

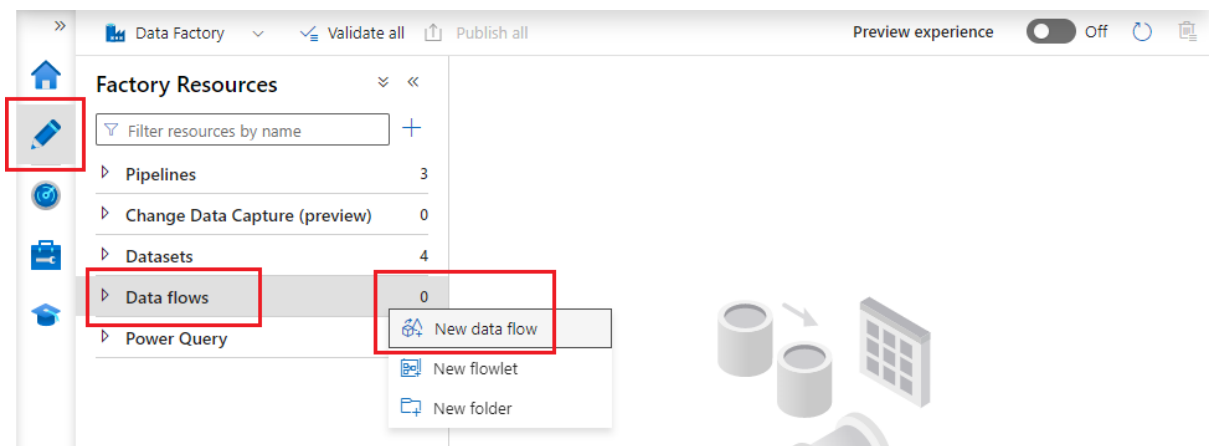
## Lab 4 – Author a data flow

In this lab you will use an Azure Data Factory **data flow** to implement a familiar data warehousing process: maintaining a dimension table.

### Lab 4.1 – Enable data flow debugging

Data flows are debugged using on-demand Apache Spark clusters. Provisioning a cluster takes several minutes, so start this lab by switching “Data flow debug” on for your ADF Studio session.

1. Navigate to ADF Studio’s Authoring experience, then use the “Data flows Actions” menu to create a new data flow.



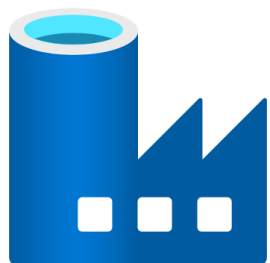
2. The data flow authoring canvas appears, with the usual “Properties” flyout on the right. Name the data flow “UpdateProductDimension” and dismiss the flyout.
3. Above the data flow canvas a “Data flow debug” toggle switch is visible. When you the move the switch to the right, the “Turn on data flow debug” flyout appears – click “OK” to accept the default options and to start provisioning the debug cluster.

While the debug cluster is warming up, continue with Lab 4.2.

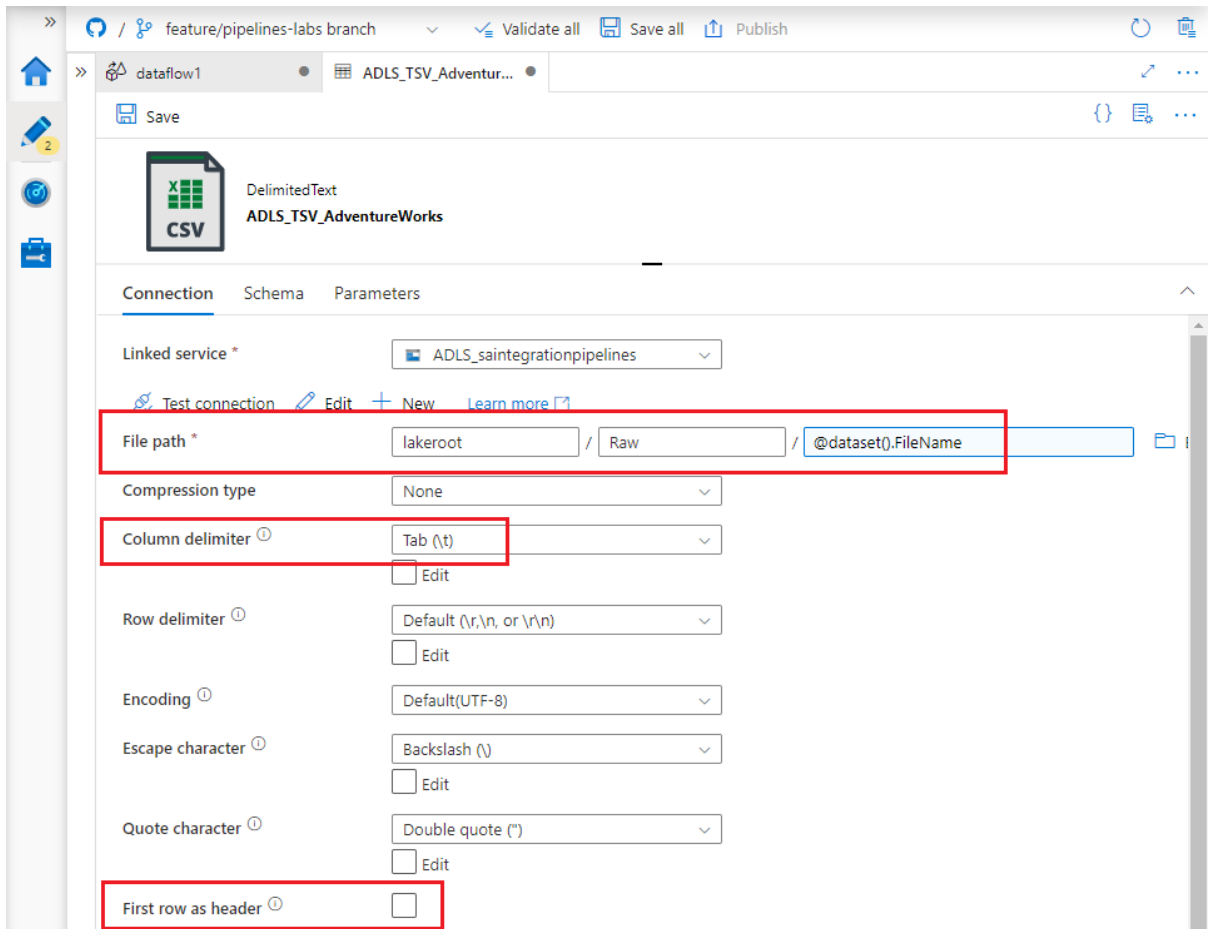
### Lab 4.2 – Create a delimited data lake dataset

The product dimension we will be building later in the lab will combine information from the Product, ProductSubcategory and ProductCategory files into a three-tier hierarchy. Reading data from a file in the data lake requires a dataset, but the binary dataset we have been using is not suitable in this case – ADF must split each file row into columns to make data available as a stream of records.

