

STATEMENT 2

Every time the pipeline runs it will overwrite the data. Ensure that all the data goes into Folder like (Customer/Year/Month/Day) to avoid any overwriting. Please rewrite the pipeline.

In the previous statement/problem we built a pipeline, which copied data from Customer table, Azure SQL into ADLS(Azure Data Lake Storage) storage in JSON format. However, each time the operation is performed, the data is overwritten in the file. Therefore, now, we have to store the data in folder hierarchy: Customer/Year/Month/Day, where, within the 'Custoemr' folder we have to dynamically create the Year, Month and Day folder corresponding to the time at which the copy operation is executed.

For eg, If we are performing the operation on 20th July, 2025, the data will be saved in **Customer/2025/July/20**.

Most of the operations and activities to be performed are similar to the previous statement, we just need to create dynamic subdirectories according to the time of pipeline execution.

Prerequisites:

1. Threshold file in ADLS location

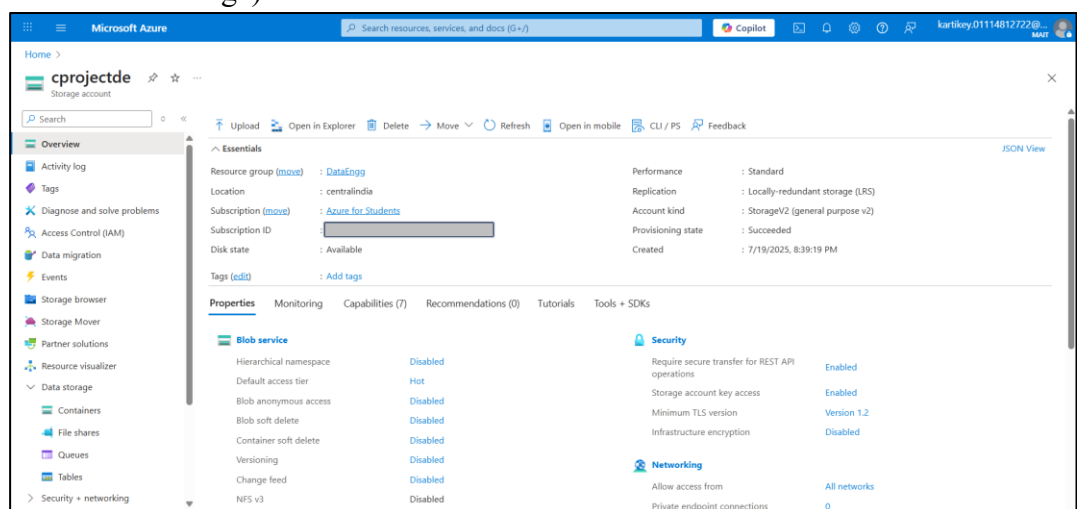
For storing the file in ADLS location, we need to create a **storage account** in Azure. Provide the

Subscription

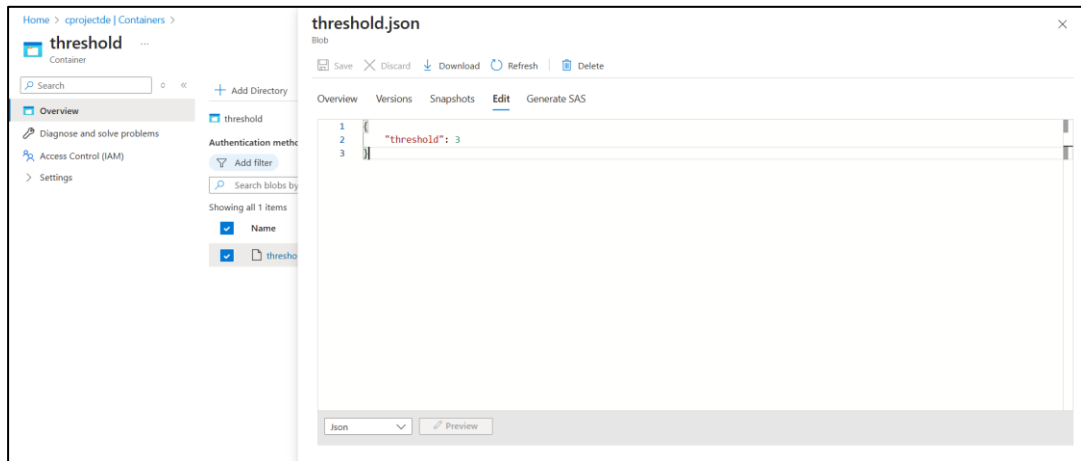
Resource Group

Storage Account Name

Set **Redundancy** accordingly.(account used in this project has Locally-Redundant Storage)



Create a threshold file in Containers



The image shows **threshold.json** file in threshold directory of **cprojectde** storage account, with a threshold value of 3.

