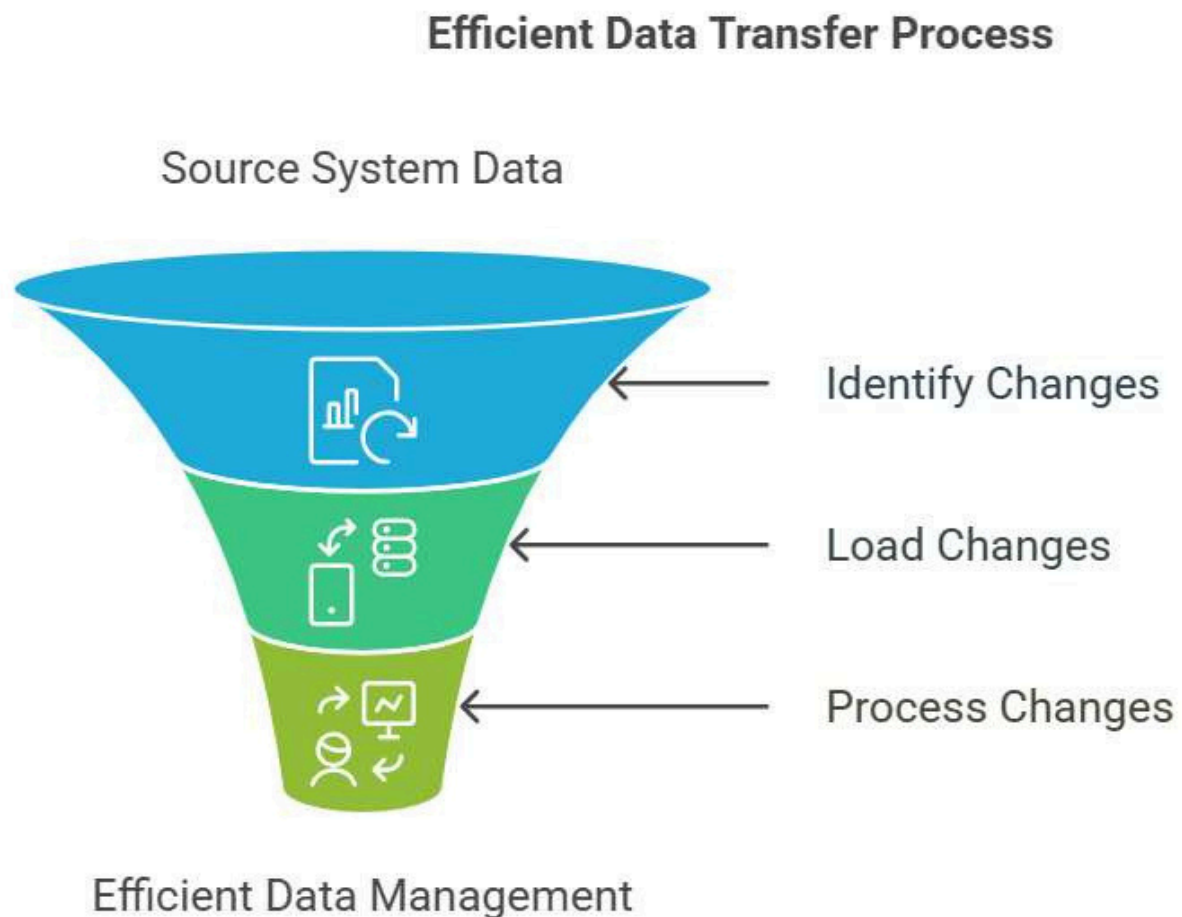
Azure  
Synapse  
Analytics

Incremental load is a crucial concept in data integration and ETL processes, particularly when working with large datasets in **Azure Synapse Analytics**.



# What is the Incremental Load?

Incremental load refers to the process of loading only the new or changed data from a source system into a target system, rather than reloading the entire dataset. This approach is particularly beneficial in scenarios where data volumes are large, and the cost of transferring and processing data can be significant. By focusing on just the changes, organizations can save time, reduce resource consumption, and improve overall efficiency.



# Scenario

Assume you have an online retail store, and you have a lot of data being stored, which needs to be analyzed. There is an example case of understanding the changes in the data over period. Step 01 – Create 5 sample tables I have created multiple tables as follows:

Customer data table (To get the details of the customer) Create Customer id, Name , Phone number, Customer\_datetime

SQL code-

```
--Customer data table (To get thedetailsofthecustomer    )
CREATE TABLE Customer (
    CustomerID INT,
    Name VARCHAR(100) NOT NULL,
    Phone VARCHAR(20),
    Customerupdateddate DATETIME    ---DeltaColumn
);
```

Customer login table (To check the time spent online and on what products) Create Login id, Username, password and login\_datetime

SQL code-

```
--Customer login table (To check the time spentonlineandonwhatproducts)
CREATE TABLE Login_id (
    LoginID INT,
    Username VARCHAR(50) UNIQUE NOT NULL,
    Password VARCHAR(255) NOT NULL,
    Updatedlogindata DATETIME    ---DeltaColumn
);
```

Payment table (To get the list of the transactions/payments completed) Create Transaction\_id, Customer ID, Product ID and Transaction\_datetime

SQL code-

```
--Payment table (To get the list of the transactions/payments completed)
CREATE TABLE Transactions (
    TransactionID INT PRIMARY KEY,    ---DeltaColumn
    CustomerID INT NOT NULL,
    ProductID INT NOT NULL,
    TransactionDate DATE NOT NULL,
);
```

Inventory table (To get the list of items in the inventory) Create Product\_id, Product\_name, Price, Quantity.

SQL code-

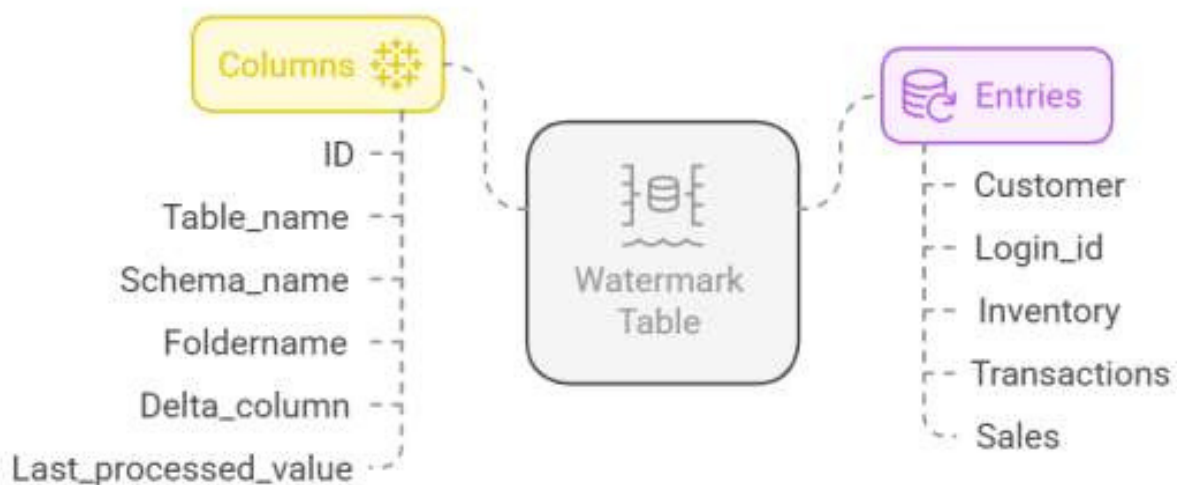
```
--Inventory table (To get the list of items in the inventory)
CREATE TABLE Inventory (
    ProductID INT PRIMARY KEY,          ---DeltaColumn
    ProductName VARCHAR(200) NOT NULL,
    Quantity INT NOT NULL CHECK (Quantity >= 0),
    Price DECIMAL(10, 2) NOT NULL
);
```

Sales table (To check the sales of the items) Create Sales ID, Product ID, Sales\_log, Revenue

SQL code-

```
--Sales table (To check the salesoftheitems)
CREATE TABLE Sales (
    SalesID INT PRIMARY KEY,
    ProductID INT NOT NULL,
    Sales_log DATETIME,                ---DeltaColumn
    Revenue DECIMAL(15, 2) NOT NULL
);
```

## Watermark Table Structure and Initialization



Insert values/data into the tables

```
INSERT INTO Customer (CustomerID, Name, Phone, Customerupdateddate)
VALUES
```

```
(1, 'John Doe', '123-456-7890', '2023-01-01 10:00:00'),
(2, 'Jane Smith', '987-654-3210', '2023-02-01 14:30:00'),
(3, 'Bob Johnson', NULL, '2023-03-01 09:15:00');
```

```
INSERT INTO Login_id (LoginID, Username, Password, Updatedlogindata)
VALUES
```

```
(1, 'johndoe', 'password123', '2023-01-01 10:00:00'),
(2, 'janesmith', 'securepass', '2023-02-01 14:30:00'),
(3, 'bobjohnson', 'secret123', '2023-03-01 09:15:00');
```

```
INSERT INTO Inventory (ProductID, ProductName, Quantity, Price)
VALUES
```

```
(1, 'Laptop', 50, 999.99),
(2, 'Smartphone', 100, 699.99),
(3, 'Headphones', 75, 149.99);
```

```
INSERT INTO Transactions (TransactionID, CustomerID, ProductID, TransactionDate)
VALUES
```

```
(1, 1, 1, '2023-01-05'),
(2, 2, 2, '2023-02-10'),
(3, 3, 3, '2023-03-15');
```

```
INSERT INTO Sales (SalesID, ProductID, Sales_log, Revenue)
```

```
VALUES (1, 1, '2023-01-05 10:00:00', 999.99),
(2, 2, '2023-02-10 14:30:00', 699.99),
(3, 3, '2023-03-15 09:15:00', 299.98);
```

## Step 02 – Create a watermark Table

This step helps to monitor the changes in the data i.e. it may be data entries, data modifications, etc.