1. Create the new dataset as follows:
  - Choose the “Azure Data Lake Storage Gen2” data store.
  - Choose the “DelimitedText” file format. (The dataset must be created from scratch rather than cloned, because file format is not editable in a dataset clone).
  - Name the dataset “ADLS\_TSV\_AdventureWorks”, and select your data lake linked service.
  - Leave the other options with their default values – we need to define a dataset parameter before we can properly complete them. Click OK.



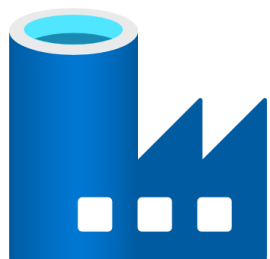
2. Define a dataset parameter called “FileName”.
3. Configure options on the dataset’s “Connection” tab:
  - Set the **File path** components to use the “lakeroot” file system, the “Raw” directory and the filename specified by the dataset parameter created in step 2.
  - Set **Column delimiter** to “Tab (\t)” – recall from Lab 2.4 that the files downloaded from GitHub are tab-separated, despite their “.csv” extensions.
  - Ensure that **First row as header** is unchecked – the files have no header row.



4. Verify that the dataset has been correctly configured by using the “Preview data” button to the right of the file path (out of view in the screenshot). You will be prompted for a filename parameter value – remember that, using the expression given in the screenshot, this must **include** the extension e.g. “Product.tsv”.

## Lab 4.3 – Combine data lake datasets

The product dimension combines product, subcategory and category information to support different aggregations of facts that have a product attribute. We are seeking to create the same effect as joining the tables together if they were in SQL database:

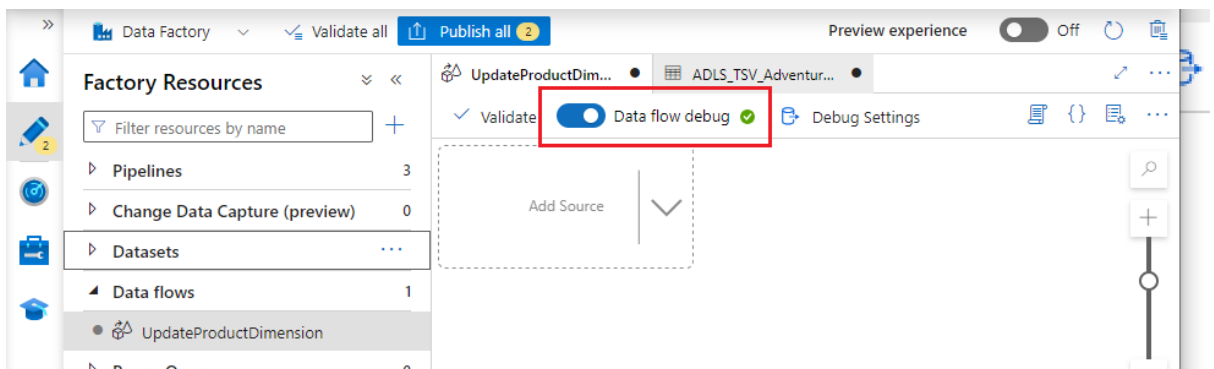


```

SELECT
  p.ProductID
, p.[Name] AS Product
, sc.[Name] AS SubCategory
, c.[Name] AS Category
FROM dbo.Product p
  LEFT JOIN dbo.ProductSubcategory sc
    ON sc.ProductSubcategoryId = p.ProductSubcategory
  LEFT JOIN dbo.ProductCategory c ON c.ProductCategoryId = sc.ProductCategoryId
    
```

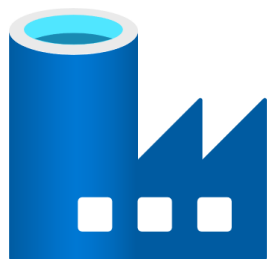
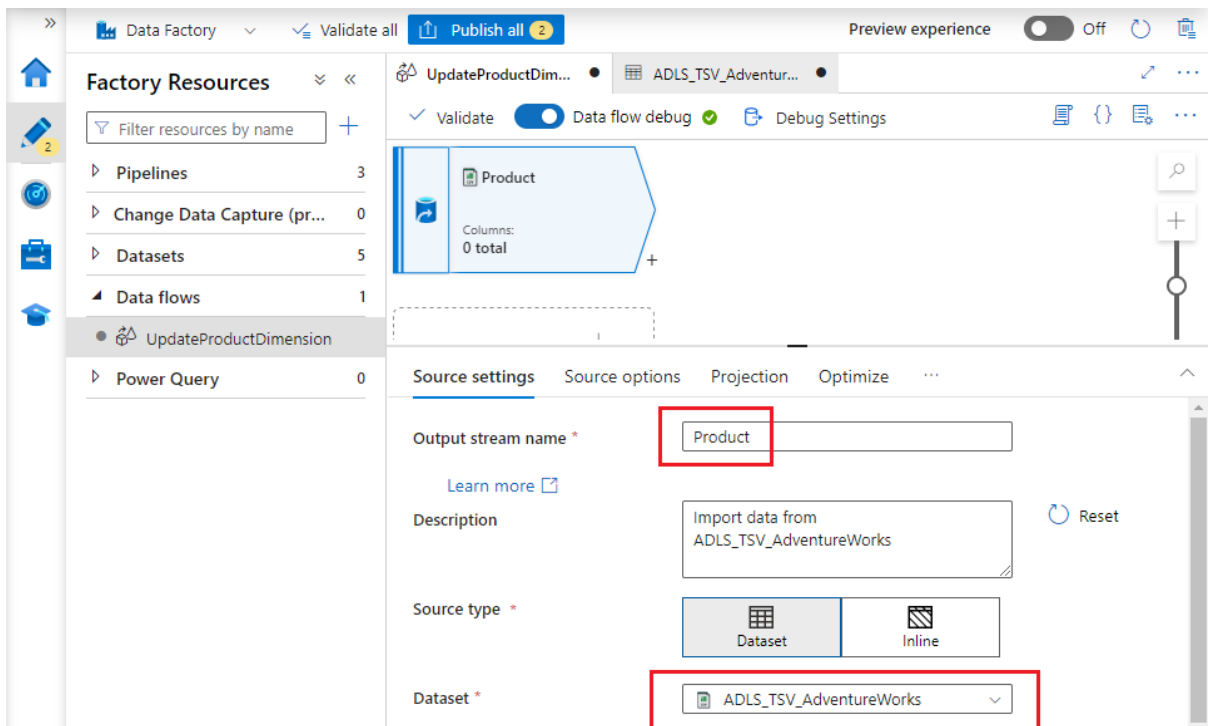
In this lab you will use a data flow to produce this result in Azure Data Factory.