2. Customer table in Azure SQL

Creating table:

```
CREATE TABLE customer (  
    customer_id INT PRIMARY KEY,  
    first_name VARCHAR(50),  
    last_name VARCHAR(50),  
    email VARCHAR(100),  
    created_date DATETIME  
);
```

Inserting values:

```
INSERT INTO customer VALUES  
(1, 'Rahul', 'Gupta', 'rahul@mail.com', '2025-07-11'),  
(2, 'Perna', 'Aggarwal', 'perna@mail.com', '2025-07-12'),  
....  
(1009, 'Manish', 'Gupta', 'manish.gupta@mail.in', '2025-08-25'),  
(1010, 'Divya', 'Singh', 'divya.singh@mail.in', '2025-08-26'),
```

The full scripts are present in **insertCustomers.sql** and **customerTable.sql**

Total records in Customer table: **71**

Add the files and table created as datasets in your Azure Data Factory

NOTE: If no data factory has been created yet, see Page 5

1. Threshold File Dataset and target file dataset

Create **Azure Data Lake Storage Gen2** dataset.

The format of the dataset should be **JSON**

The screenshot shows the configuration page for a JSON dataset named 'P_DataLakeJSON'. The 'Connection' tab is selected, showing the linked service 'P_AzureDataLakeStorage'. The 'File path' is configured with '@dataset().FolderPath' for the directory and '@dataset().FileName' for the file name. The 'Compression type' is set to 'No compression' and the 'Encoding' is set to 'Default(UTF-8)'.

Property	Value
Linked service	P_AzureDataLakeStorage
File path	@dataset().FolderPath / Directory / @dataset().FileName
Compression type	No compression
Encoding	Default(UTF-8)

Create two parameters:

The screenshot shows the 'Parameters' tab for the dataset 'P_DataLakeJSON'. It displays two parameters: 'FolderPath' and 'FileName', both of type 'String' with a default value of 'Value'.

Name	Type	Default value
FolderPath	String	Value
FileName	String	Value

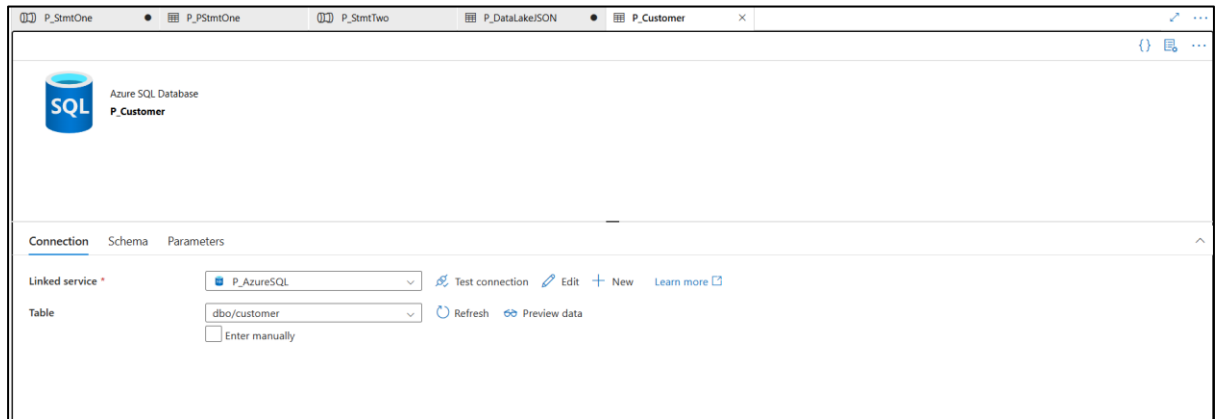
Populate the **File Path**:

File System: @dataset().FolderPath

File Name: @dataset().FileName

2. SQL Table Dataset

Create Azure SQL dataset



Provide:

