# Exercise aim

The aim of this exercise is to construct a pipeline using Azure Data Factory that enables a CSV file to be uploaded to Azure Blob Storage, cleaned up and the columns remapped in Italian. It is essential that the process maintains only the columns 'Films', 'Genres' and 'Ratings', filters out films with ratings higher than 7 and ensures efficient and parallel data transfer. Furthermore, it is crucial to retain the metadata for any additional useful information.

# Exercise performance

1. **Setting up the Resource Group, Managed Identity and Key Vault**

   Once the subscription had been activated, a Resource Group was created by entering the desired subscription and location (Fig. 1). This was followed by the creation of a Managed Identity and Key Vault.
   A Managed Identity was created with the objective of facilitating the management of access to Azure services (Fig. 2). A Key Vault was configured, granting access to the Managed Identity via access policies and setting a key for secure access to resources (Fig. 3).
   The Managed Identity was assigned the role of 'Contributor' in the subscription, thereby ensuring access and management of all resources (Fig. 4).
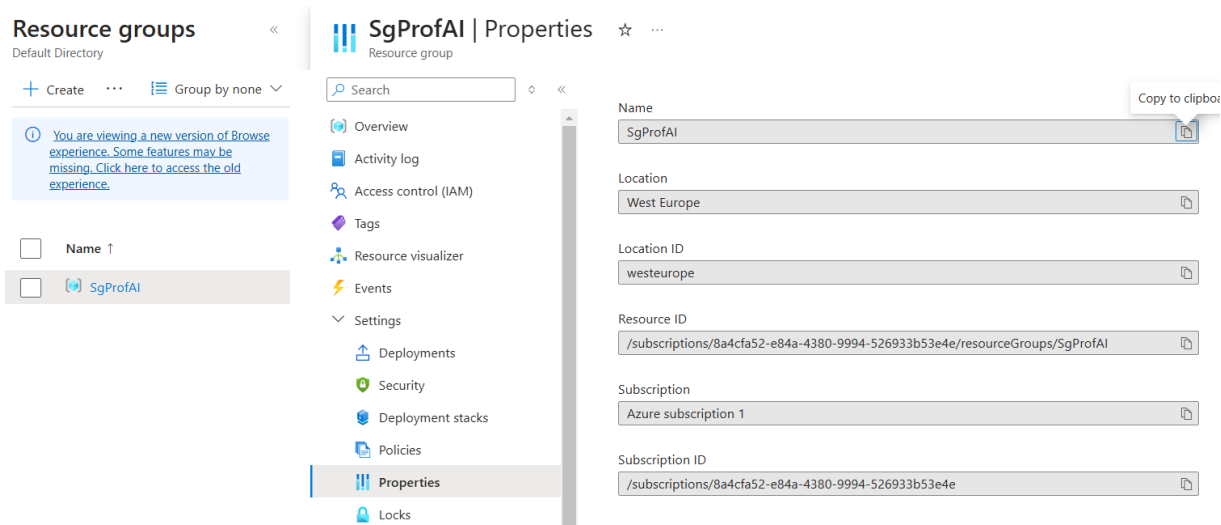