1. Return to the data flow you created in Lab 4.1 and check that the debug cluster has been successfully provisioned. When the cluster is available, a tick mark in a green circle appears to the right of the “Data flow debug” slider.

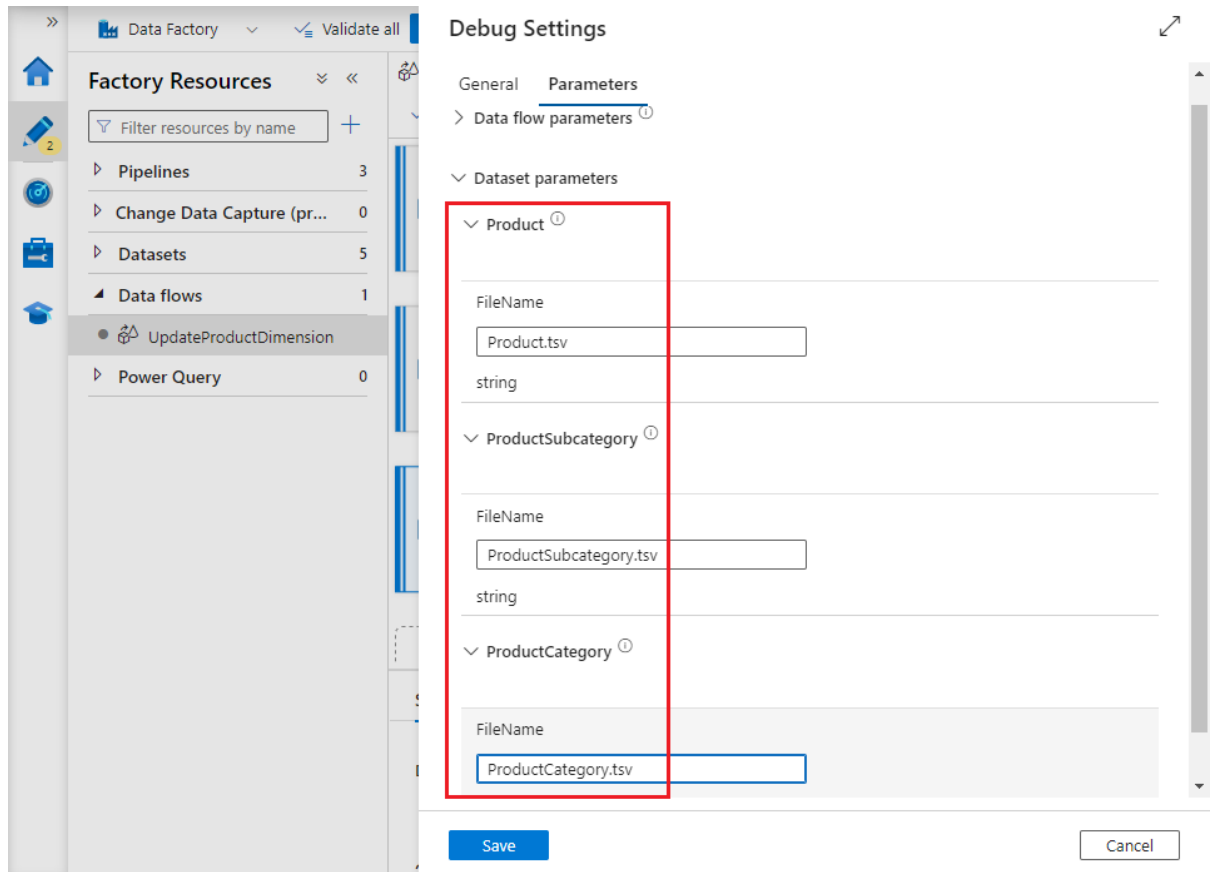


If the cluster is not ready yet, wait for it to finish warming up.

2. Click the “Add source” tile on the data flow canvas. On the source transformation’s **Source settings** tab, change its “Output stream name” to “Product” and select the “ADLS\_TSV\_AdventureWorks” dataset.