SQL Code-

--Create a Watermark table

```
CREATE TABLE Watermark (  
  
    ID INT PRIMARY KEY,                                --Can't accept similar id's or NULL -- Only  
    Unique Value  
    Table_name VARCHAR(100),  
    Schema_name VARCHAR(100),  
    Foldername VARCHAR(50),  
    Delta_column VARCHAR(100),  
    Last_processed_value VARCHAR(255) NOT NULL  
);
```

Inserted data into the watermark table

-- Initialize watermark entries

```
INSERT INTO Watermark  
VALUES  
(1, 'Customer', 'dbo', 'RetailDB/Customer_data', 'Customerupdateddate', '1900-01-01 00:00:00'),  
(2, 'Login_id', 'dbo', 'RetailDB/Login_id_data', 'Updatedlogindata', '1900-01-01 00:00:00'),  
(3, 'Inventory', 'dbo', 'RetailDB/Inventory_data', 'ProductID', '0'),  
(4, 'Transactions', 'dbo', 'RetailDB/Transactions_data', 'TransactionID', '0'),  
(5, 'Sales', 'dbo', 'RetailDB/Sales_data', 'Sales_log', '1900-01-01 00:00:00')
```

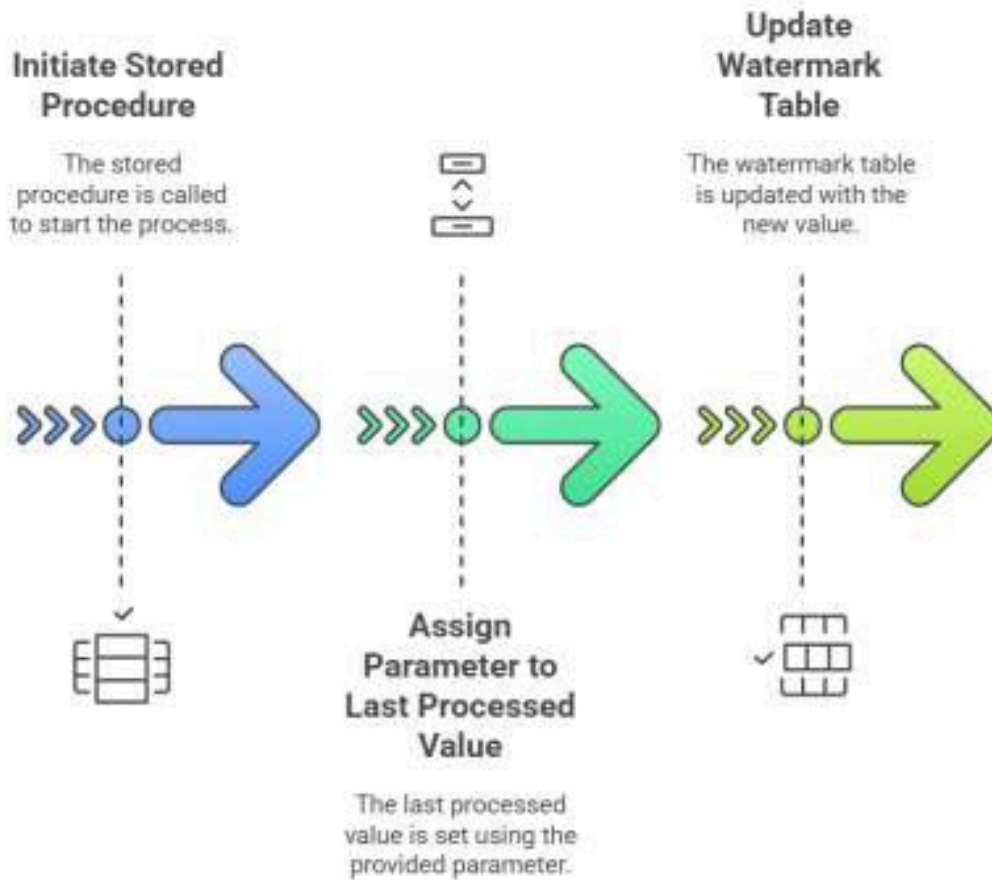
Now, the watermark table is ready to be used.

### Step 03 – Create a stored procedure

This helps to assign a parameter to the changes in the last processed values corresponding to the table names.

```
--CreateStored procedure
CREATE PROC USP_Watermark_RetailDB
@lpv VARCHAR(100),
@TBname VARCHAR(100)
AS
BEGIN
    UPDATE Watermark
    SET Last_processed_value=@lpv WHERE Table_name=@TBname
END
```

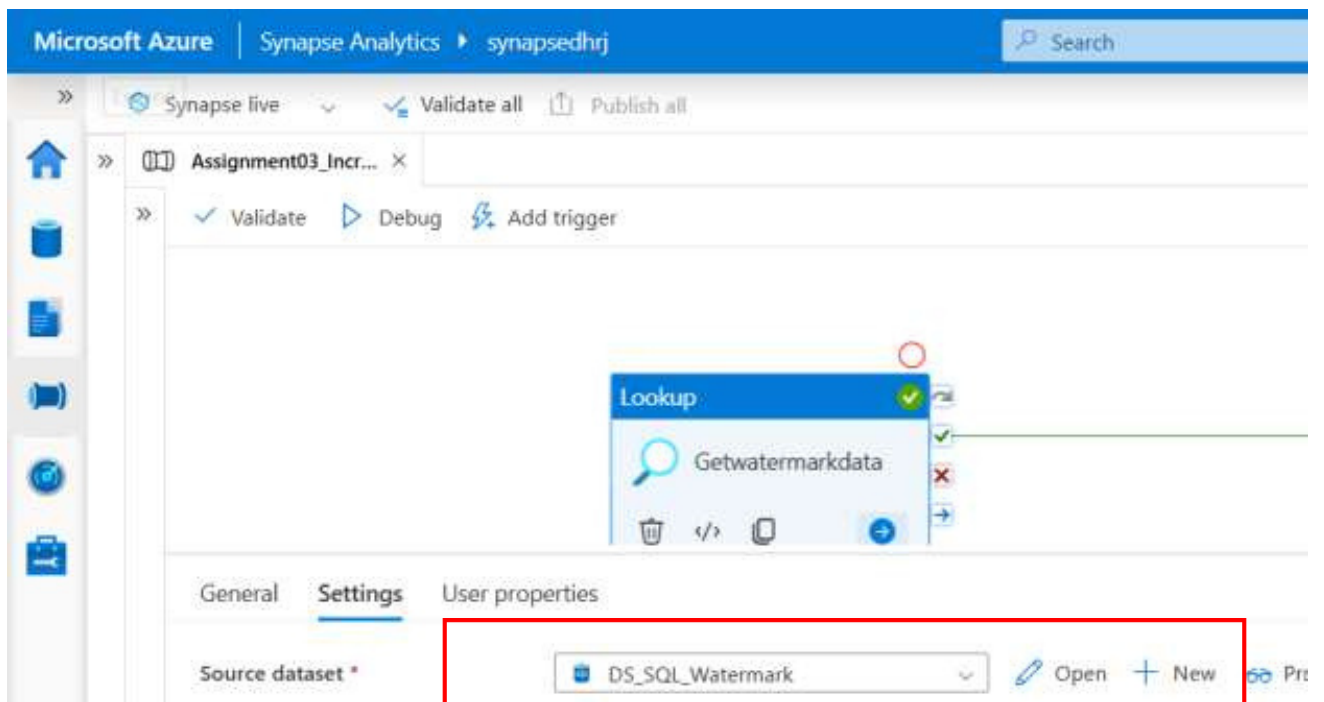
#### Stored Procedure Execution Sequence





## Step 04 – Creation of pipeline in Synapse

- Create a lookup activity and named it as Getwatermarkdata. Click on source dataset and create a new linked service with SQL server and name
- Its as DS\_SQL\_Watermark.



