



Cache Sink Implementation

In Azure Data Factory, the Cache Sink is a component used within Mapping Data Flows to optimize data processing performance by caching intermediate data during transformations. Here's a brief overview of its implementation:

1. **Purpose:** The Cache Sink is used to store intermediate results of data transformations within a Mapping Data Flow. By caching data, it reduces the need to reprocess the same data multiple times, improving overall performance and reducing data movement.
2. **Configuration:** To implement a Cache Sink in Azure Data Factory, you add it to your Mapping Data Flow and configure it to cache the data at a specific point in your data transformation logic. You specify the columns to cache and the cache storage settings.
3. **Usage:** Once configured, the Cache Sink temporarily stores the data during the execution of the data flow. Subsequent transformation steps can then access this cached data without needing to reprocess the original source data.
4. **Benefits:** Using a Cache Sink can significantly improve performance, especially in scenarios where multiple transformations are applied to the same dataset. By caching intermediate results, it reduces the computational overhead and data movement within the data flow, leading to faster processing times.
5. **Considerations:** While Cache Sink can enhance performance, it's essential to use it judiciously, as caching large datasets can consume additional storage resources. It's also crucial to manage cache expiration and refresh policies to ensure that the cached data remains up-to-date.



Use cases of Cache Sink:

The Cache Sink in Azure Data Factory's Mapping Data Flows finds application in various scenarios where optimizing data processing performance is crucial. Here are some common use cases:

1. **Complex Transformations:** When implementing complex data transformations involving multiple steps or iterations, the Cache Sink can be used to store intermediate results. This prevents redundant processing of the same data, reducing computational overhead and improving overall efficiency.
2. **Aggregations and Joins:** In scenarios where aggregations or joins are performed on large datasets, caching intermediate results using the Cache Sink can significantly improve performance. It allows subsequent transformations to access pre-computed results instead of recalculating them, leading to faster processing times.
3. **Repetitive Data Access:** If certain data subsets are accessed multiple times within a data flow, caching them using the Cache Sink can eliminate the need for repeated data reads from the source. This is particularly useful in scenarios where source data is located in remote or slow-access systems.
4. **Iterative Development:** During the development and testing phases of a data integration project, the Cache Sink can be used to temporarily store intermediate data

for debugging and validation purposes. It allows developers to inspect and analyze intermediate results without rerunning the entire data flow repeatedly.

5. **Real-Time Data Processing:** In real-time data processing pipelines, where low-latency processing is essential, caching intermediate results using the Cache Sink can help meet stringent performance requirements. It enables faster data processing by reducing the computational overhead associated with complex transformations.
6. **Batch Processing Optimization:** When dealing with batch processing of large volumes of data, caching intermediate results using the Cache Sink can streamline processing and reduce overall execution times. It ensures that previously processed data is readily available for subsequent transformations, minimizing redundant work.
7. **Data Deduplication:** In scenarios where duplicate records need to be identified and removed from datasets, caching intermediate results using the Cache Sink can facilitate efficient deduplication processes. By storing unique identifiers or hash values of processed records, it enables faster comparison and elimination of duplicates.
8. **Incremental Loading:** When loading data incrementally into a target system, caching intermediate results using the Cache Sink can improve performance by reducing the need to reprocess unchanged data. It allows for efficient delta processing by caching previously processed records and identifying new or modified ones.

In this guide, we're implementing the Cache Sink feature in Azure Data Factory's Mapping Data Flows to optimize data processing performance. The end goal is to reduce computational overhead and improve efficiency by caching intermediate results during data transformations. By following the provided steps, users can configure Cache Sink, apply it to various use cases such as complex transformations and aggregations, and ultimately streamline their data integration processes for faster and more efficient results.

To begin with the Lab:

1. In the last lab we added surrogate keys to our Dimension Product Table.
2. This time we will be doing the same for the Dimension Customer Table.
3. First, we are going to drop our customer table and create a new one with the Product SK column in it.
4. At the same time since we are going to run our pipeline again, we have to delete the data from the product table also.

```
SQLQuery1.sql - dat... (sqladminuser (0)) * X
DROP TABLE [dbo].[DimCustomer]

CREATE TABLE [dbo].[DimCustomer](
    [CustomerSK] [int] NOT NULL,
    [CustomerID] [int] NOT NULL,
    [CompanyName] varchar(200) NOT NULL,
    [SalesPerson] varchar(300) NOT NULL
)
WITH
(
    DISTRIBUTION = REPLICATE
)

150 % <
Messages
Commands completed successfully.

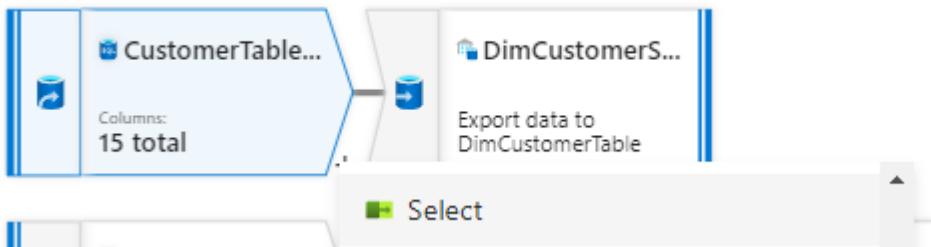
Completion time: 2024-05-01T15:20:50.6433869+05:30
```

```
Delete From [dbo].[DimProduct]

150 % <
Messages
(295 rows affected)

Completion time: 2024-05-01T15:22:24.2218086+05:30
```

5. Now we are going to do the same implementation as before but this time it will be our customer table.
6. So, head towards Data Factory Wizard and open your dimension dataflow.
7. Now choose your customer table stream and then click on the plus icon and choose select column modifier just to select the columns of our interest.



8. Then you need to give it name and after that scroll down to the mappings.

Select settings Optimize Inspect Data preview

Output stream name *	<input type="text" value="CustomerSelectColumns"/>	Learn more
Description	Renaming CustomerTableStream to CustomerSelectColumns with columns 'CustomerID, NameStyle, Title,'	
Incoming stream *	<input type="text" value="CustomerTableStream"/>	
Options	<input checked="" type="checkbox"/> Skip duplicate input columns ? <input checked="" type="checkbox"/> Skip duplicate output columns ?	

9. Here we just need the customer ID, company name, and salesperson. You have to remove the other mappings.

<input checked="" type="checkbox"/> CustomerTableStream's column	Filter	Name as	Filter
<input type="checkbox"/> 123 CustomerID	▼	<input type="text" value="CustomerID"/>	+ Delete
<input type="checkbox"/> abc CompanyName	▼	<input type="text" value="CompanyName"/>	+ Delete
<input type="checkbox"/> abc SalesPerson	▼	<input type="text" value="SalesPerson"/>	+ Delete

10. Then we need to add the surrogate key. For the customer select columns click on plus icon and choose surrogate key.



11. Now you need to give your stream a name then choose customer select columns as your incoming stream and in the key column give Customer SK.

Surrogate key settings Optimize Inspect Data preview

Output stream name * [Learn more](#)

Description [Reset](#)

Incoming stream *

Key column *

Start value *

Step value

12. Now you need to go to the sink and click on open.

```
graph LR; A[CustomerTable...] --> B[CustomerSelectC...]; B --> C[CustomerSKSurrogateKeyStream]; C --> D[DimCustomerS...]
```

Sink Settings Errors Mapping Optimize Inspect Data preview

Output stream name * [Learn more](#)

Description [Reset](#)

Incoming stream *

Sink type * Dataset Inline Cache

Dataset * [Test connection](#) [Open](#) [New](#)

Options Allow schema drift [?](#) Validate schema [?](#)

13. Then go to schema and click on import schema.

Connection Schema Parameters

Column name	Type
CustomerSK	int
CustomerID	int
CompanyName	varchar
SalesPerson	varchar

14. After that go to mappings and click on reset.

Sink Settings Errors Mapping Optimize Inspect Data preview

⚠ At least one incoming column is mapped to a column in the sink dataset schema with a conflicting type, which can cause NULL values or runtime errors.

Options Skip duplicate input columns ⓘ

Skip duplicate output columns ⓘ

Auto mapping ⓘ ↻ Reset + Add mapping Delete Output format

Input columns	Output columns
121 CustomerSK	123 CustomerSK
123 CustomerID	123 CustomerID
abc CompanyName	abc CompanyName
abc SalesPerson	abc SalesPerson

15. Then you have to click on Publish All.

16. Now you have to go to your CSV container in your storage account. Here you have to click on Add directory. Then name it and save it.