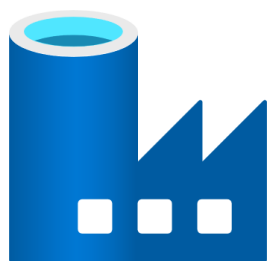


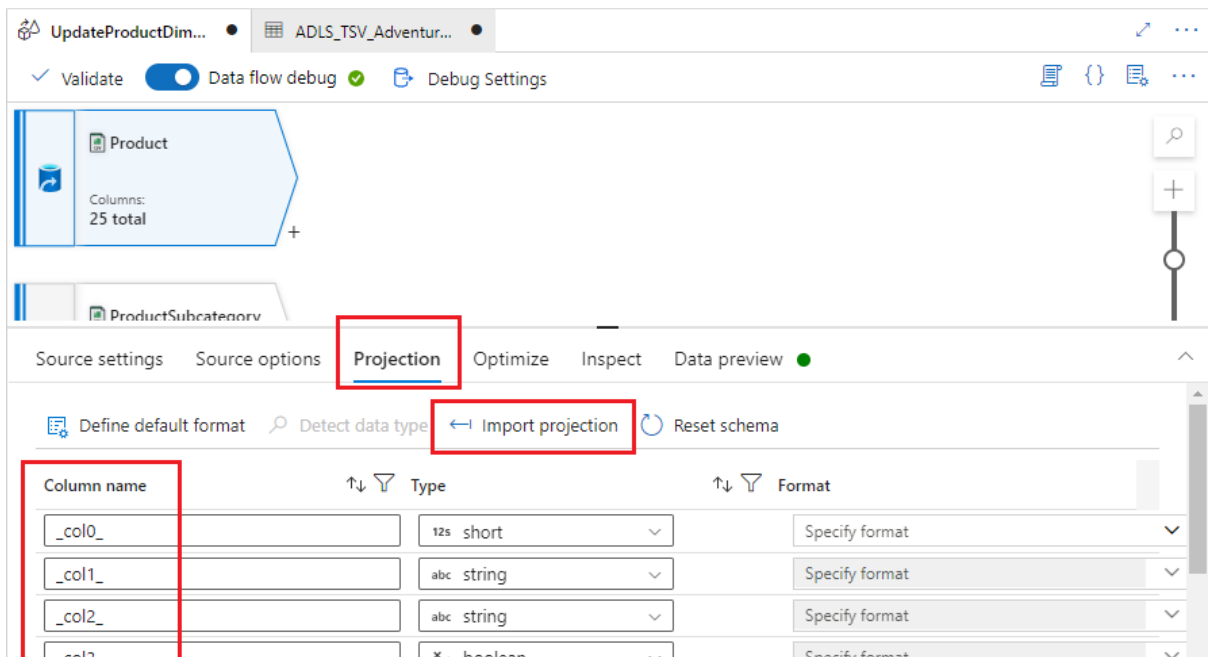
- Repeat step 2 to add two more source transformations below “Product”. Name one “ProductSubcategory” and the other “ProductCategory” – both should also use the “ADLS\_TSV\_AdventureWorks” dataset.
- At execution time, source dataset parameter values are specified from **outside** a data flow. To configure them during debug, click “Debug settings” in the data flow canvas header bar and choose the “Parameters” tab. Each dataset parameter appears in the “Dataset parameters” section – set each one appropriately and click “Save”.



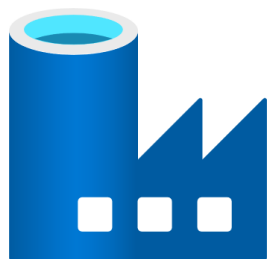
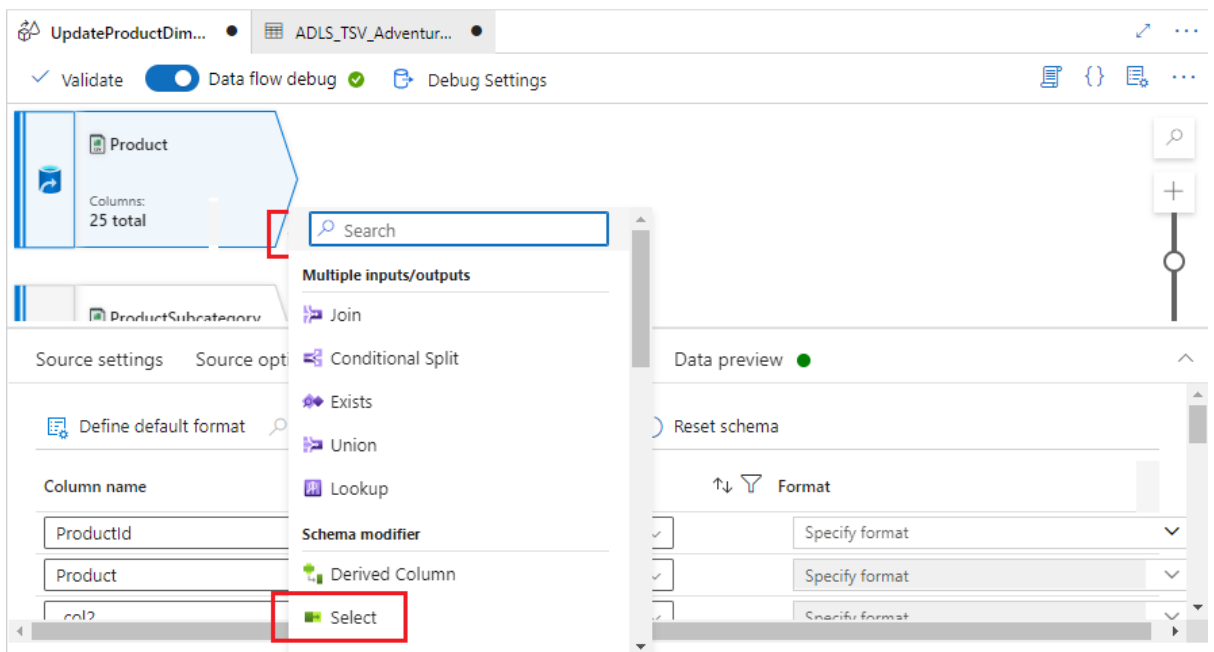
The screenshot shows the 'Debug Settings' dialog box in Azure Data Factory. The 'Parameters' tab is selected, and the 'Dataset parameters' section is expanded. Three parameters are listed: 'Product', 'ProductSubcategory', and 'ProductCategory'. Each parameter has a 'FileName' field and a 'string' type. The 'Product' parameter has a value of 'Product.tsv', 'ProductSubcategory' has 'ProductSubcategory.tsv', and 'ProductCategory' has 'ProductCategory.tsv'. A red box highlights the 'Product' and 'ProductSubcategory' parameters, and a blue box highlights the 'ProductCategory' parameter. The 'Save' button is visible at the bottom left, and the 'Cancel' button is at the bottom right.

- Return to the “Product” source transformation, select its “Projection” configuration tab and click “Import projection”. After a few moments, a schema will be inferred from the source file and presented in the tab.





- Imported column names are autogenerated and are not meaningful, because the source file contains no column headers. The projection tab is editable, so you can change column names and types as required. In the case of Product
  - Rename “\_col0\_” to “ProductId” and ensure its type is “integer”
  - Rename “\_col1\_” to “Product” and ensure its type is “string”
  - Rename “\_col18\_” to “SubcategoryId” and ensure its type is “integer”
- Use the “+” button at the bottom right of the “Product” transformation to append a “Select” transformation. The transformation allows you to remove and rename columns, without changing their types – remove all of the “\_col...” named columns, leaving only the three columns you renamed in step 6.



If you don't see a list of columns on the Select transformation's "Select settings" tab, make sure that the schema is present in the Source transformation and that the "Auto mapping" checkbox is **not** ticked.