- Open the DS\_SQL\_Watermark dataset.
- Create 2 parameters i.e. Schemaname and Tablename for dynamically assigning the Schemaname and Tablename.

Microsoft Azure | Synapse Analytics ▶ synapsedhrj

Search

>> Synapse live Validate all Publish all

>> Assignment03\_Incr... DS\_SQL\_Watermark

Data

SQL Azure SQL Database DS\_SQL\_Watermark

Connection Schema **Parameters**

+ New Delete


<input type="checkbox"/> Name	Type	Default value	
<input type="checkbox"/> Schemaname	String	Value	
<input type="checkbox"/> Tablename	String	Value	

- Assigning dynamic variables to the Schema and Table name for the table

Microsoft Azure | Synapse Analytics | synapsedhvj

Synapse live | Validate all | Publish all

Assignment03\_Incr... | DS\_SQL\_Watermark

 Azure SQL Database  
DS\_SQL\_Watermark

Connection | Schema | Parameters

Linked service \* | AzureSqlDatabase1 | Test connection | Edit | New | Learn more

Integration runtime \* | AutoResolveIntegrationRuntime | Edit

Table | @dataset().Schemaname | @dataset().Tablename | review data

☒ Enter manually

- Connect the source dataset to DS\_SQL\_Watermark ( Watermark table created in SQL database)
- Enter the values of fields Schemaname and Tablename as dbo and Watermark respectively.

The screenshot displays the Microsoft Azure Synapse Analytics interface. At the top, the header shows 'Microsoft Azure | Synapse Analytics' and a search bar. Below the header, there are tabs for 'Synapse live', 'Validate all', and 'Publish all'. The main workspace shows a 'Lookup' activity named 'Getwatermarkdata' with a green checkmark indicating it is successful. The 'Settings' tab is selected, showing the 'Source dataset' as 'DS\_SQL\_Watermark'. Below this, the 'Dataset properties' section is expanded, showing a table with 'Name' and 'Value' columns. The 'Schemaname' is set to 'dbo' and the 'Tablename' is set to 'Watermark'. Other settings include 'First row only' (unchecked), 'Use query' (radio button selected for 'Table'), 'Query timeout (minutes)' set to 120, and 'Isolation level' set to 'Select...'. A red box highlights the 'Source dataset' and 'Dataset properties' sections.

Microsoft Azure | Synapse Analytics | synapsedhrj

Synapse live | Validate all | Publish all

Assignment03\_Incr... x

Validate | Debug | Add trigger

Lookup

Getwatermarkdata

General | **Settings** | User properties

Source dataset \*

DS\_SQL\_Watermark | Open | + New | Pr

Dataset properties ⓘ

Name	Value
Schemaname	dbo
Tablename	Watermark

First row only ☐

Use query ☒ Table ☐ Query ☐ Stored procedure

Query timeout (minutes) ⓘ 120

Isolation level ⓘ Select... v

- Create a Foreach activity
- Connect the lookup activity i.e. Getwatermarkdata with the Foreach activity

The screenshot displays the Microsoft Azure Synapse Analytics interface. The top navigation bar shows "Microsoft Azure" and "Synapse Analytics" with a search bar. The main workspace shows a workflow diagram with two activities: a "Lookup" activity named "Getwatermarkdata" and a "Foreach" activity named "Foreach1". The "Lookup" activity is connected to the "Foreach" activity. The "Foreach" activity has a "ForEach1" sub-activity and an "Activities" section. Below the workflow diagram, the "General" tab is selected, showing the "Name" field with the value "Getmaxvalue", a "Description" field, and the "Activity state" set to "Activated".

Microsoft Azure | Synapse Analytics | synapsedhrj

Synapse live | Validate all | Publish all

Assignment03\_Incr...

Validate | Debug | Add trigger

Lookup

Getwatermarkdata

ForEach

ForEach1

Activities

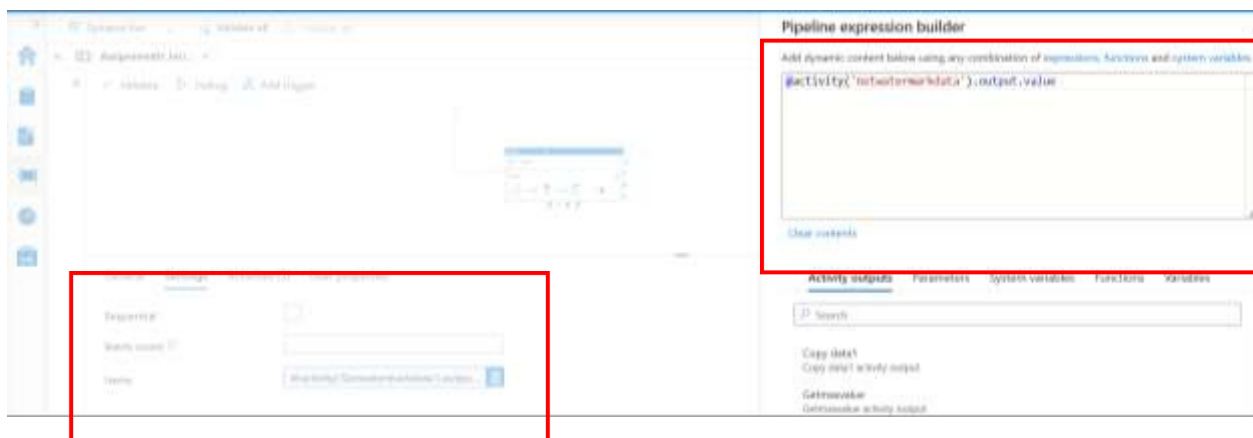
General | Settings | User properties

Name \* Getmaxvalue [Learn more](#)

Description

Activity state ☒ Activated ☐ Deactivated

- Provide an expression in the items field in the Settings tab to connect with the lookup activity.
- Go to the pipeline expression builder-@activity('Lookuptablename').output.value i.e. in current scenario @activity('Getwatermarkdata').output.value



- Created a Lookup activity inside the Foreach activity for deriving the maximum value from the watermark table.
- Name it as lookup activity as Getmaxvalue

The screenshot displays the Microsoft Azure Synapse Analytics interface. The top navigation bar shows 'Microsoft Azure' and 'Synapse Analytics' with a search bar. Below the navigation bar, there are tabs for 'Synapse live', 'Validate all', and 'Publish all'. The main workspace shows a pipeline named 'Assignment03\_Incr...' with a 'Foreach1' loop. Inside the loop, a 'Lookup' activity named 'Getmaxvalue' is highlighted. The 'Lookup' activity is shown in a context menu with options to delete, edit, copy, and paste. Below the workspace, the 'General' tab is selected, showing the activity's configuration. The 'Name' field is set to 'Getmaxvalue', the 'Description' field is empty, the 'Activity state' is 'Activated', and the 'Timeout' is set to '0.12:00:00'.

Microsoft Azure | Synapse Analytics | synapsedhrj

Synapse live | Validate all | Publish all

Assignment03\_Incr... x

Validate | Debug | Add trigger

Assignment03\_Incremental data loading > ForEach1

Lookup

Getmaxvalue

General | Settings | User properties

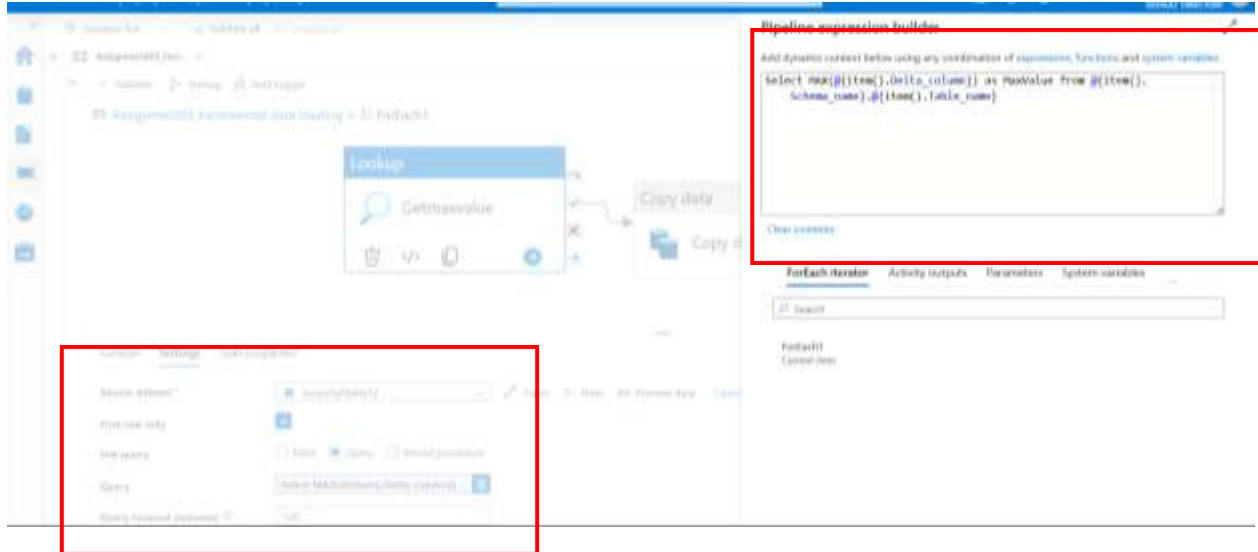
Name \* | Getmaxvalue | [Learn more](#)