Dataset Name

Linked Service

Table Name

Set properties

Name
DatasetName

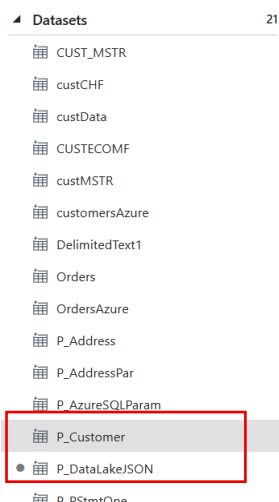
Linked service *
P_AzureSQL

Table name
dbo.customer

☐ Enter manually

Import schema
☒ From connection/store ☐ None

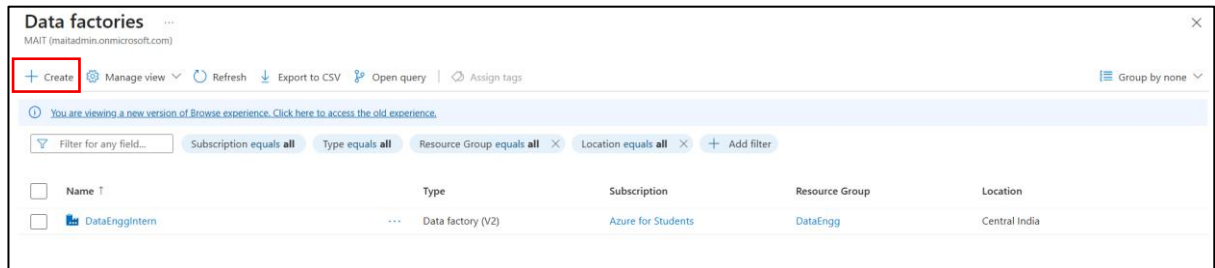
Final



Steps taken to achieve the required goal:

We have to create a pipeline in Azure Data Factory.

NOTE: Create an Azure Data Factory if not already created.



Home > Data factories >

Create Data Factory

Basics | Git configuration | Networking | Advanced | Tags | Review + create

One-click to create data factory with sample pipeline and datasets. [Try it](#)

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription * ⓘ

Resource group * ⓘ

[Create new](#)

Instance details

Name * ⓘ

Region * ⓘ

Version * ⓘ

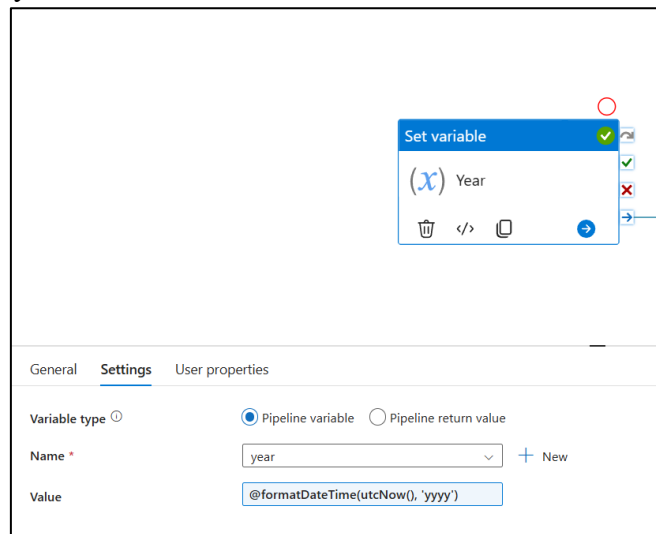
Azure for Students

East US

V2

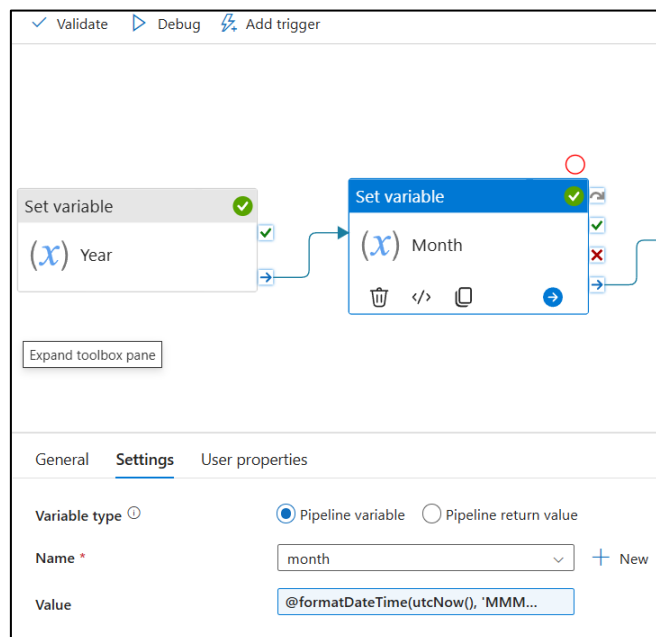
After filling the details, click, **Review + Create**

1. Create new pipeline, name it accordingly.
2. We will create three variables, that store the year, month and date of the pipeline execution. Activity to be used is **Set Variable**