- Repeat steps 5, 6 & 7 for the "ProductSubcategory" and "ProductCategory" source transformations. The names and types of columns to be retained for product subcategory are:

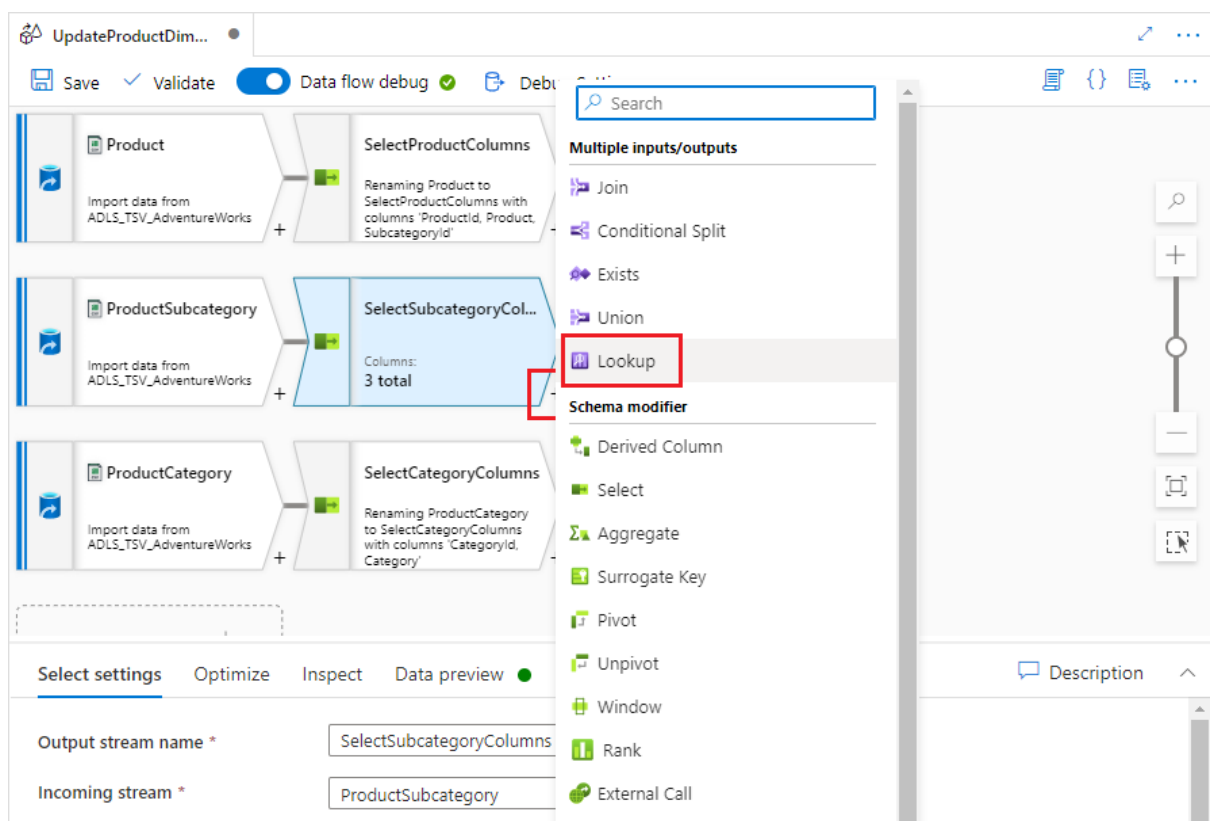
- Rename "\_col0\_" to "SubcategoryId" and ensure its type is "integer"
- Rename "\_col1\_" to "CategoryId" and ensure its type is "integer"
- Rename "\_col2\_" to "Subcategory" and ensure its type is "string"

The names and types of columns to be retained for product category are:

- Rename "\_col0\_" to "CategoryId" and ensure its type is "integer"
- Rename "\_col1\_" to "Category" and ensure its type is "string"

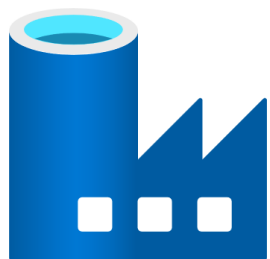
You now have three parallel streams, loading and modifying data from the three source files.

- You can combine data into the Product stream from the other two streams using the data flow "Lookup" transformation. First, click the small "+" button below the Product Subcategory stream's Select transformation, then select "Lookup" from the popup menu.



The screenshot shows the Azure Data Factory interface for a data flow named "UpdateProductDim...". It features three source streams: "Product", "ProductSubcategory", and "ProductCategory", each importing data from "ADLS\_TSV\_AdventureWorks". Each source stream is followed by a "Select" transformation. The "ProductSubcategory" stream's "SelectSubcategoryColumns" transformation is highlighted with a red box, and a small "+" button below it is also highlighted. A popup menu is open, showing various transformation options, with "Lookup" selected and highlighted by a red box. The "Lookup" option is under the "Schema modifier" section. The "Select settings" tab is active, showing the "Output stream name" as "SelectSubcategoryColumns" and the "Incoming stream" as "ProductSubcategory".

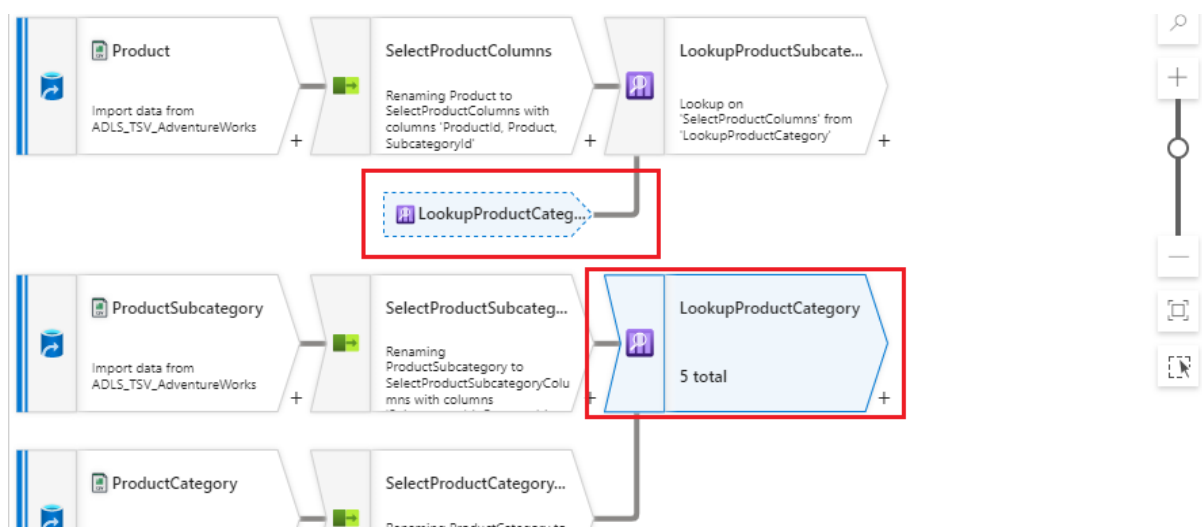
- On the **Lookup settings** tab, name the transformation then set "Lookup stream" to use product category columns. (You can choose from any of the other five previous transformations, so take care to pick the ProductCategory stream's Select transformation, and not the earlier Source for the stream).