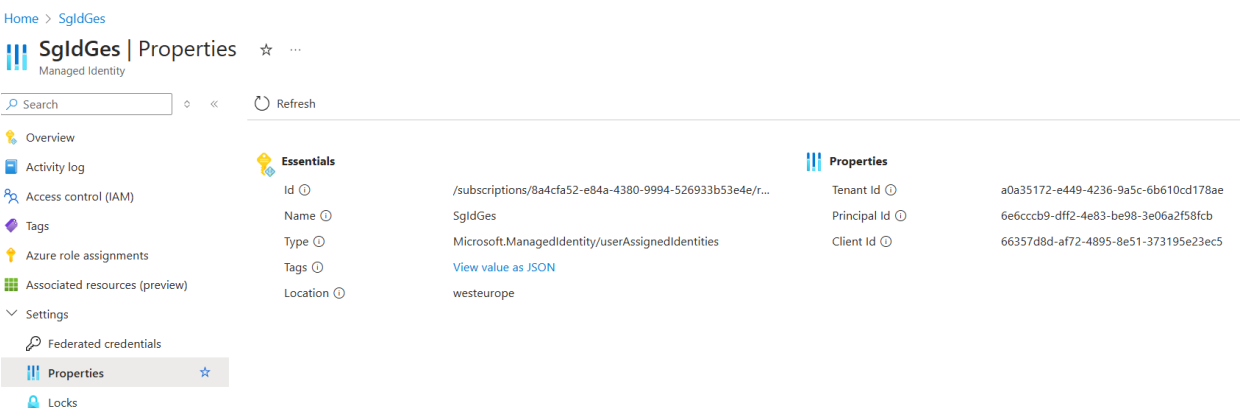
*Fig. 1 – Resource Group properties.*
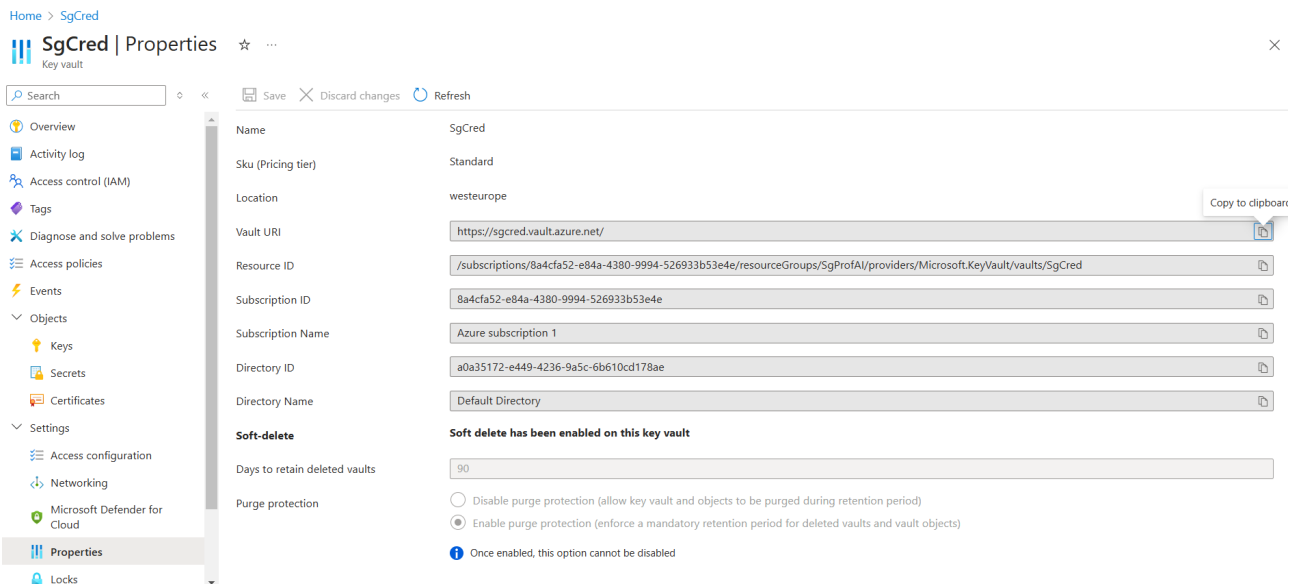