CSV Container

+ Add Directory Upload ⟳ Refresh ⟳ Rename Delete ⟲ Change tier ⚡ Acquire lease ⚡ Break lease ↗ Give feedback

Authentication method: Access key ([Switch to Microsoft Entra user account](#))
Location: CSV

Search blobs by prefix (case-sensitive)

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
Log.csv	1/5/2024, 12:36:53 pm	Hot (Inferred)		Block blob	6.64 MB	Available

17. After that you have to go inside that directory and upload two files. So, these two files have the data from the customer table which we have seen in our dimension table. The data is divided into both of them.

Upload + Add Directory ⟳ Refresh ⟳ Rename Delete ⟲ Change tier ⚡ Acquire lease ⚡ Break lease ↗ Give feedback

Authentication method: Access key ([Switch to Microsoft Entra user account](#))
Location: csv / Customer

Search blobs by prefix (case-sensitive)

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]						***
Customer01.csv	1/5/2024, 3:38:04 pm	Hot (Inferred)		Block blob	61.08 KiB	Available
Customer02.csv	1/5/2024, 3:38:04 pm	Hot (Inferred)		Block blob	113.44 KiB	Available

18. Now if you see in your customer table stream the dataset is pointing towards the customer table now we are going to create a new dataset for it. For that click on new.

The screenshot shows a data flow pipeline in the Azure Data Factory interface. The pipeline consists of four main activities connected sequentially:

- CustomerTableStream**: A source activity with 15 columns.
- CustomerSelectC...**: A Rename activity.
- CustomerSKSur...**: An activity with the note "Adding new key CustomerSK starting from 1 with step 1".
- DimCustomerS...**: A sink activity with the note "Export data to DimCustomerTable".

Below the pipeline, the "Source settings" tab is selected. The configuration includes:

- Output stream name ***: CustomerTableStream
- Description**: Import data from Customertable
- Source type ***: Dataset (selected)
- Dataset ***: Customertable (selected)
- Options**:
 - Allow schema drift
 - Infer drifted column types
 - Validate schema
- Sampling * ①**:
 - Enable
 - Disable

19. Then you have to choose Azure data lake gen2 storage account.

New dataset

In pipeline activities and data flows, reference a dataset to specify the location and structure of your data within a data store. [Learn more](#)

Select a data store

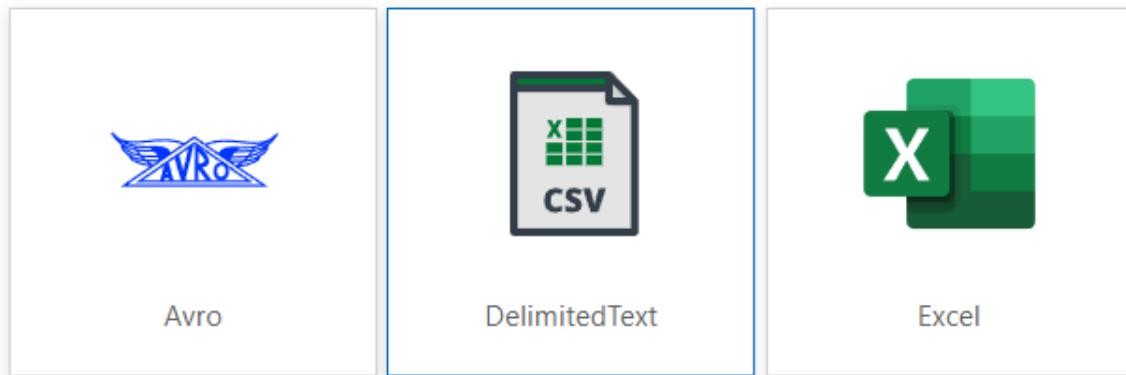
The screenshot shows the "New dataset" dialog in the Azure Data Factory interface. The search bar at the top contains the text "Data Lake". Below the search bar, there are tabs for "All", "Azure", "Database", "File", "Generic protocol", "NoSQL", and "Services and apps". The "All" tab is selected. The results section displays one item:

- Azure Data Lake Storage Gen2**: Represented by a thumbnail icon of a storage unit.

20. After that in the format choose delimited text.

Select format

Choose the format type of your data



21. Then in the set properties option you have to give it a name, browse for the CSV container, in that the new directory, and then choose Customer01.csv

Set properties

Name

CustomerCSV_Source

Linked service *

sqlstorage1010_service



File path

csv / Customer / Customer01.csv



First row as header



Import schema

From connection/store From sample file None

> Advanced

22. Once it is uploaded, you have to click on open, then move to schemas click on import schemas, and choose import schemas from the connection store.