“Lookup conditions” specifies the fields to compare from each transformation and the operator to compare them with – choose the streams’ respective CategoryId fields. Notice that the lookup relationship is also displayed on the data flow canvas.

The screenshot shows the 'Lookup settings' panel for a 'LookupProductCategory' activity. The 'Lookup stream' is set to 'SelectProductCategoryColumns'. The 'Lookup conditions' are configured to match 'CategoryId' from 'SelectProductSubcategoryColumns' with 'CategoryId' from 'SelectProductCategoryColumns'.

11. Add a second “Lookup” activity, this time on the Product stream, performing a lookup against the ProductCategory lookup’s output, based on matching SubcategoryId. This time the canvas displays a “reference node” instead of showing a direct link between the two transformations.



This is just for readability – when you hover over either the reference node or the transformation it refers to, both light up in blue to indicate that they mean the same thing.

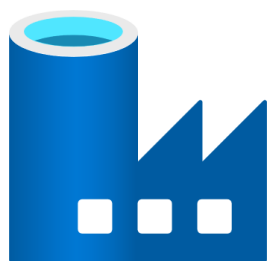
- Open the new Lookup transformation’s “Inspect” tab to view the set of columns present in the combined stream – notice it includes two copies of each of the join fields, one from each stream participating in the lookup. Add another Select transformation to clean this up, deleting one of each pair.

The screenshot shows the Azure Data Factory interface. At the top, there's a pipeline diagram with nodes: 'Product' (Import data from ADLS\_TSV\_AdventureWorks), 'SelectProductColumns' (Renaming Product to SelectProductColumns with columns 'Productid, Product, Subcategoryid'), 'LookupProductSubcate...' (Lookup on 'SelectProductColumns' from 'LookupProductCategory'), and 'RemoveDuplicateColu...' (Columns: 6 total). A 'LookupProductCateg...' node is also visible below the main flow.

The 'Inspect' tab is active, showing a table of mappings. A red box highlights the 'Delete' button in the toolbar. Another red box highlights the first two rows of the table, which are selected.

	LookupProductSubcategory's column	Name as		
<input type="checkbox"/>	123 Productid	Productid	+	🗑️
<input type="checkbox"/>	abc Product	Product	+	🗑️
<input type="checkbox"/>	123 SelectProductColumns@SubcategoryId	SubcategoryId	+	🗑️
<input checked="" type="checkbox"/>	123 SelectProductSubcategoryColumns@Subcat.v	SubcategoryId	+	🗑️
<input type="checkbox"/>	125 SelectProductSubcategoryColumns@Catego.x	CategoryId	+	🗑️
<input type="checkbox"/>	123 Subcategory	Subcategory	+	🗑️
<input checked="" type="checkbox"/>	123 SelectProductCategoryColumns@CategoryId	CategoryId	+	🗑️
<input type="checkbox"/>	abc Category	Category	+	🗑️

- Finally, write the transformed dimension data back to the data lake using a “Sink” transformation. Add the transformation in the usual way, using the small “+” button following the Select transformation that removes duplicate columns. Sink is at the bottom of the list on the popup menu.





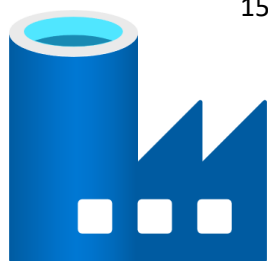
The screenshot shows the 'Settings' tab for a 'WriteToDataLake' sink. The 'Sink type' is set to 'Inline' (highlighted with a red box). The 'Inline dataset type' is set to 'Delta' (also highlighted with a red box). The 'Linked service' is 'AzureDataLakeStorage1'. Other settings include 'Output stream name' as 'WriteToDataLake' and 'Description' as 'Write results to data lake'.

You haven't yet created a dataset to use as the data flow sink – we will use an inline dataset instead. Inline datasets aren't reusable by other data flows or pipelines, but avoid clutter when reusability is not required. Select "Sink type" = "Inline" and choose "Delta" from the "Inline dataset type" dropdown. Delta tables are built on Parquet files, a column-oriented, highly compressible file format, offering significant performance and other benefits for data lakes.

- Specify the sink "Folder path" on the transformation's "Settings" tab (Delta is a multi-file storage format, so folder path identifies the folder **containing** those files). Choose the "lakeroot" file system and enter "Prepared/DimProduct" as the folder path. Select the "Overwrite" radio button to replace previous versions at each execution of the pipeline.

The screenshot shows the 'Settings' tab for a 'WriteToDataLake' sink. The 'Folder path' is set to 'lakeroot / Prepared/DimProdu...' (highlighted with a red box). The 'Table action' is set to 'Overwrite' (also highlighted with a red box). Other settings include 'Compression type' as 'None' and 'Vacuum' as '0'.

- Click "Publish all" to publish your source dataset and data flow to the ADF service.



## Lab 4.4 – Run the data flow

Data flows are executed by ADF pipelines. To run your data flow, create a pipeline for it.

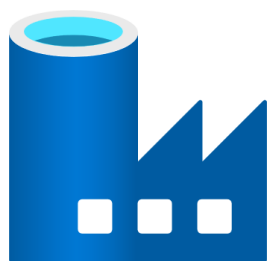
1. Create a new ADF pipeline.
2. Expand the “Move & transform” group in the activity toolbox, then drag a “Data flow” activity onto the pipeline canvas. Name the activity appropriately, then on its “Settings” tab:
  - choose the data flow you created in Lab 4.3
  - provide dataset parameter values (identifying files to be loaded by each source transformation).