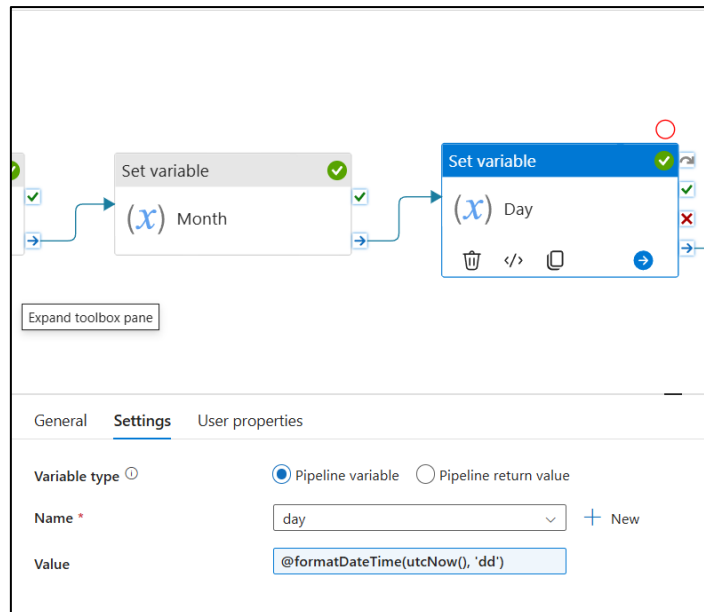
Variable: year

Value: @formatDateTime(utcNow(), 'yyyy')



Variable: month

Value: @formatDateTime(utcNow(), 'MMMM')



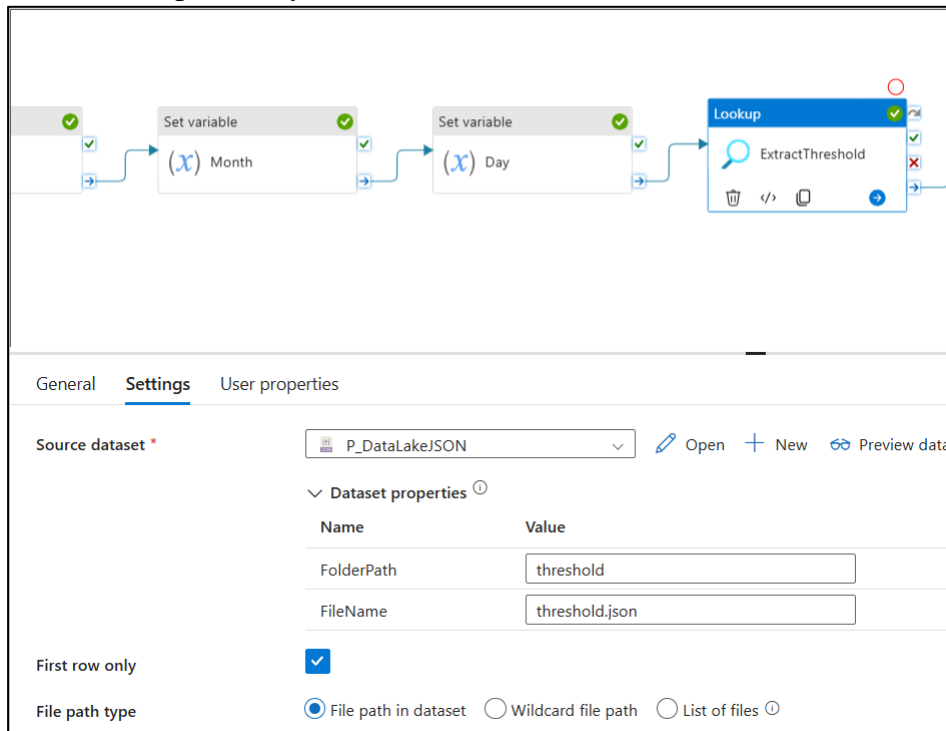
Variable: day

Value: @formatDateTime(utcNow(), 'dd')

Pipeline Variables Initialised as:

Parameters	Variables	Settings	Output
<div>+ New Delete</div>			
<input type="checkbox"/>	Name	Type	Default value
<input type="checkbox"/>	year	String	Value
<input type="checkbox"/>	month	String	Value
<input type="checkbox"/>	day	String	Value
<input type="checkbox"/>	threshold	Integer	Value

3. Insert **Lookup** Activity



The name of the **Lookup** activity as shown above is ExtractThreshold.

The source dataset is the dataset, where threshold.json is present.