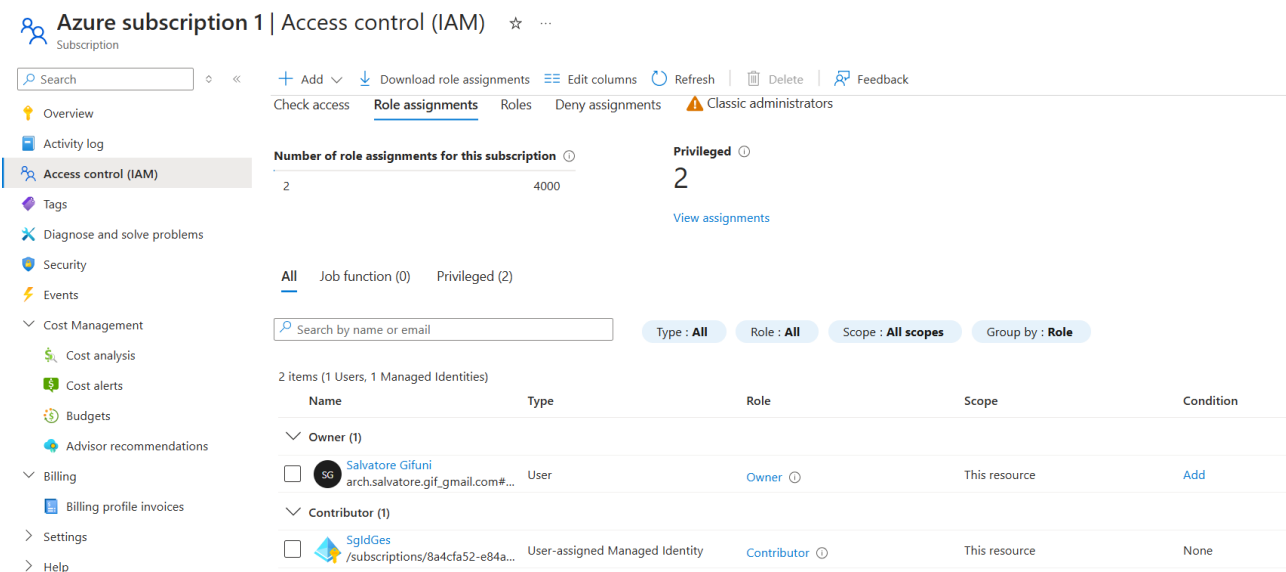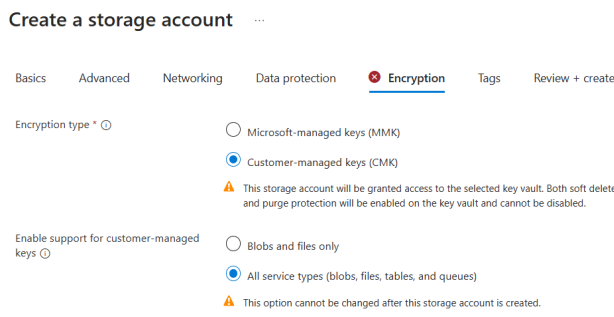*Fig. 2 - Managed Identity properties.*

Fig. 3 - Key Vault properties.



Fig. 4 - Setting the 'Contributor' role to the Managed Identity in the subscription.

2. **Azure Storage Account creation**

A storage account was created, paying particular attention to the security settings in the 'Encryption' section (Fig. 5).

Three containers were created: one for the input file, one for files with intermediate changes and one for the final file.

Fig. 5 - Encryption' settings for storage account creation.

### 3. Creating the Azure Translator

A translator was created, keeping the keys and URL for use in the pipeline web activity (Fig. 6).



Fig. 6 – Page containing keys and URLs to be passed to the pipeline web activity.

### 4. Creation of the Azure Data Factory

A Data Factory has been created. All previously configured security settings have been entered in the 'Advanced' section (Fig. 7).



Fig. 7 – Data Factory Encryption Settings

## 5. Creation of the pipeline

The pipeline created (Fig. 8) takes the original dataset as input via a data flow - 'DataWrangling' (Fig. 9) - and selects only the columns 'Movies'[1], 'Genres' and 'Ratings'.  Genres' and 'Ratings'[2] columns, filters the films, keeping only those with a rating higher than 7[3] (Fig. 11), and sorts them in descending order (optional). The resulting file is saved in an intermediate file for subsequent pipeline activities.
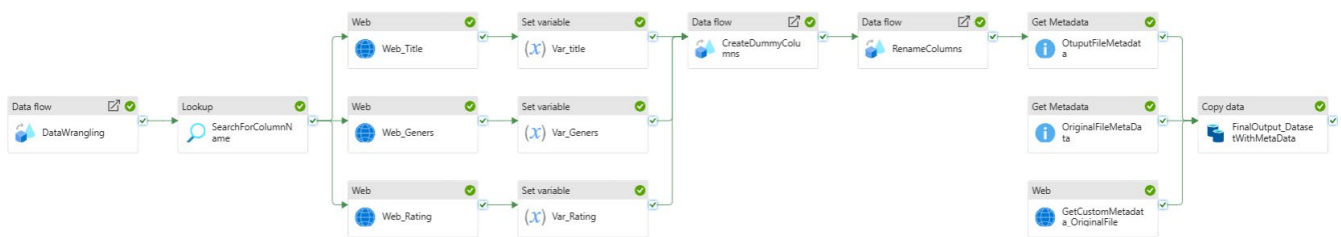