Description

Activity state ☒ Activated ☐ Deactivated

Timeout

- Providing a query dynamically to get the maximum value from the watermark table  
General expression ( For getting the maximum value from the watermark table) `SELECT MAX('Columnname') as MaxValue FROM dbo.Watermark`
- Dynamic expression (Converting dynamically the above expression) `SELECT MAX(@item().Delta_Column) as MaxValue FROM @item().Schema_name).@item().Table_name)`





- Adding a Copy data activity
- Connect with the lookup activity i.e. Getmaxvalue activity.
- Select Source and write a dynamic query to copy only the modified values from the watermark table

The screenshot displays the Microsoft Azure Synapse Analytics interface. At the top, the header shows 'Microsoft Azure | Synapse Analytics | synapsiedhrj'. Below the header, the pipeline is named 'Assignment03\_Incr...' and is in a 'Synapse live' state. The pipeline contains two activities: 'Lookup' (with a 'Getmaxvalue' operation) and 'Copy data' (labeled 'Copy data1'). A red box highlights the 'Copy data' activity. Below the pipeline canvas, the 'Source' tab is selected, showing the configuration for the 'Copy data1' activity. The 'Source dataset' is set to 'AzureSqlTable12'. The 'Use query' option is selected, and the 'Query' field contains the dynamic query: 'SELECT \* FROM @(item().Table\_name...)'. Other settings include 'Query timeout (minutes)' set to 120, 'Isolation level' set to 'Select...', and 'Partition option' set to 'None'.

Microsoft Azure | Synapse Analytics | synapsiedhrj

Synapse live | Validate all | Publish all

Assignment03\_Incr... | Validate | Validate copy runtime | Debug | Add trigger

Assignment03\_Incremental data loading > ForEach1

Lookup  
Getmaxvalue

Copy data  
Copy data1

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

Source dataset \* | AzureSqlTable12 | Open | New | Preview data | Learn more

Use query | Table | **Query** | Stored procedure

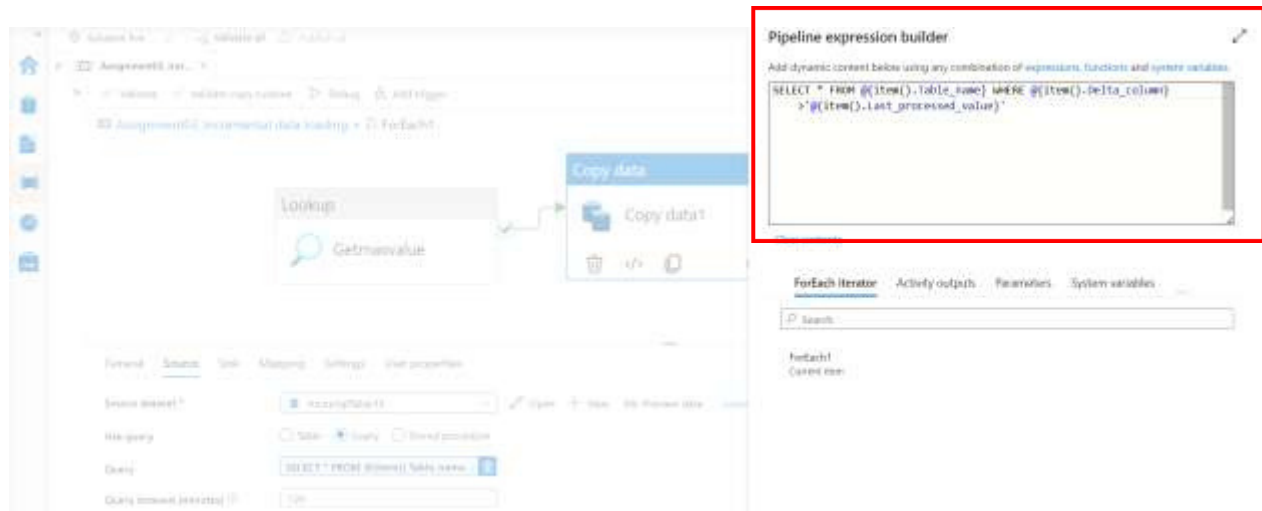
Query | SELECT \* FROM @(item().Table\_name...)

Query timeout (minutes) | 120

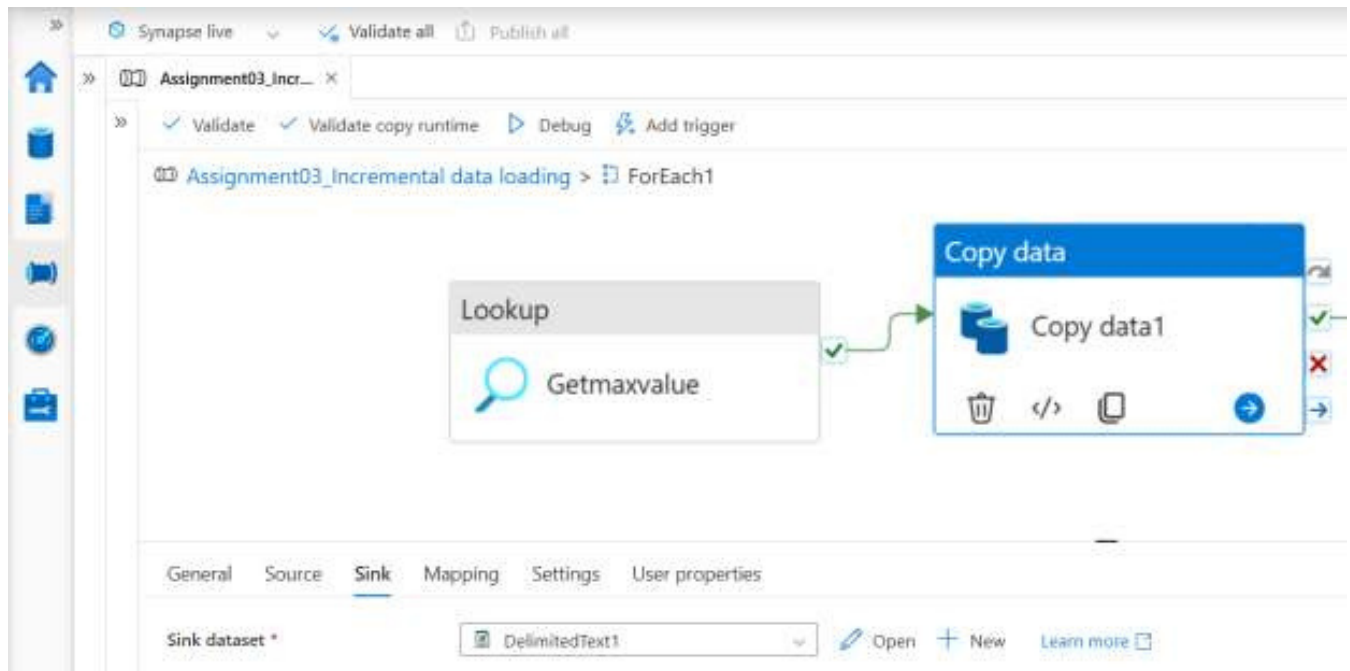
Isolation level | Select...