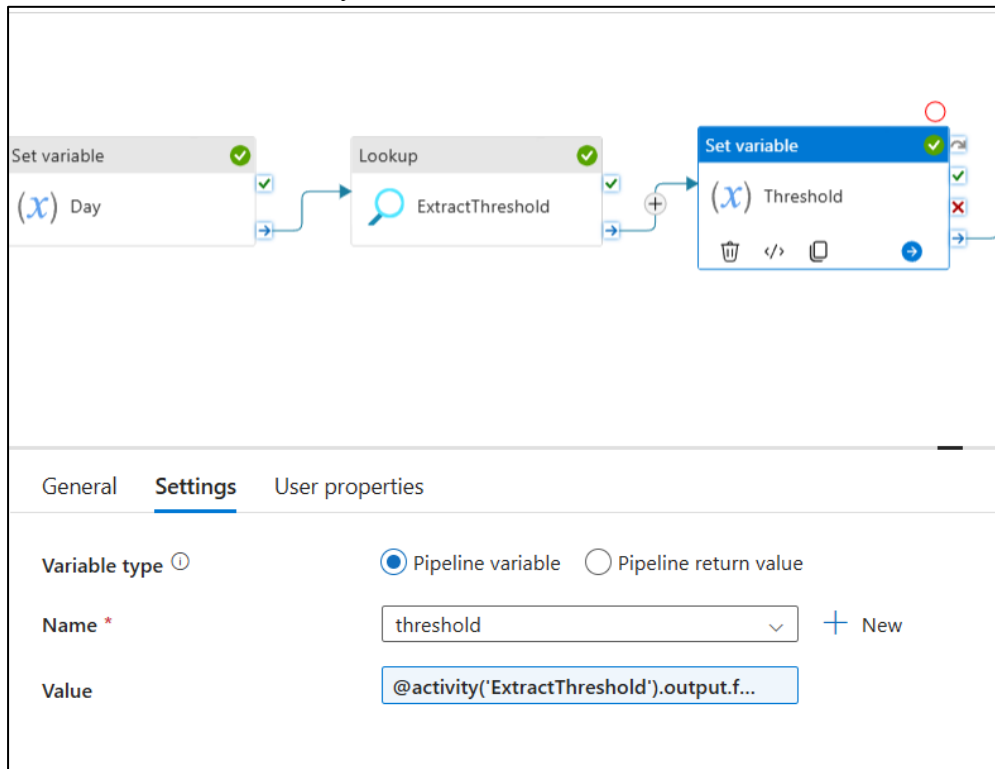
We can extract the file using the parameter we created during dataset creation([Page 3](#))

The parameters

FolderPath: Name of the directory in which json file is present(can be left blank if file is not in directory)

FileName: Name of the file(eg. threshold.json)

4. Insert **Set Variable** activity



Value of variable: @activity('ExtractThreshold').output.firstRow.threshold
This variable extracts the threshold value from the output of Lookup Activity.
The value is according to the file format, in our case the file is:

```
{  
  "threshold": 3  
}
```

Where, the first row contains the **threshold** parameter, we extract this using firstRow.threshold (see **Value**)

For this, initialise variable 'threshold' in the pipeline of type integer. (as shown on [Page 7](#))

5. Insert **Lookup** Activity

The screenshot shows a pipeline with three activities: a 'Lookup' activity named 'ExtractThreshold', a 'Set variable' activity named 'Threshold', and another 'Lookup' activity named 'ExtractTotalRecords'. The 'ExtractTotalRecords' activity is selected, and its settings are shown in the bottom panel. The settings include:

- Source dataset: P_Customer
- First row only: ☒
- Use query: ☒ Query
- Query: `SELECT count(*) AS 'totalrecords' FR...`
- Query timeout (minutes): 15
- Isolation level: Select...

This Lookup activity extracts the total number of records present in the Customer table in Azure SQL.

Where,

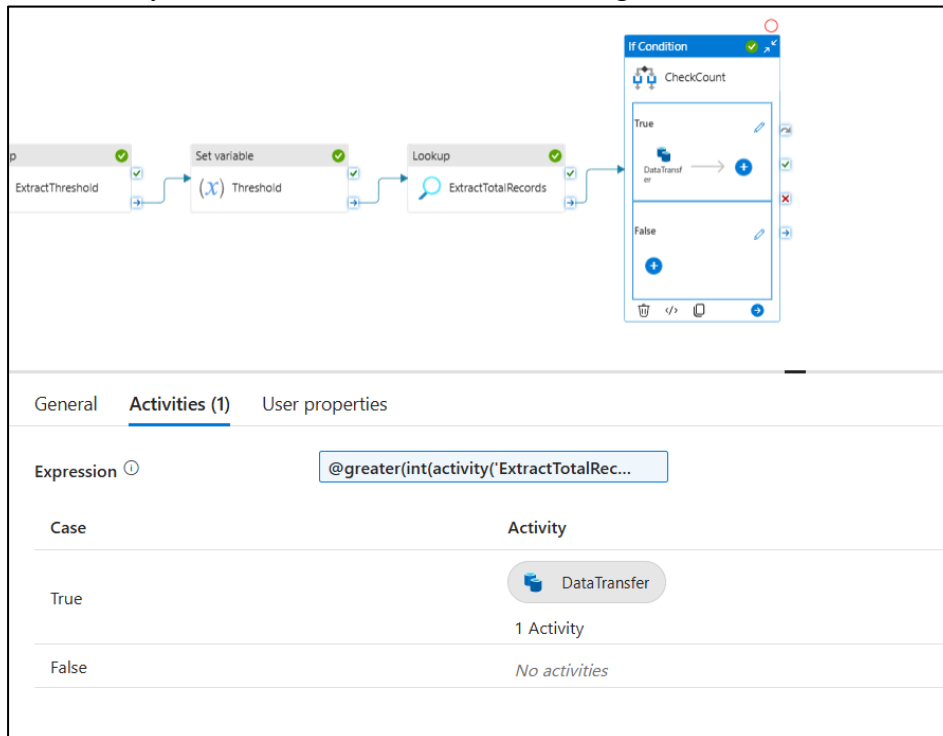
Source Dataset is the customer table(we created earlier, [Page 4](#))

Query: To count the total number of records

`SELECT count(*) AS 'totalrecords' FROM dbo.Customer;`

6. Insert **If Condition** Activity

This activity checks, if the total record count is greater than the threshold value or not.