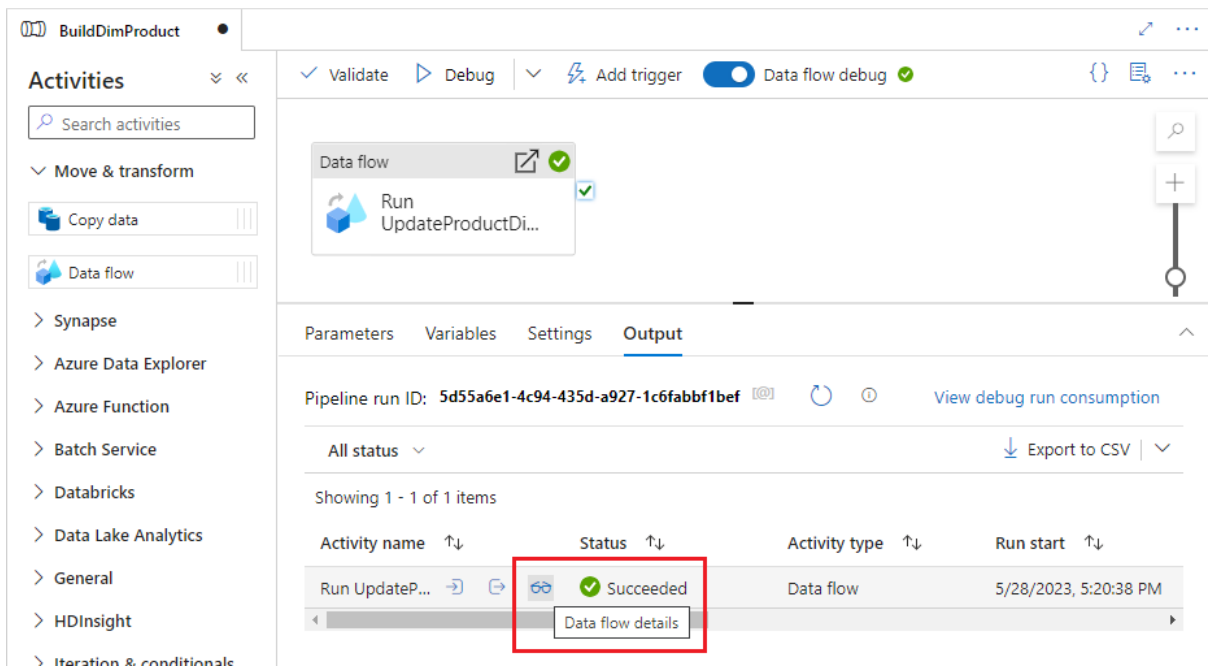
The screenshot shows the Azure Data Factory interface for a pipeline named 'BuildDimProduct'. The 'Activities' pane on the left shows the 'Data flow' activity selected. The main canvas displays a 'Data flow' activity named 'Run UpdateProductDi...'. The 'Settings' tab is active, showing the 'Data flow \*' dropdown set to 'UpdateProductDimension'. Below this, the 'Product parameters' section is expanded, displaying a table with the following content:

Name	Value	Type
FileName	Product.tsv	string

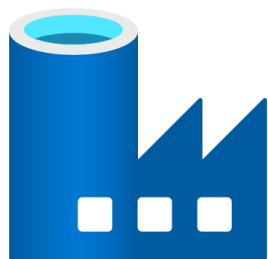
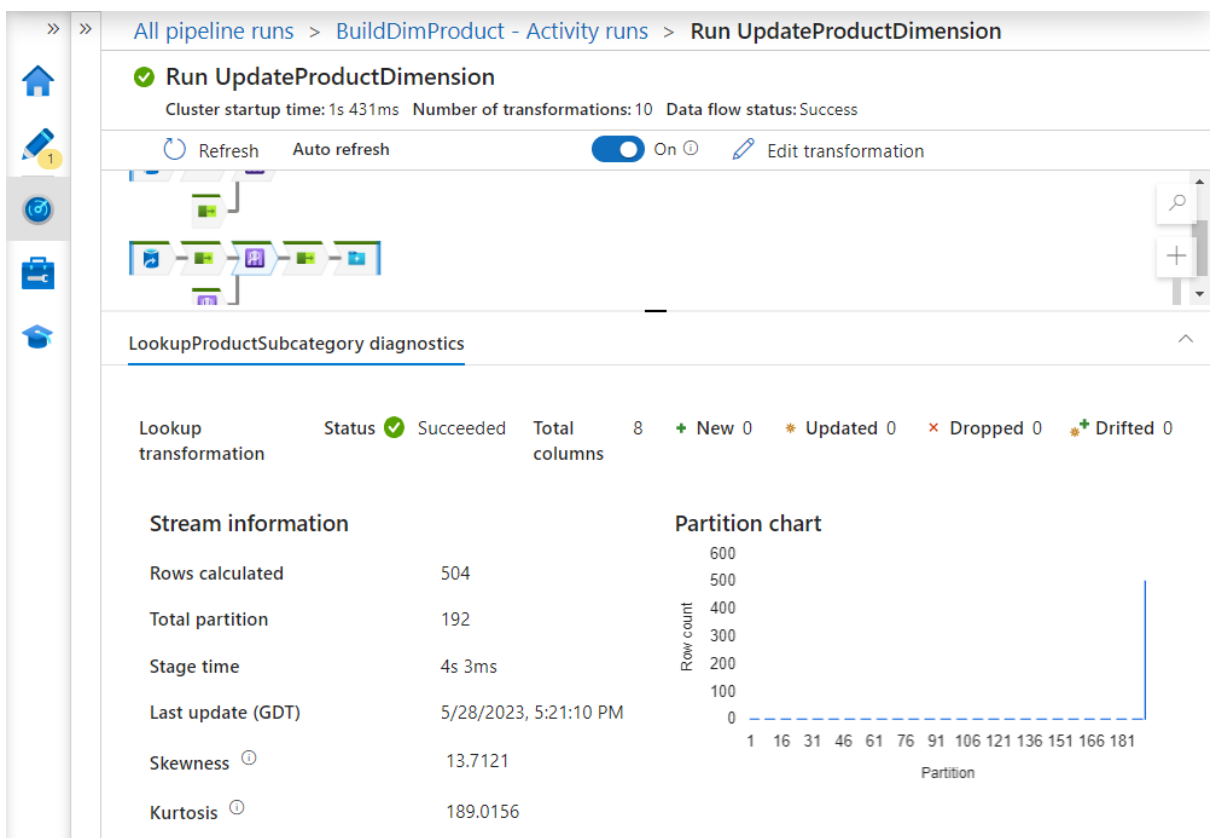
Below the table, there are expandable sections for 'ProductSubcategory parameters' and 'ProductCategory parameters'. The 'Run on (Azure IR)' dropdown is set to 'AutoResolveIntegrationRuntime'.