Partition option | **None** | Physical partitions of table | Dynamic range

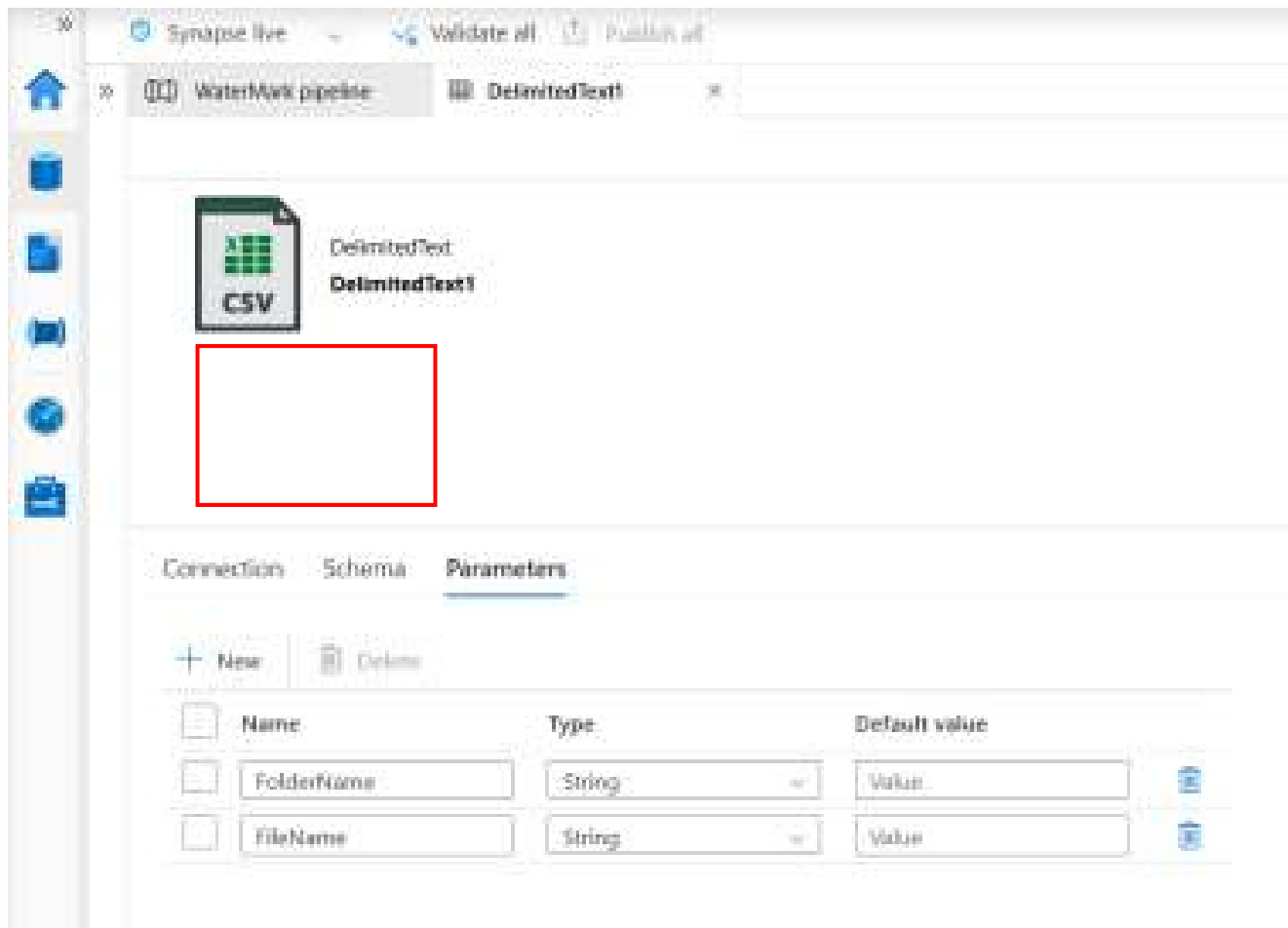
- General expression ( To check for modified values or changes in the data entries)  
SELECT \* FROM TABLE WHERE Delta\_Column > 'Last processed value column'
- Dynamic expression ( Converting dynamically the above expression) SELECT \*  
FROM @ (item().Table\_name) WHERE  
@ (item().Delta\_column) > '@ (item().Last\_processed\_value)



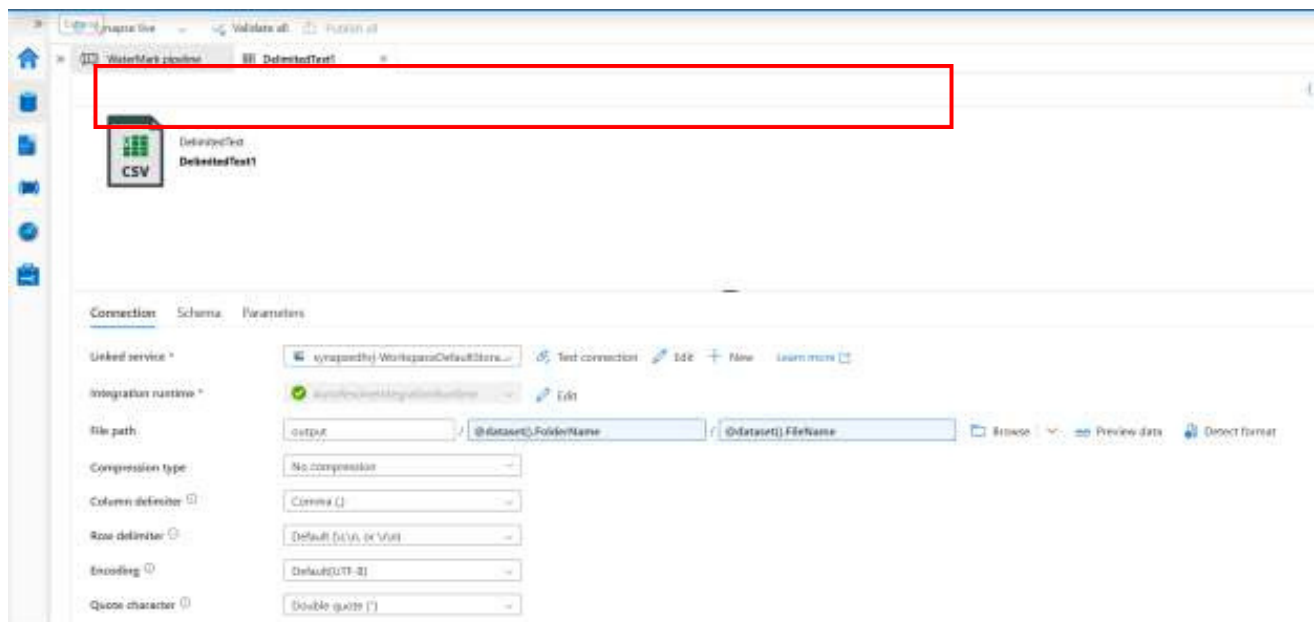
- Select Sink and create a dataset name DelimitedText1 to store the output into a Azurelakegen2 or Blob storage .



- Create 2 parameters named FolderName and FileName to store the folder name and file names respectively



- Assign the parameters dynamically to the folder name and file name in the File path.



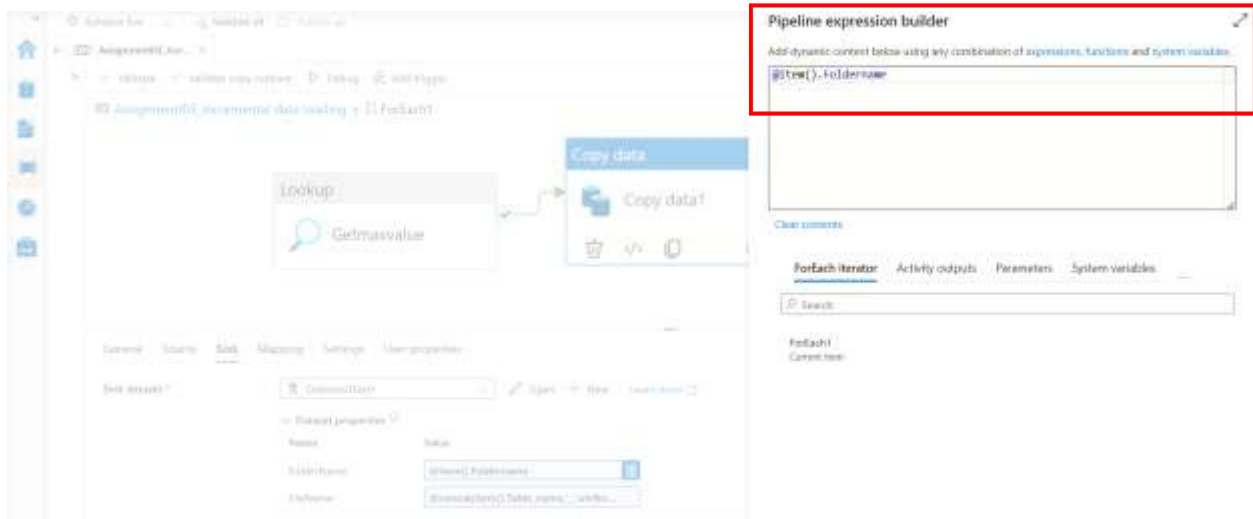
- Provide expressions as dynamic queries in the Foldername field and Filename fields.

The screenshot displays the Synapse IDE interface. At the top, there are tabs for 'Assignment03\_Incr...' and 'Assignment03\_Incremental data loading > ForEach1'. Below the tabs, there are buttons for 'Validate', 'Validate copy runtime', 'Debug', and 'Add trigger'. The main workspace shows a data pipeline with a 'Lookup' activity (containing 'Getmaxvalue') connected to a 'Copy data' activity (containing 'Copy data1').