*Fig. 8 - Complete Pipeline.*



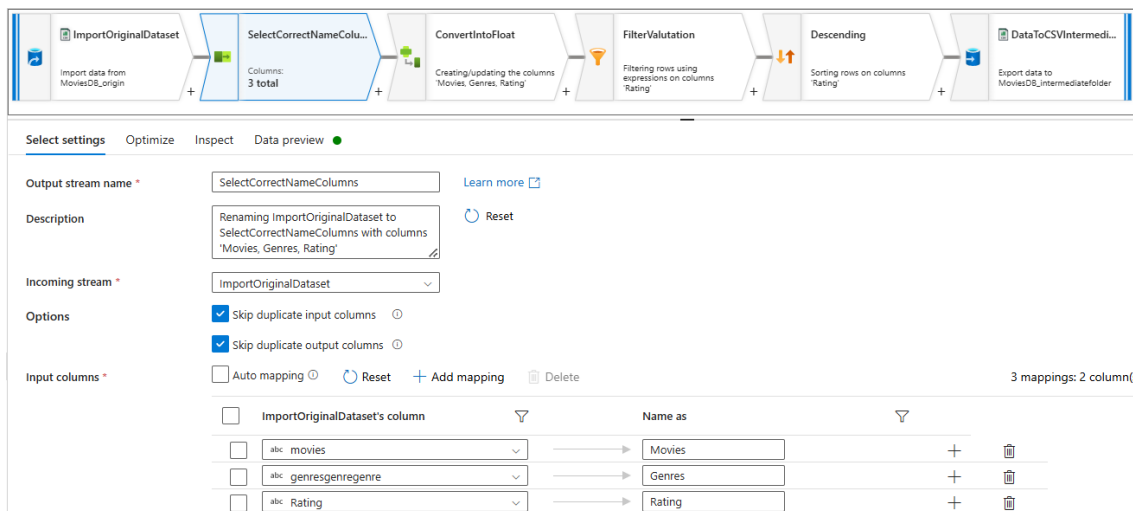*Fig. 9 - Data Flow 'DataWrangling'.*



*Fig.  10 – The node allows the selection of only the 3 required columns, instead of the original 5 columns, and is used to correct the format and errors in the names in order to facilitate subsequent machine translation activities.*

---

[1]  It was considered to keep the original 'Movies' column, but the process works equally well with the others. Also, using Azure Translator, you can see that the translation of 'Movies' from English to Italian is 'Cinema'. Subsequently, I replaced the translation from Italian to French, resulting in 'Movies', as required by the exercise.

[2] The column names could have been transformed directly in this Data Flow at the 'SelectCorrectColumnName' node, but we wanted to try to automate the translation of the column names with the subsequent use of Azure Translator.

[3] The values of the Rating column were transformed into floats to enable the subsequent filter to function (Fig. 11)
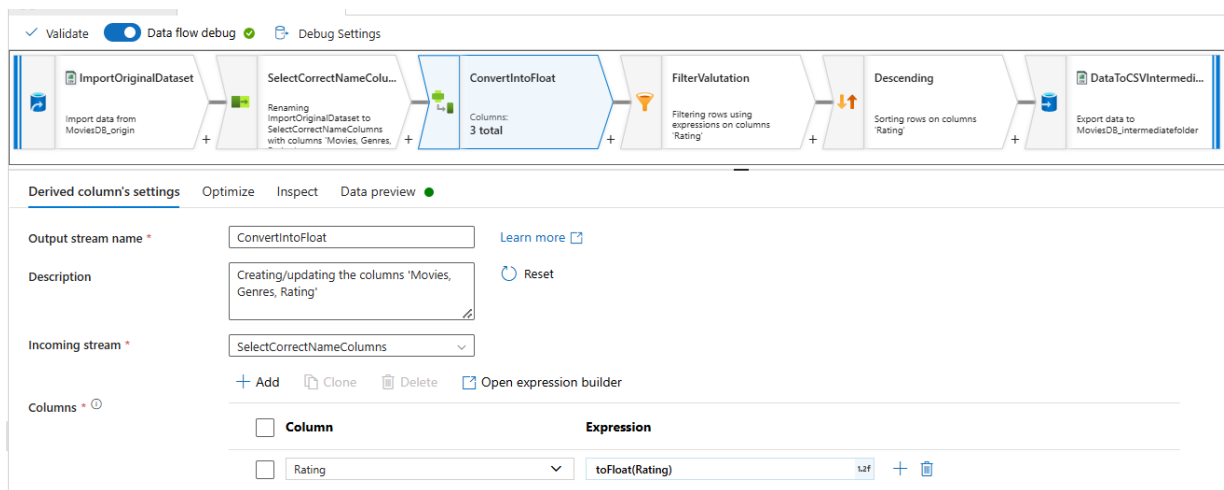
Fig. 11 – *Casting Column Values Rating.*