3. Click “Debug” to run the pipeline in debugging mode. A Spark cluster (Data flow debug enabled) is required to debug pipelines containing data flows, just as when you are developing them – if your debug session has timed out, you will need to enable a new one as in Lab 4.1.
4. When the pipeline has finished running, a row of data appears in its “Output” pane for the Data flow activity. Hover over the activity’s name to reveal the “Data flow details” button (“eyeglasses” icon).



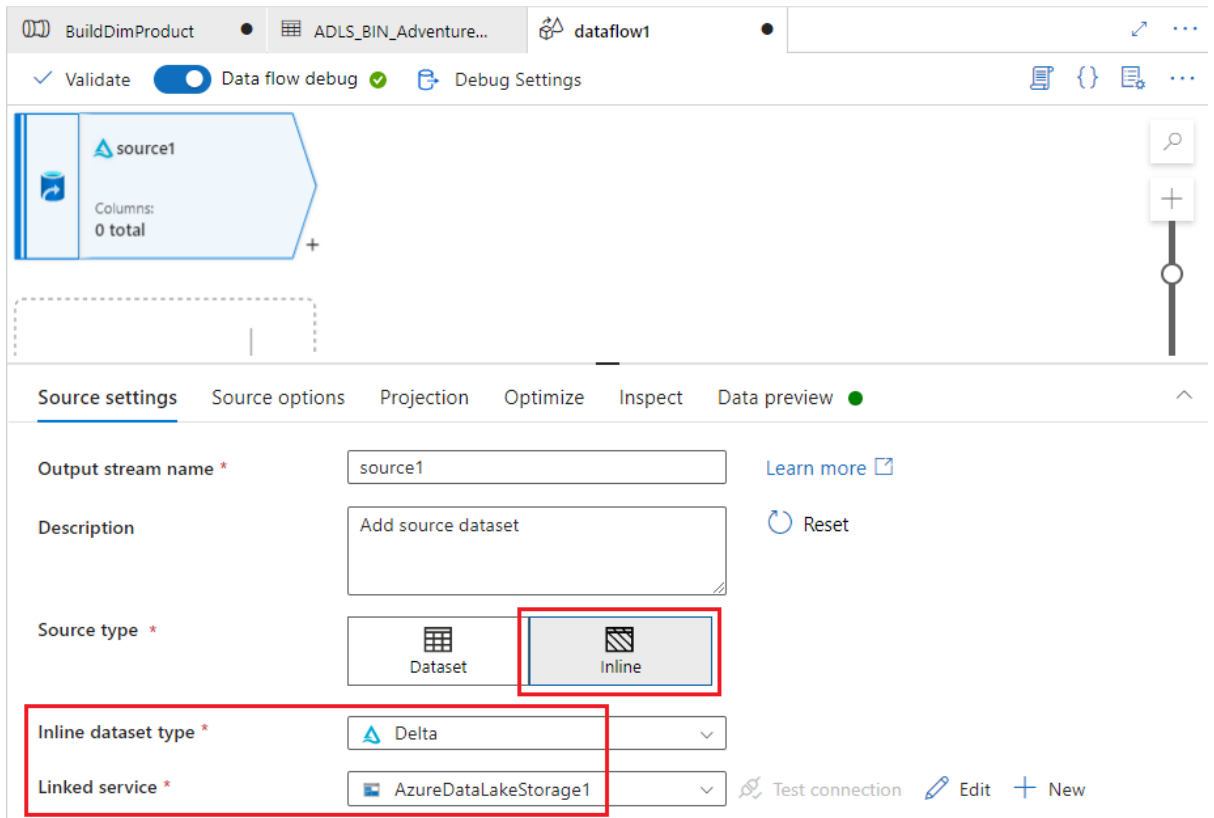


- Click the “Data flow details” button to open an interactive visualisation of more detailed data flow performance information. Selecting different transformations allows you to see the number of rows processed by a transformation, how quickly, and how the Spark cluster partitioned data for parallel processing. (For larger datasets, you can configure dataset distribution yourself to optimise Spark executor partitioning, via each transformation’s “Optimize” tab in the data flow editor).



6. Finally, you may wish to inspect the product dimension data written to your data lake. You can view the collection of Delta table files in the Prepared/DimProduct data lake folder, but you cannot read them directly. To inspect dimension contents:

- create a new data flow and add a Source transformation
- on the **Source settings** tab, specify an inline dataset of type “Delta” (shown below)
- on the **Source options** tab, browse to the lakeroor/Prepared/DimProduct folder
- on the **Data preview** tab, click “Refresh”.



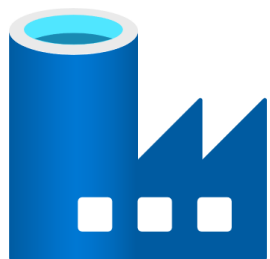
**Note:** When you preview data in this way, you may see very few subcategory and category values – this is because not every product in the source data is linked to a subcategory. To assure yourself that the lookups are correctly implemented, add a Sort transformation to your UpdateProductDimension data flow, ordering the output by Subcategory **without** nulls first.

## Lab 4.5 – Further work

This lab introduced concepts essential to the creation of a basic ADF data flow, but there is much more to learn. When using the popup menu to add Select, Lookup and Sink transformations you will have noticed that many more transformations exist. Data flows have their own expression language which supports powerful, complex data transformations. Flowlets allow you to create reusable data flow components which can be used by multiple data flows.

Start to broaden your knowledge by:

- using some of the other transformation types
- using the “Derived Column” transformation to begin to explore the data flow expression language.



## Recap

In Lab 4 you have:

- created an ADF data flow
- used Data flow debug to import file schemas (using “Import projection”)
- created a pipeline to execute your data flow
- used Data flow debug to run the pipeline in ADF Studio.