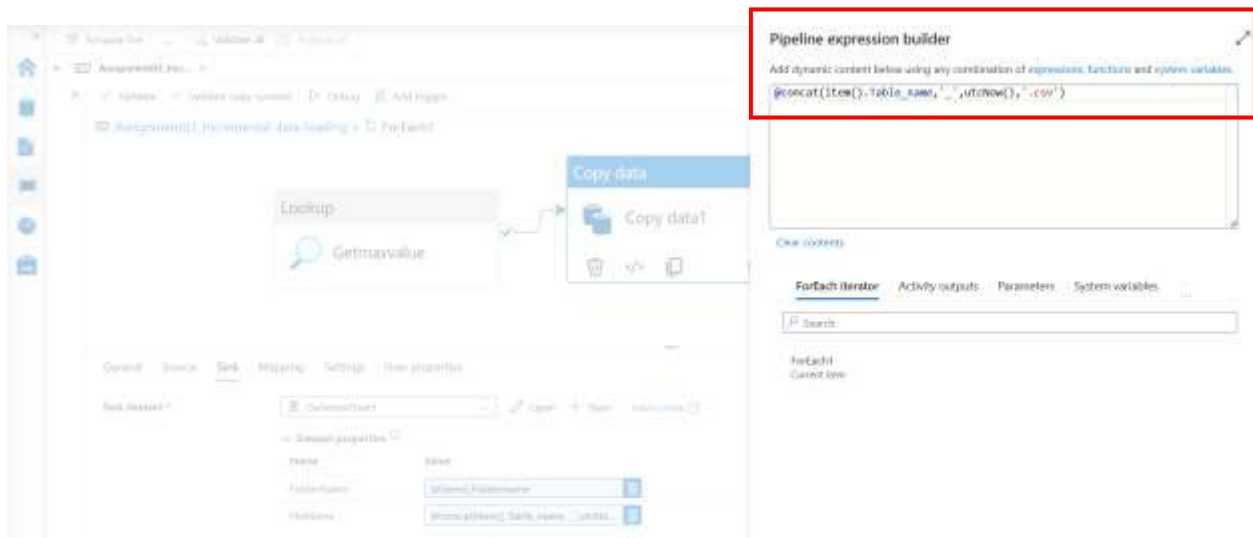
Below the workspace, the 'Sink' tab is selected in the 'Dataset properties' section. The 'Sink dataset' is set to 'DelimitedText1'. The 'Dataset properties' table is shown below:

Name	Value
FolderName	@item().Foldername
FileName	@concat(item().Table_name,'_utcNo...

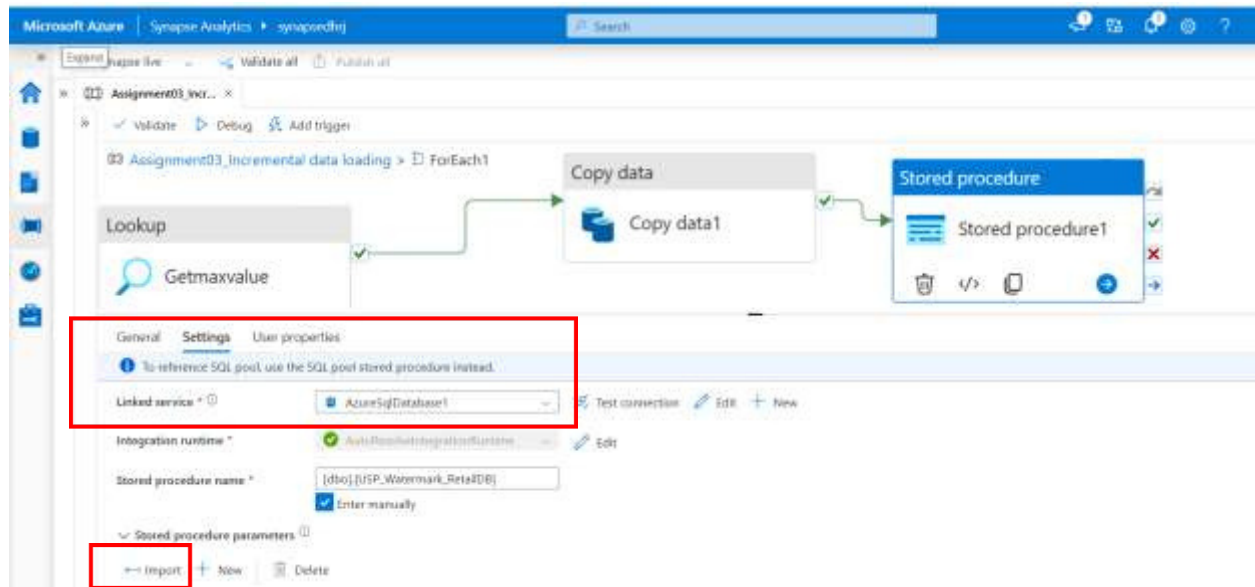
- Expression for the Foldername field would be `@(item().Foldername)`
- This expression will output the values with corresponding foldernames.



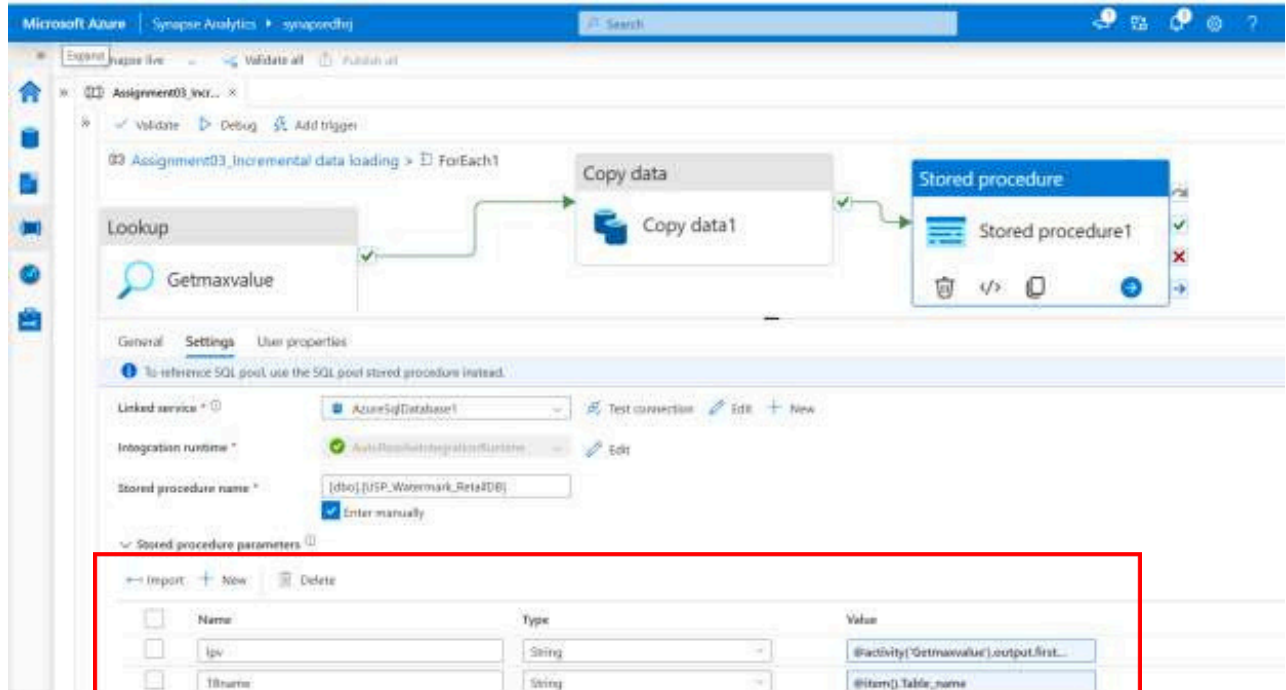
- Expression for the Foldername field would be `@(item().Foldername)` and `@concat(item().Tablename,'_',utcNow(),'.csv')`
- This expression will output the values as a .csv file format with current timestamp.