Source settings Source options Projection Optimize Inspect Data preview

Output stream name * CustomerTableStream [Learn more](#)

Description Import data from CustomerCSV_Source [Reset](#)

Source type * Dataset Inline

Dataset * CustomerCSV_Source [Test connection](#) [Open](#) (boxed) [New](#)

Options Allow schema drift Infer drifted column types Validate schema

Skip line count

Sampling * Enable Disable

Connection Schema Parameters

[Import schema](#) [Clear](#)

[From sample file](#)

[From connection/store](#) (boxed)

23. Now go back to customer select column stream then scroll down and reset your schema. You should have the same schema as shown below.

Select settings Optimize Inspect Data preview

Output stream name * CustomerSelectColumns [Learn more](#)

Description Renaming CustomerTableStream to CustomerSelectColumns with columns 'CustomerId, CompanyName' [Reset](#)

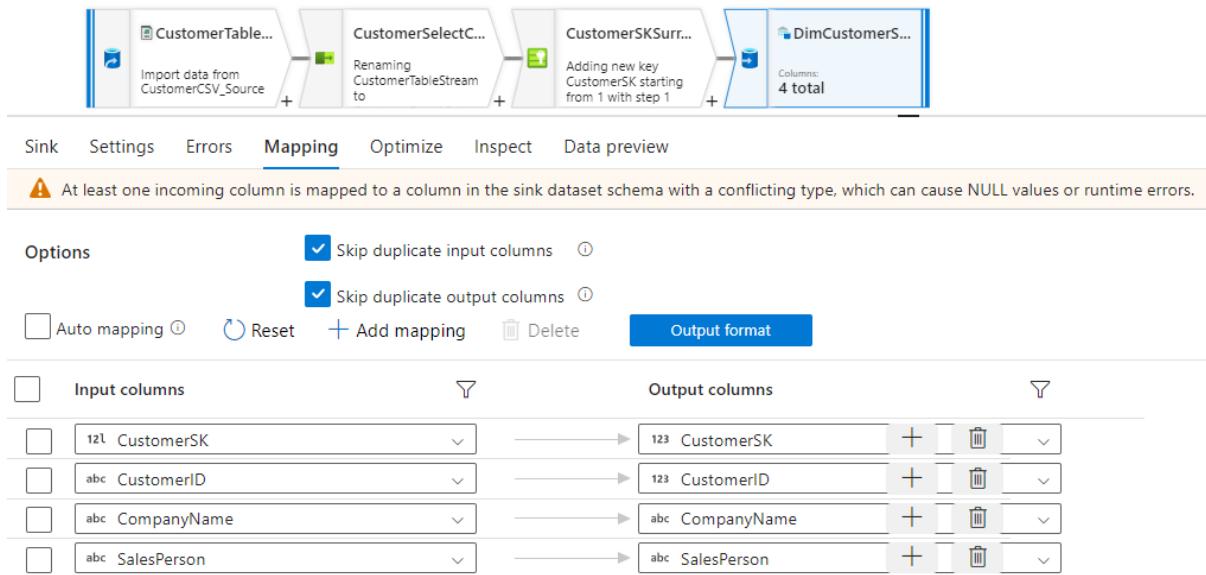
Incoming stream * CustomerTableStream

Options Skip duplicate input columns Skip duplicate output columns

Input columns * [Auto mapping](#) [Reset](#) [Add mapping](#) [Delete](#) 3 mappings: 11 column(s) from the inputs left unmapped

CustomerTableStream's column	Name as
abc CustomerID	CustomerId
abc CompanyName	CompanyName
abc SalesPerson	SalesPerson

24. Then go back to sink and in the mapping click on reset. You should have the same mappings as shown below. Then click on publish all.



25. After that go back to your pipeline and trigger it to run your pipeline. Below you can see that your pipeline run was successful.

Run Succeeded

Successfully ran 05_mapping_dataflow_dimtable (Pipeline).