The expression is:

**@greater(int(activity('ExtractTotalRecords').output.firstRow.totalrecords),
int(variables('threshold')))**

It compares the variable threshold and the output from previous activity, i.e, totalRecords from Lookup activity.

If the record count is greater than the threshold value, we need to perform data copy operation. Therefore we need to add **Copy Data** activity in the True condition.

7. Insert **Copy Data** activity in True Condition

Validate Debug Add trigger

Lookup ExtractTotalRecords

If Condition

CheckCount

True

DataTransfer

False

General Source Sink Mapping Settings User properties

Source dataset * P_Customer Open + New Preview data

Use query ☒ Table ☐ Query ☐ Stored procedure

Query timeout (minutes) 15

Isolation level Select...

Partition option ☒ None ☐ Physical partitions of table ☐ Dynamic range

Where,

Source dataset is the customer table(created on [Page 4](#))

ExtractThreshold Set variable (x) Threshold Lookup ExtractTotalRecords If Condition CheckCount True DataTransfer False

General Source Sink Mapping Settings User properties

Sink dataset * P_DataLakeJSON Open + New Learn more

Dataset properties

Name	Value
FolderPath	@concat('customer/', variables('year'...
FileName	customer_data.json

Copy behavior Preserve hierarchy

Max concurrent connections

Where,

Sink Dataset is the target file we created earlier(see [Page 5](#))

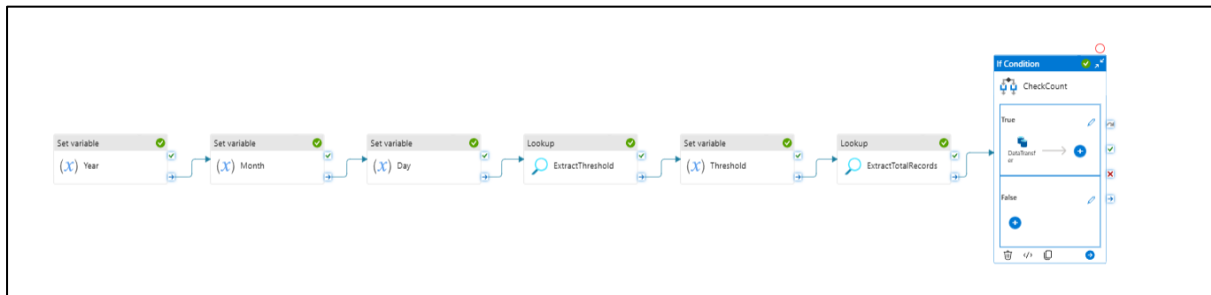
Parameters:

FolderPath: @concat('customer/', variables('year'), '/', variables('month'), '/', variables('day'), '/')

FileNname: customer_data.json

Preserve the hierarchy and create the file as **Array of objects**

Final Pipeline:





Pipeline Execution

Activity runs

Pipeline run ID 179b8004-dba0-4762-94d4-d98054711c46

All status ▾ List ▾

Monitor in Azure Metrics  

Showing 1 - 8 of 8 items

Activity name	Activity st...	Activit...	Run start	Duration	Integration runtime	User prop...
DataTransfer	✔ Succeeded	Copy data	7/20/2025, 5:44:16 PM	13s	AutoResolveIntegrationRuntime (Central India)	
CheckCount	✔ Succeeded	If Condition	7/20/2025, 5:44:14 PM	18s		
ExtractTotalRecords	✔ Succeeded	Lookup	7/20/2025, 5:43:44 PM	29s	AutoResolveIntegrationRuntime (Central India)	
Threshold	✔ Succeeded	Set variable	7/20/2025, 5:43:43 PM	Less than 1s		
ExtractThreshold	✔ Succeeded	Lookup	7/20/2025, 5:43:28 PM	14s	AutoResolveIntegrationRuntime (Central India)	
Day	✔ Succeeded	Set variable	7/20/2025, 5:43:27 PM	Less than 1s		
Month	✔ Succeeded	Set variable	7/20/2025, 5:43:26 PM	Less than 1s		
Year	✔ Succeeded	Set variable	7/20/2025, 5:43:26 PM	Less than 1s		

Step-Wise Output:

1. Set Variables Activity:
- Variable **year**:

Output

 Copy to clipboard

```
{
  "name": "year",
  "value": "2025"
}
```

Variable **month**:

Output

 Copy to clipboard

```
{
  "name": "month",
  "value": "July"
}
```

Variable **day**:

Output ↗ ✕

📄 Copy to clipboard

```
{
  "name": "day",
  "value": "20"
}
```

2. Lookup Activity(ExtractThreshold)

Output ↗ ✕

📄 Copy to clipboard

```
{
  "firstRow": {
    "threshold": 3
  },
  "effectiveIntegrationRuntime":
  "AutoResolveIntegrationRuntime (Central India)",
  "billingReference": {
    "activityType": "PipelineActivity",
    "billableDuration": [
```

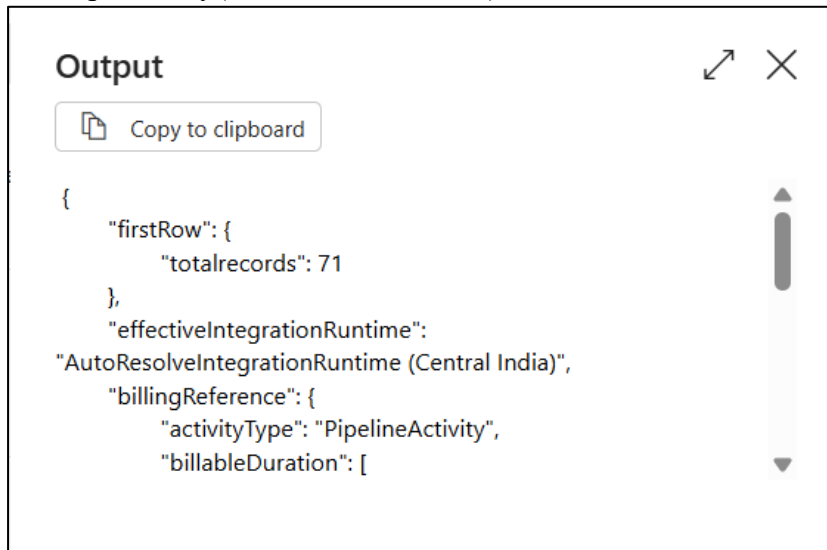
3. Set Variable Activity(Threshold)

Output ↗ ✕

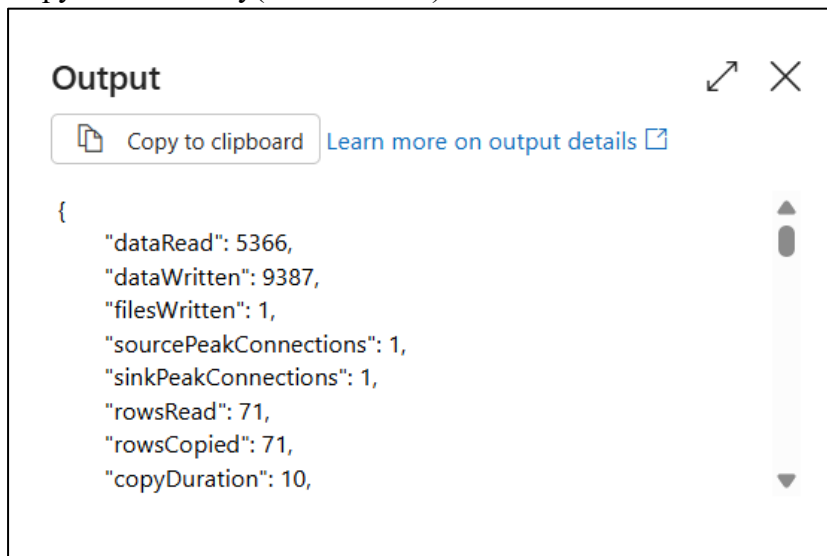
📄 Copy to clipboard

```
{
  "name": "threshold",
  "value": 3
}
```

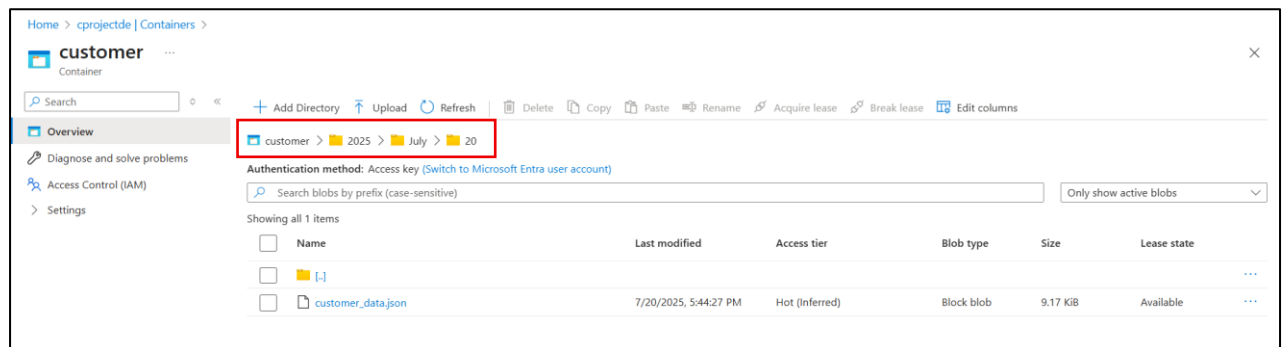
4. Lookup Activity(ExtractTotalRecords)