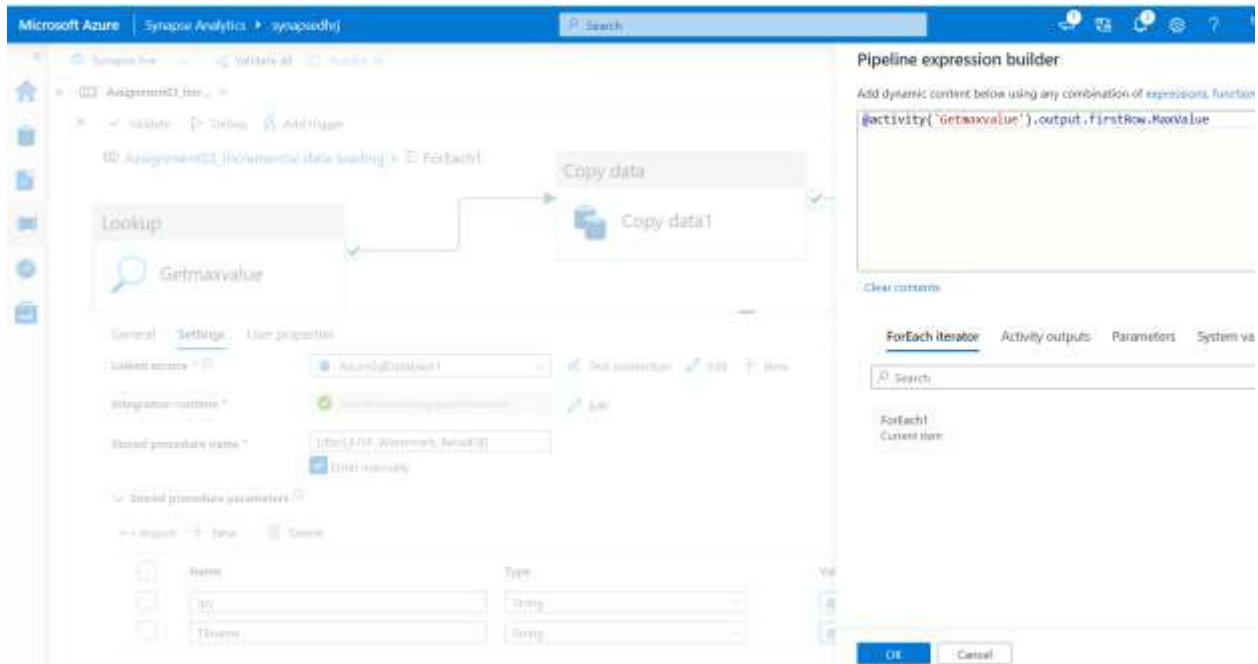
- Next Step - Adding a Stored procedure activity.
- Linking the source to the SQL database.



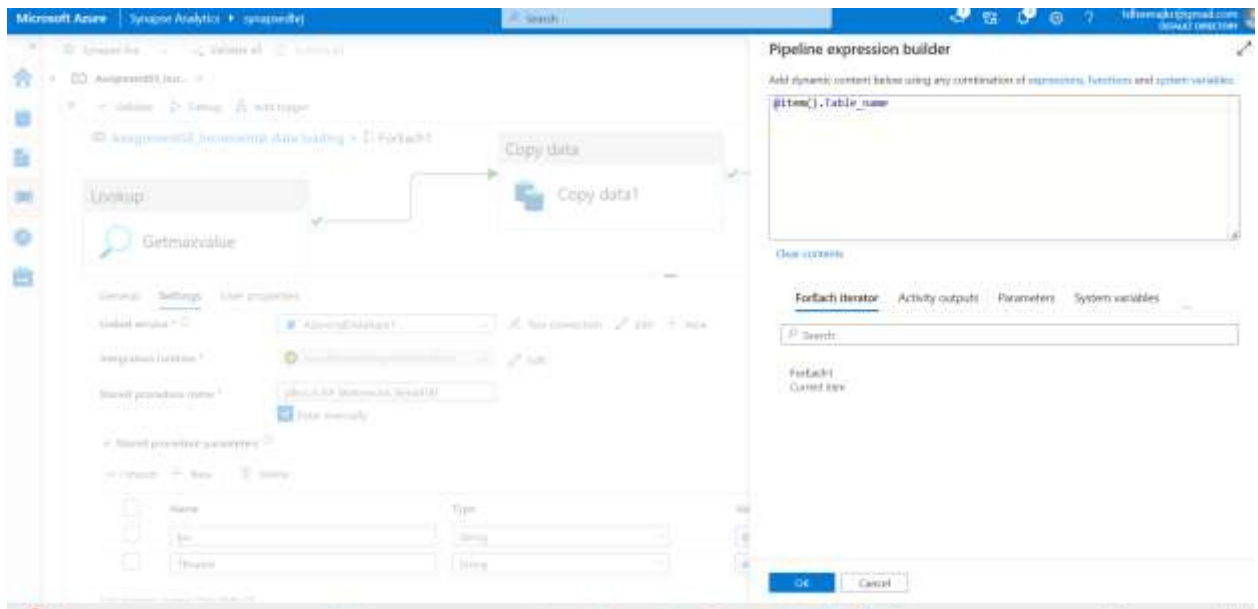
- Next click on Import parameters to import all the stored procedures created .



- Providing the dynamic expression for the lpv  
@activity('Getmaxvalue').output.firstRow.Maxvalue
- This will only input the modified values from the dataset.



- Providing the dynamic expression for the TBname .@item().Table\_name.
- This will only input the modified values from the respective Table names.





- Publishing and checking the pipeline if its working

Microsoft Azure | Synapse Analytics | synapseedge

Search

Assignment01, Inc...

Validate Debug Add trigger

Parameters Variables Settings **Output**

Pipeline run ID: 3b5e2666-44ac-4543-8530-3f8b798e264e Pipeline status: ✔ Succeeded [View debug run consumption](#)

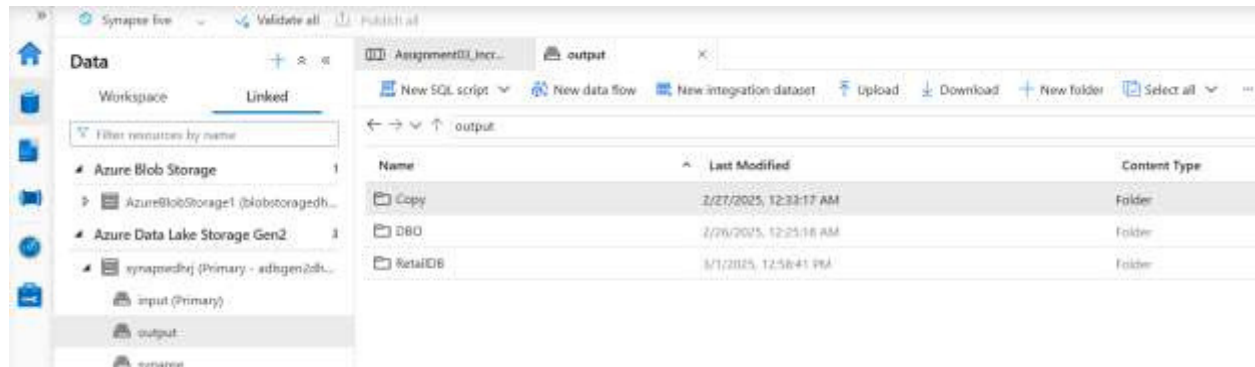
All status: ✔ List [Monitor in Azure Monitor](#) [Export to CSV](#)