[View pipeline run](#)

4 minutes ago

26. And below you can see that when we run the select statement with order, we can see the data accordingly. Here we get the entries added till customer ID 468 because our CSV file has this much information.

```
Select * from [dbo].[DimCustomer]  
Order by [CustomerSK]
```

150 %

Results Messages

	CustomerSK	CustomerID	CompanyName	SalesPerson
271	271	421	Fitness Cycling	adventure-works\shu0
272	272	424	Bikes for Two	adventure-works\jose1
273	273	425	Bikes for Kids and Adults	adventure-works\jose1
274	274	426	Custom Accessories Company	adventure-works\garrett1
275	275	430	Uttermost Bike Shop	adventure-works\jae0
276	276	431	Small Cycle Store	adventure-works\michael9
277	277	433	Thorough Parts and Repair ...	adventure-works\pamela0
278	278	434	Road-Way Mart	adventure-works\david8
279	279	435	Satin Finish Company	adventure-works\jillian0
280	280	436	Sheet Metal Manufacturing	adventure-works\jillian0
281	281	437	Professional Cycle Store	adventure-works\shu0
282	282	438	Remarkable Bike Store	adventure-works\linda3
283	283	439	Regional Manufacturing	adventure-works\shu0
284	284	442	Original Bicycle Supply Com...	adventure-works\jose1
285	285	443	Courteous Bicycle Specialists	adventure-works\jose1
286	286	444	Finer Cycle Shop	adventure-works\garrett1
287	287	448	Action Bicycle Specialists	adventure-works\jae0
288	288	451	Roadway Bike Emporium	adventure-works\pamela0
289	289	452	A Cycle Shop	adventure-works\david8
290	290	453	Unique Bikes	adventure-works\jillian0
291	291	454	Sleek Bikes	adventure-works\jillian0
292	292	455	Timely Shipping Service	adventure-works\shu0
293	293	456	Riding Excursions	adventure-works\linda3
294	294	457	Retail Sporting Equipment	adventure-works\shu0
295	295	460	Price-Cutter Discount Bikes	adventure-works\jose1
296	296	461	Active Life Toys	adventure-works\jose1
297	297	462	General Riding Supplies	adventure-works\garrett1
298	298	466	Central Bicycle Specialists	adventure-works\jae0
299	299	468	Blue Bicycle Company	adventure-works\michael9

Query executed successfully.

27. Now come back to our customer table Stream and click on open the dataset.

Source settings

Output stream name * CustomerTableStream [Learn more](#)

Description Import data from CustomerCSV_Source [Reset](#)

Source type * [Dataset](#) [Inline](#)

Dataset * [CustomerCSV_Source](#) [Test connection](#) [Open](#) [New](#)

Options

- Allow schema drift [①](#)
- Infer drifted column types [①](#)
- Validate schema [①](#)

28. Here in the connection, we need to change the location for the CSV file. This time we will choose Customer02.csv to get the other half of the data.

CustomerCSV_Source

Connection [Edit](#) [New](#) [Learn more](#)

File path * csv / Customer / Customer02.csv [Browse](#) [Preview data](#)

Compression type Select...

Column delimiter Comma (,)

Row delimiter Default (\r\n or \r\n)

Encoding Default(UTF-8)

Quote character Double quote ("")

Escape character Backslash (\)

First row as header

Null value

29. Now to our existing dimension data flow we are going to add a new source.
 30. First give it a name then for the dataset we are going to choose our Dimension Customer Table which we created in the dataset for our customer sink.

Source settings Source options Projection Optimize Inspect Data preview

Output stream name * Learn more [?](#)

Description [Reset](#)

Source type * Dataset Inline

Dataset * **DimCustomerTable** [Test connection](#) [Open](#) [New](#)

Options Allow schema drift [?](#)
 Infer drifted column types [?](#)
 Validate schema [?](#)

Sampling * [?](#) Enable Disable

31. Then move to the source option choose the query option then paste the query as shown below.

This SQL query selects the maximum value from the column "CustomerSK" in the table "DimCustomer" from the database schema "dbo" (which typically stands for "database owner"). Here's a breakdown:

SELECT: This keyword is used to specify the columns that you want to retrieve data from.

MAX([CustomerSK]): This function calculates the maximum value of the column "CustomerSK". The MAX function is an aggregate function that returns the highest value in a set of values.

as CustomerSK: This part of the query renames the result of the MAX([CustomerSK]) calculation to "CustomerSK" in the output.

FROM [dbo].[DimCustomer]: This specifies the table from which the data is being retrieved. [dbo].[DimCustomer] indicates the schema (dbo) and table (DimCustomer) in the database.

In summary, this query retrieves the highest value of the "CustomerSK" column from the "DimCustomer" table in the specified database schema. This type of query is often used to find the maximum or minimum value in a column for analytical or reporting purposes.

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

Input

Table Query Stored procedure

Query ⓘ

```
SELECT MAX([CustomerSK]) as CustomerSK FROM [dbo].[DimCustomer]
```

Import projection

Enable staging

Batch size

Incremental column

Isolation level ⓘ

32. Now from our get max customer Stream we need to click on the plus icon and choose sink.

33. Then in the Sink first you need to give it a name and the incoming stream is get max customer SK. After that, you need to choose the Cache option.