Next, a dataset was created from the intermediate file (Fig. 12) by disabling the 'first row as header' option to allow the Lookup – 'SearchForColumnName' activity to extract the column names. These were sent via a web activity (fig. 13) for automatic translation using Azure Translator (Fig. 14).
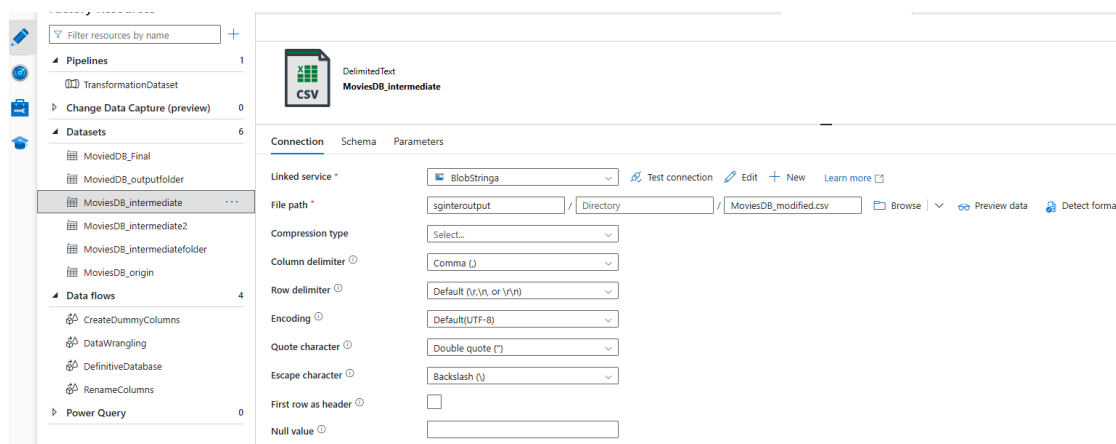


Fig. 12 – *Creating the dataset linked to the output file of the 'Data Wrangling' activity.*



Fig. 13 – *On the left the output of the Look up activity, on the right the body of the Web activity sent for translation via Azure Translator.*
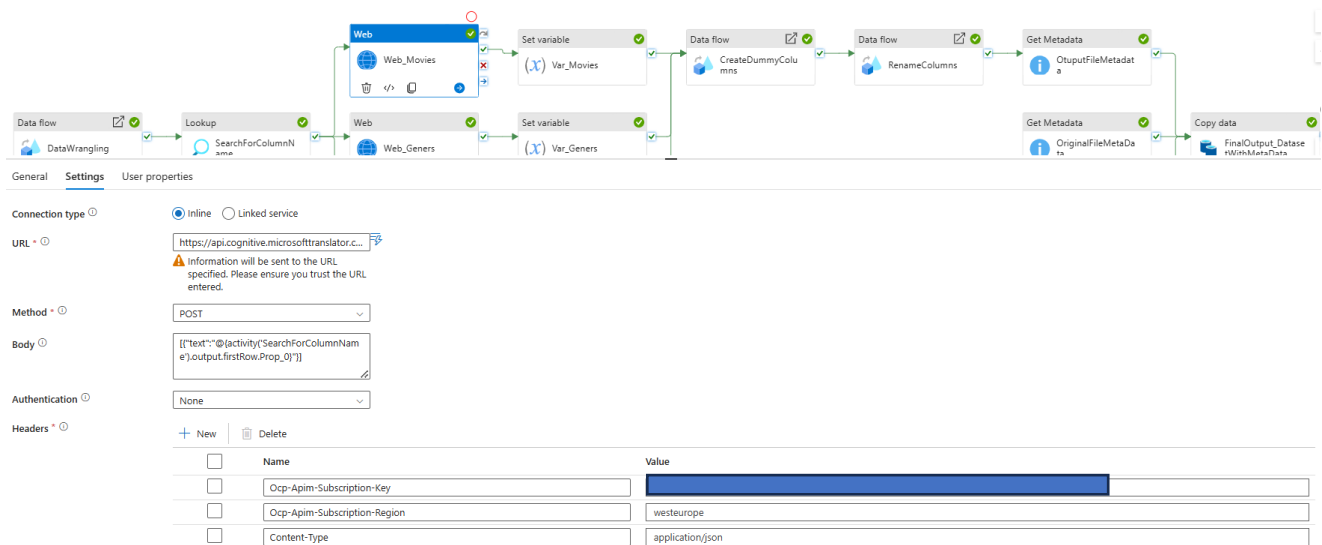
General | Settings | User properties

**Connection type** ⓘ  ◉ Inline  ○ Linked service

**URL** * ⓘ  https://api.cognitive.microsofttranslator.c...
⚠ Information will be sent to the URL specified. Please ensure you trust the URL entered.