5. Copy Data Activity(DataTransfer)

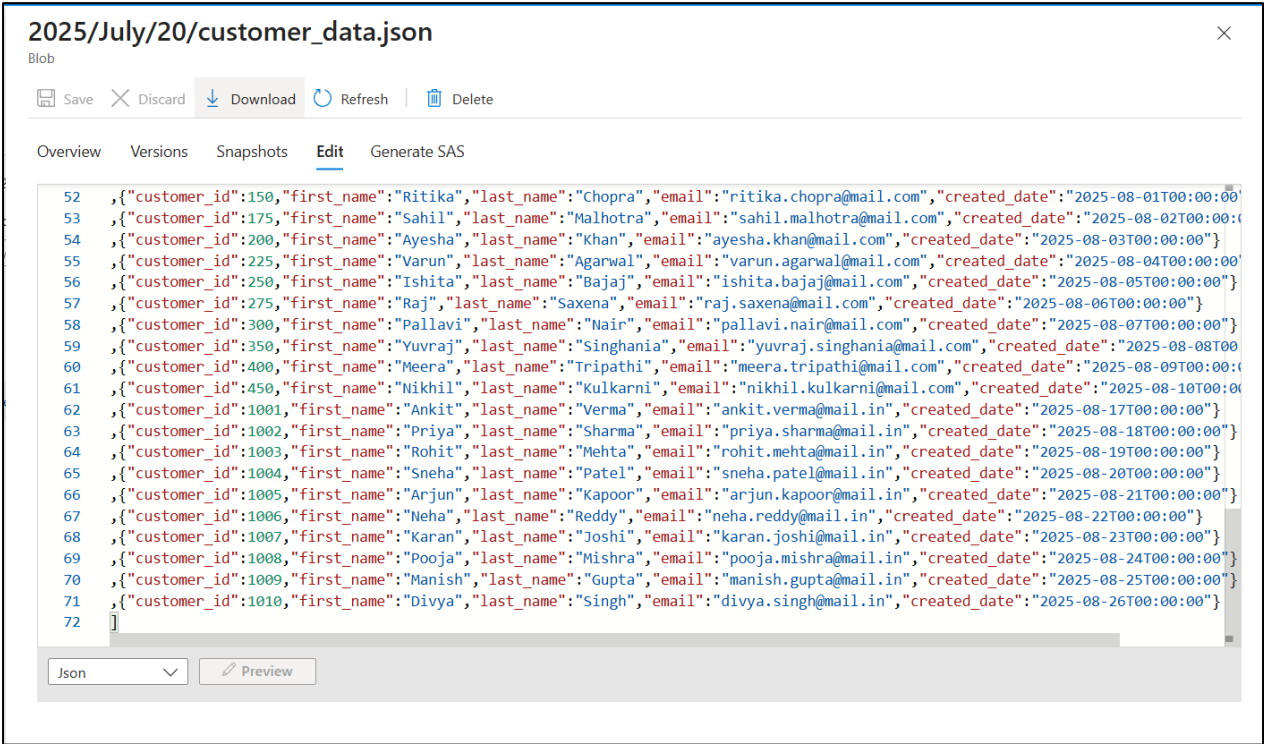


Target JSON file:



File has been created in **customer/2025/July/20** directory of **cprojectde** storage account.

customer_data.json File:



The pipeline was also executed on 19th July, 2025. After executing pipeline on 20th July 2025, the old data was not overwritten.




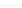

customer > 2025 > July

Authentication method: Access key (Switch to Microsoft Entra user account)

Search blobs by prefix (case-sensitive)

Only show active blobs

Showing all 4 items

<input type="checkbox"/>	Name	Last modified	Access tier	Blob type	Size	Lease state	
<input type="checkbox"/>	 [-]						...
<input type="checkbox"/>	 19						...
<input type="checkbox"/>	 20						...
<input type="checkbox"/>	 19	7/19/2025, 11:11:14 PM	Hot (Inferred)	Block blob	0	Available	...
<input type="checkbox"/>	 20	7/20/2025, 5:44:27 PM	Hot (Inferred)	Block blob	0	Available	...



There were 5 records when executed on 19th July, 2025.