Sink Settings Errors Mapping Optimize Inspect Data preview

Output stream name * Learn more

Description Reset

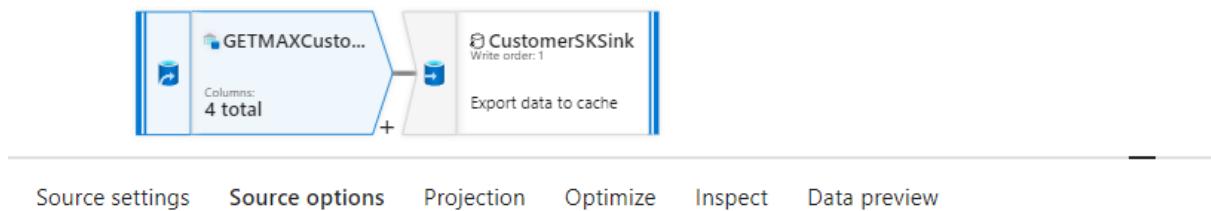
Incoming stream *

Sink type *

Dataset	Inline	Cache
---------	--------	-------

Options Write to activity output ⓘ

34. Then come back to get max customer stream and click on import projection.



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

Input

Table Query Stored procedure

Query ⓘ

```
SELECT MAX([CustomerSK]) as CustomerSK FROM [dbo].[DimCustomer]
```

[Import projection](#)