Showing 1 - 17 of 17 items

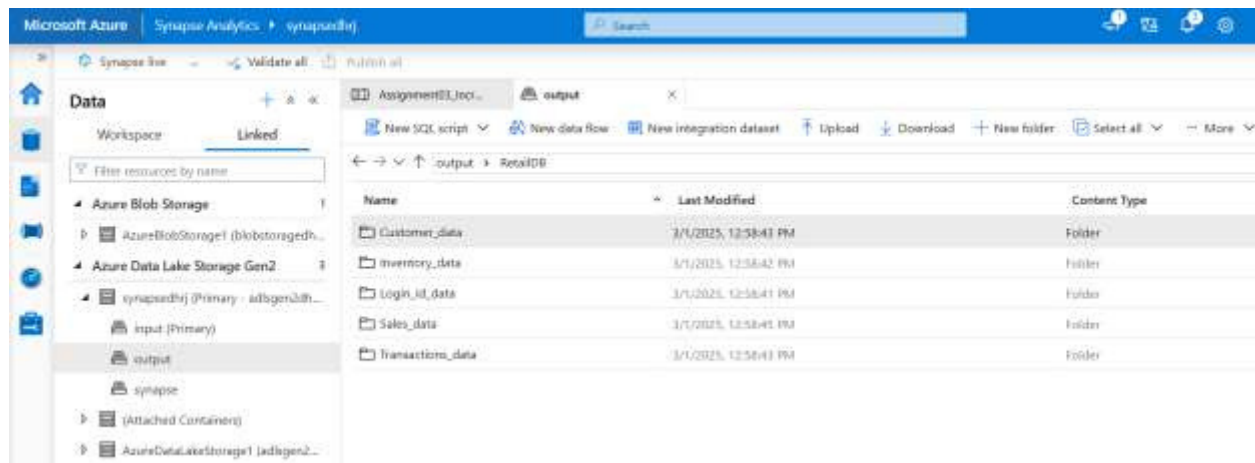
Activity name	Activity st...	Activit...	Run start	Duration	Integration runtime	User prop...	Activity run ID
Stored procedure1	<span style="color: green;">✔</span> Succeeded	Stored procedu...	3/1/2023, 1:00:20 PM	0s	AzureDataLakeIntegrationRuntime (Australia East)		887b4733-3461-4856-8389-c35286765099
Stored procedure1	<span style="color: green;">✔</span> Succeeded	Stored procedu...	3/1/2023, 1:00:19 PM	3s	AzureDataLakeIntegrationRuntime (Australia East)		D8682646-582a-47ac-a438-ca4177176e49
Stored procedure1	<span style="color: green;">✔</span> Succeeded	Stored procedu...	3/1/2023, 1:00:19 PM	5s	AzureDataLakeIntegrationRuntime (Australia East)		4c2b0281-a094-4316-8725-31e32a3a3978
Stored procedure1	<span style="color: green;">✔</span> Succeeded	Stored procedu...	3/1/2023, 1:00:18 PM	0s	AzureDataLakeIntegrationRuntime (Australia East)		97fa4895-7c1a-4946-8623-e6796128488
Stored procedure1	<span style="color: green;">✔</span> Succeeded	Stored procedu...	3/1/2023, 1:00:17 PM	3s	AzureDataLakeIntegrationRuntime (Australia East)		ea3b1f52-87fa-4180-9236-2b6e41f9b3d
Copy data1	<span style="color: green;">✔</span> Succeeded	Copy data	3/1/2023, 1:00:01 PM	17s	AzureDataLakeIntegrationRuntime (Australia East)		d850e02b-7800-4952-b502-82e7b5a02wa5
Copy data1	<span style="color: green;">✔</span> Succeeded	Copy data	3/1/2023, 1:00:00 PM	19s	AzureDataLakeIntegrationRuntime (Australia East)		3ada7229-b8c7-4449-8417-8a2a5709443a
Copy data1	<span style="color: green;">✔</span> Succeeded	Copy data	3/1/2023, 1:00:00 PM	10s	AzureDataLakeIntegrationRuntime (Australia East)		64d612a4-6887-400b-95c3-7375784788b3
Copy data1	<span style="color: green;">✔</span> Succeeded	Copy data	3/1/2023, 1:00:58 PM	19s	AzureDataLakeIntegrationRuntime (Australia East)		7e993e60-8c5e-46c5-9447-90e55c7a2a3c
Copy data1	<span style="color: green;">✔</span> Succeeded	Copy data	3/1/2023, 1:00:58 PM	17s	AzureDataLakeIntegrationRuntime (Australia East)		4704d02d-0413-4a4a-b279-4231b67e1d3a
Getmanuscript	<span style="color: green;">✔</span> Succeeded	Lookup	3/1/2023, 1:00:53 PM	7s	AzureDataLakeIntegrationRuntime (Australia East)		05d84943-2d53-4547-9420-1c5c0e6476d4
Getmanuscript	<span style="color: green;">✔</span> Succeeded	Lookup	3/1/2023, 1:00:55 PM	0s	AzureDataLakeIntegrationRuntime (Australia East)		a4b43648-199a-4639-683e-44b088b076c

Checking the output folders...

- RetailDB folder is created



- As mentioned in the Watermark table all the corresponding folders have been created i.e. Customer\_data folder, Inventory\_data folder, Login\_id\_data folder, Sales\_data folder and Transactions\_data folder.



- Cross checking the folder names in the SSMS.
- Checking in SSMS whether the name are matching with the Foldername

Assignment 03.sql - sqlserver-dheeraj/database/windows.net/sqlserver-dheeraj (Server (71)) - Microsoft SQL Server Management Studio

File Edit View Query Project Tools Window Help

sqlserver-dheeraj Execute

```

85 CREATE TABLE Watermark (
86     ID INT PRIMARY KEY,
87     Table_name VARCHAR(100),
88     Schema_name VARCHAR(100),
89     Foldername VARCHAR(50),
90     Delta_column VARCHAR(100),
91     Last_processed_value VARCHAR(255) NOT NULL
92 );
93
94
95 SELECT * FROM Watermark
96
97 INSERT INTO Watermark VALUES
98 (1,'Customer','dbo','RetailDB/Customer_data','Customerupdateddate','1900-01-01 00:00:00'),
99 (2,'Login_id','dbo','RetailDB/Login_id_data','Updatedlogindata','1900-01-01 00:00:00'),
100 (3,'Inventory','dbo','RetailDB/Inventory_data','ProductID','0'),
101 (4,'Transactions','dbo','RetailDB/Transactions_data','TransactionID','0'),
102 (5,'Sales','dbo','RetailDB/Sales_data','Sales_log','1900-01-01 00:00:00')

```

100 %

Results Messages

ID	Table_name	Schema_name	Foldername	Delta_column	Last_processed_value
1	Customer	dbo	RetailDB/Customer_data	Customerupdateddate	2023-06-01T10:00:00
2	Login_id	dbo	RetailDB/Login_id_data	Updatedlogindata	2023-03-01T08:15:00
3	Inventory	dbo	RetailDB/Inventory_data	ProductID	0
4	Transactions	dbo	RetailDB/Transactions_data	TransactionID	3
5	Sales	dbo	RetailDB/Sales_data	Sales_log	2023-06-15T10:00:00

- Checking for the Incremental load.
- Adding more values into Customer, Inventory and Sales tables and publish in Synapse to check whether the incremental load is working or not ?

Assignment 03.sql - ...heeraj (Server (71))\*

```

124 END
125
126
127 INSERT INTO Customer (CustomerID, Name, Phone, Customerupdateddate)
128 VALUES
129 (4, 'Alice Brown', '555-123-4567', '2023-04-01 12:00:00'),
130 (5, 'Charlie Davis', '555-789-0123', '2023-05-01 11:00:00'),
131 (6, 'Emily Wilson', NULL, '2023-06-01 10:00:00');
132
133 INSERT INTO Inventory (ProductID, ProductName, Quantity, Price)
134 VALUES
135 (4, 'Tablet', 30, 299.99),
136 (5, 'Smartwatch', 40, 199.99),
137 (6, 'Earbuds', 60, 89.99);
138
139
140 INSERT INTO Sales (SalesID, ProductID, Sales_log, Revenue)
141 VALUES
142 (4, 4, '2023-04-05 12:00:00', 299.99),
143 (5, 5, '2023-05-10 11:00:00', 199.99),
144 (6, 6, '2023-06-15 10:00:00', 89.99);
145

```

108 %

Results Messages

	ID	Table_name	Schema_name	Foldername	Delta_column	Last_processed_value
1	1	Customer	dbo	RetailDB/Customer_data	Customerupdateddate	2023-06-01T10:00:00
2	2	Login_id	dbo	RetailDB/Login_id_data	Updatedlogindata	2023-03-01T09:15:00
3	3	Inventory	dbo	RetailDB/Inventory_data	ProductID	6
4	4	Transactions	dbo	RetailDB/Transactions_data	TransactionID	3
5	5	Sales	dbo	RetailDB/Sales_data	Sales_log	2023-06-15T10:00:00

- Verifying the output values in the corresponding file locations.
- **Inventory table**

The screenshot shows a data management interface. On the left is a sidebar with a 'Data' section containing a tree view of folders and files. The main panel on the right displays details for a selected file: 'Inventory\_2025-03-01T18:09:00.3575607Z.csv'. Below the file name, there is a 'Path' field showing a long file path, a 'Modified' timestamp of '01/03/25, 18:09:00', and a 'With column header' toggle switch that is turned on. A table with four columns is displayed: 'PRODUCTID', 'PRODUCTNAME', 'QUANTITY', and 'PRICE'. The table contains three data rows and a total row at the bottom. A progress bar at the bottom of the table indicates 100% completion. A blue button is located at the bottom left of the main panel.

PRODUCTID	PRODUCTNAME	QUANTITY	PRICE
1	Apple	10	100.00
2	Banana	20	50.00
3	Orange	30	30.00
TOTAL		60	180.00

- **Customer table**

Microsoft Dynamics CRM - My Account

Home My Account My Recent Items

My Account

Customer\_2025-05-01T18:08:59.3637496Z.crm

Name: Customer\_2025-05-01T18:08:59.3637496Z.crm

Modified: 5/1/2025, 3:08:59 PM

With external provider: ☒ On

CUSTOMER	NAME	PHONE	CUSTOMER
A	Apple Store	555-123-4567	2025-05-01T18:08:59.3637496Z.crm
F	Facebook Group	555-789-0123	2025-05-01T18:08:59.3637496Z.crm
B	Google Play	555-456-7890	2025-05-01T18:08:59.3637496Z.crm

Progress bar: 100%

Go

- Sales table

Microsoft Azure | Import Wizard | Importing

File name: Sales\_2025-03-01T18:28:39.0549193Z.csv

Path: https://adfsapoc01.blob.core.windows.net/adfsapoc01/users/01-01-01T18:28:39.0549193Z.csv

Modified: 01/01/2025 18:03:04

With column header: ☒ On

SACID80	PROXOC750	SALES_LOS	REVENUE
1	1	2025-04-01 12:1	100.00
2	2	2025-04-01 12:1	100.00
3	3	2025-04-01 12:1	100.00
4	4	2025-04-01 12:1	100.00

1 of 4 rows