Enable staging

Batch size ⓘ

Incremental column ⓘ

Isolation level ⓘ

35. It is going to turn on data flow debug.

Turn on data flow debug

Integration runtime

AutoResolveIntegrationRuntime

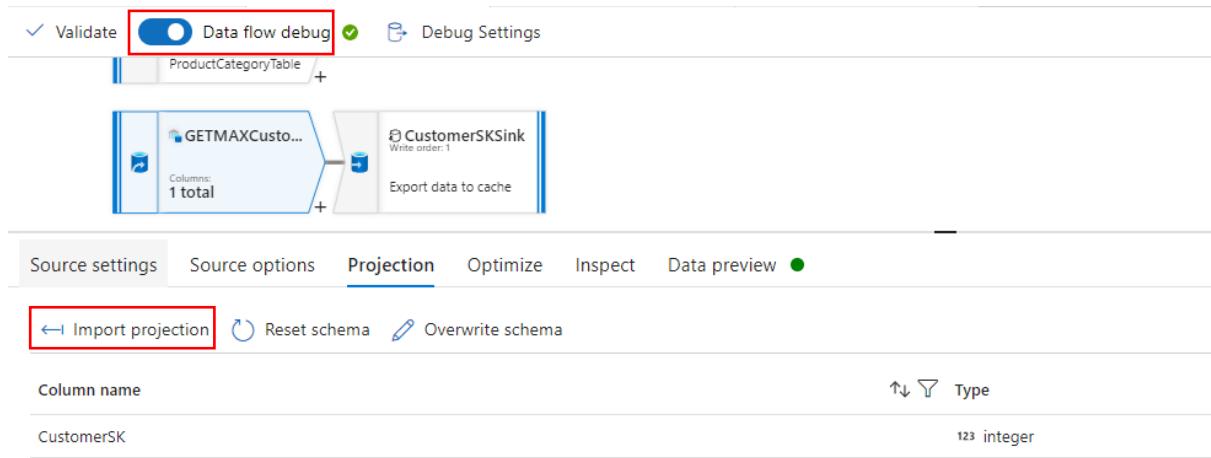
AutoResolveIntegrationRuntime

Region	AutoResolve
Compute size	Small

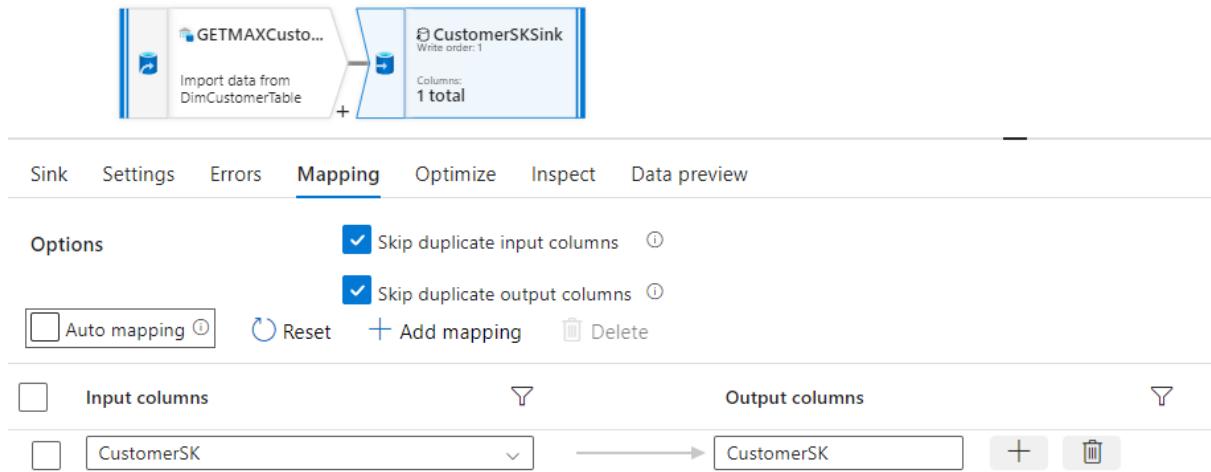
Debug time to live ⓘ

1 hour

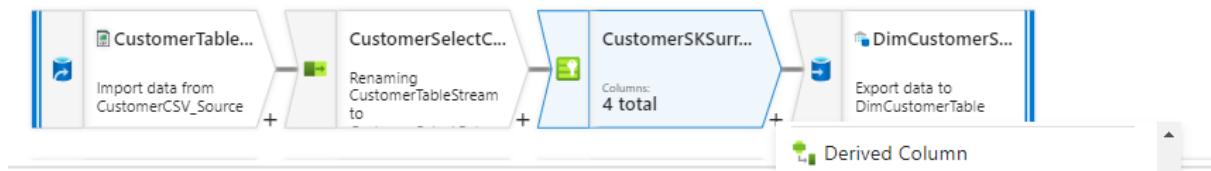
36. After that you need to come to projections and click on import projections and this will give you the below projection. Then you just need to turn off the data flow debug option.



37. Then go back to sink and go to mappings then turn off the auto-mapping feature. Then you will see that latest mapping in place.



38. Then from your dataflow diagram you need select customer SK surrogate key and click on the plus icon from it. Then choose derived column.



39. First give a name to your derived column then in the columns option choose Customer SK as your column and for the Expression you need to open the expression builder.

Derived column's settings

Output stream name *: CustomerSKDerivedColumn

Description: Creating/updating the columns 'CustomerID, CompanyName, SalesPerson, CustomerSK'

Incoming stream *: CustomerSKSurrogateKeyStream

Columns *:

Column	Expression
CustomerSK	Enter expression... ANY

[Open expression builder](#)

40. Here you have to write this expression as shown below and click on save and finish.

Dataflow expression builder

CustomerSKDerivedColumn

Derived Columns

CustomerSK

Column name *: CustomerSK

Expression: CustomerSK + CustomerSKSink#outputs()[1].CustomerSK

Save

Expression elements

All
Functions
abc CustomerID

Data preview

Save and finish

41. Once you are done with all this click on publish and then trigger your pipeline again.

42. Below you can see that our pipeline has ran successfully.

Run Succeeded

Successfully ran 05_mapping_dataflow_table (Pipeline).

[View pipeline run](#)

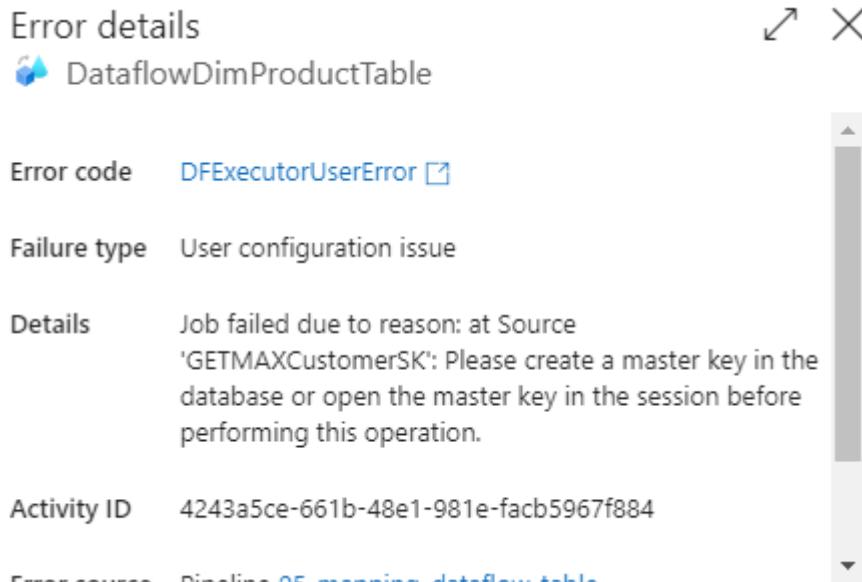
14 minutes ago

All pipeline runs > 05_mapping_dataflow_table - Activity runs

Rerun Cancel Refresh Update pipeline List Gantt

Activity runs									
Pipeline run ID: 32a11373-9989-4d92-9b92-04194958b42a									
All status									
Showing 1 - 1 items									
Activity name	Activity status	Activity type	Run start	Duration	Integration runtime	User properties	Activity run ID	Log	
DataflowDimProductTable	Succeeded	Data flow	5/2/2024, 6:43:43 PM	3m 20s	AutoResolveIntegrationR...		58fdce3e-6729-48ef-9530-f850bd87d9f8		

43. Now if in case you are getting an error like this below then you just have to create a password encryption key in you SSMS where you were creating tables all along.



44. You have to use the below statement to create an encryption key in your dedicated pooldb.

CREATE MASTER KEY ENCRYPTION BY PASSWORD = 'Password@1234';

45. Now you have to run the select command for the Customer table. This time you'll be able to see the data below customer id 468.

```
SQLQuery1.sql - dat...sqladminuser (0)* - 150% - [ ]
```

```
select * from DimCustomer
order by [CustomerSK]
```

CustomerSK	CustomerID	CompanyName	SalesPerson
468	268	Channel Outlet	adventure-works\jeo0
469	271	Alternative Vehicles	adventure-works\pamela0
470	272	Metro Cycle Shop	adventure-works\david8
471	273	A Typical Bike Shop	adventure-works\jilliano0
472	274	Active Systems	adventure-works\jilliano0
473	275	Outdoor Sporting Goods	adventure-works\shuhu0
474	276	Thrifty Parts and Sales	adventure-works\linda3
475	277	The Bicycle Accessories C...	adventure-works\shuhu0
476	280	Standard Bikes	adventure-works\jose1
477	281	Petroleum Products Distrib...	adventure-works\jose1
478	282	Quick Parts and Service	adventure-works\garnett1
479	286	Number One Bike Co.	adventure-works\jeo0
480	287	Handy Bike Services	adventure-works\michael9
481	288	Kidstand Sellers	adventure-works\michael9
482	289	Moderately-Priced Bikes St...	adventure-works\pamela0
483	290	Friendly Bike Shop	adventure-works\david8
484	291	Grand Bicycle Stores	adventure-works\jilliano0
485	292	Grease and Oil Products C...	adventure-works\jilliano0
486	293	Fashionable Bikes and Ac...	adventure-works\shuhu0
487	294	Flawless Bike Shop	adventure-works\linda3
488	295	Engineered Bike Systems	adventure-works\shuhu0
489	298	City Cycling	adventure-works\jose1
490	299	Citywide Service and Repair	adventure-works\jose1
491	300	Nice Bikes	adventure-works\garnett1
492	304	Essential Bike Works	adventure-works\jeo0
493	306	Work and Play Association	adventure-works\michael9
494	307	Riding Associates	adventure-works\pamela0
495	308	Rugged Bikes	adventure-works\david8
496	309	The Gear Store	adventure-works\jilliano0
497	310	Orange Bicycle Company	adventure-works\jilliano0
498	311	Principal Bike Company	adventure-works\shuhu0
499	312	Rezale Services	adventure-works\linda3
500	313	Metropolitan Manufacturing	adventure-works\shuhu0
501	316	Low Price Cycles	adventure-works\jose1
502	317	Finer Riding Supplies	adventure-works\jose1
503	318	First-Rate Outfit	adventure-works\garnett1
504	322	Tachometers and Accessor...	adventure-works\jeo0
505	323	Metro Bike Works	adventure-works\michael9
506	325	All Cycle Shop	adventure-works\pamela0
507	326	Year-Round Sports	adventure-works\david8
508	327	World of Bikes	adventure-works\jilliano0

Query executed successfully.

dataynapse1234.sql.azuresy... | sqladminuser (0) | demodb | 00:00:00 | 598 rows