**Method** * ⓘ  POST

**Body** ⓘ  [{"text":"@{activity('SearchForColumnName').output.firstRow.Prop_0}"}]

**Authentication** ⓘ  None

**Headers** * ⓘ  + New | 🗑 Delete

| | Name | Value |
|---|---|---|
| ☐ | Ocp-Apim-Subscription-Key | |
| ☐ | Ocp-Apim-Subscription-Region | westeurope |
| ☐ | Content-Type | application/json |

Fig. 14 – Web Activity Settings. In the URL, in addition to the 'Text Translation' link in figure 6, 'translate?api version=3.0&from=en&to=en' must also be entered to specify the version of the API you are using and that you want the text to be translated from English into Italian. In the first header, the 'key1' in figure 6 should be entered.

The outputs of the web activities were stored as pipeline variables, which were crucial for the subsequent activities; in fact, parameters were created in the Data Flow – 'RenameColumns' (Fig. 15 and Fig. 16) to link the pipeline variables to the column names (Fig. 17).
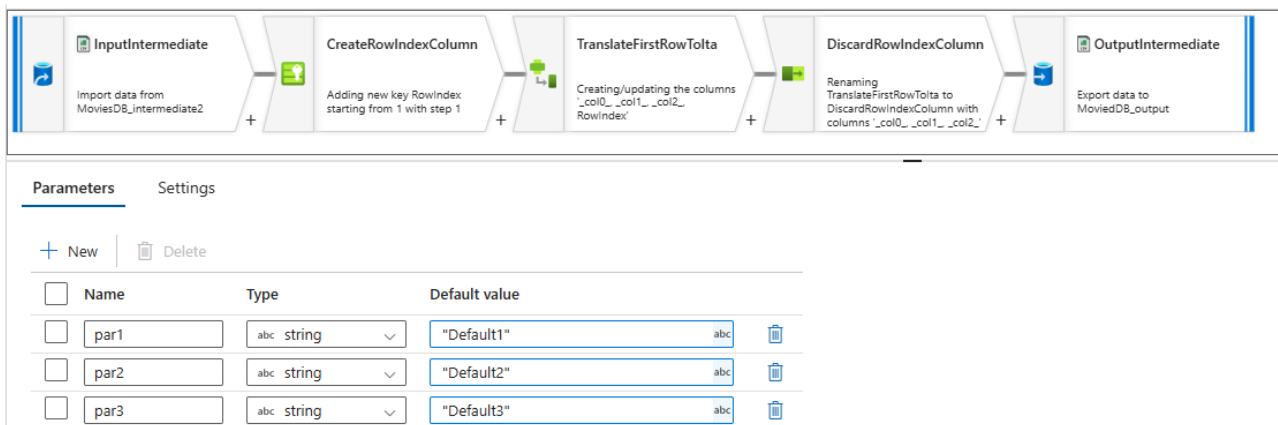


| Name | Type | Default value |
|---|---|---|
| par1 | abc string | "Default1" |
| par2 | abc string | "Default2" |
| par3 | abc string | "Default3" |

Fig. 15 – Creation of Data Flow parameters.



**Data flow parameters** ⓘ

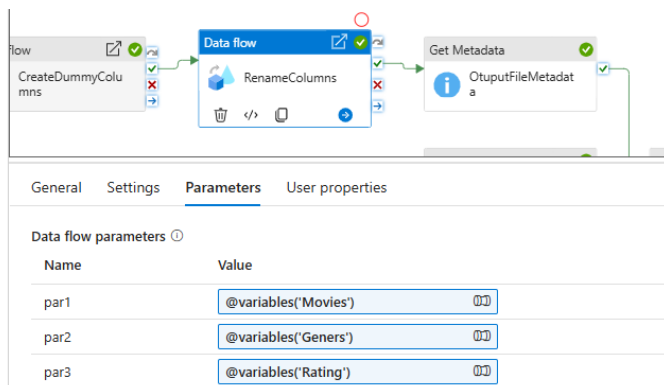| Name | Value |
|---|---|
| par1 | @variables('Movies') |
| par2 | @variables('Geners') |
| par3 | @variables('Rating') |

Fig. 16 – Linking in the pipeline of parameters to data flow variables.
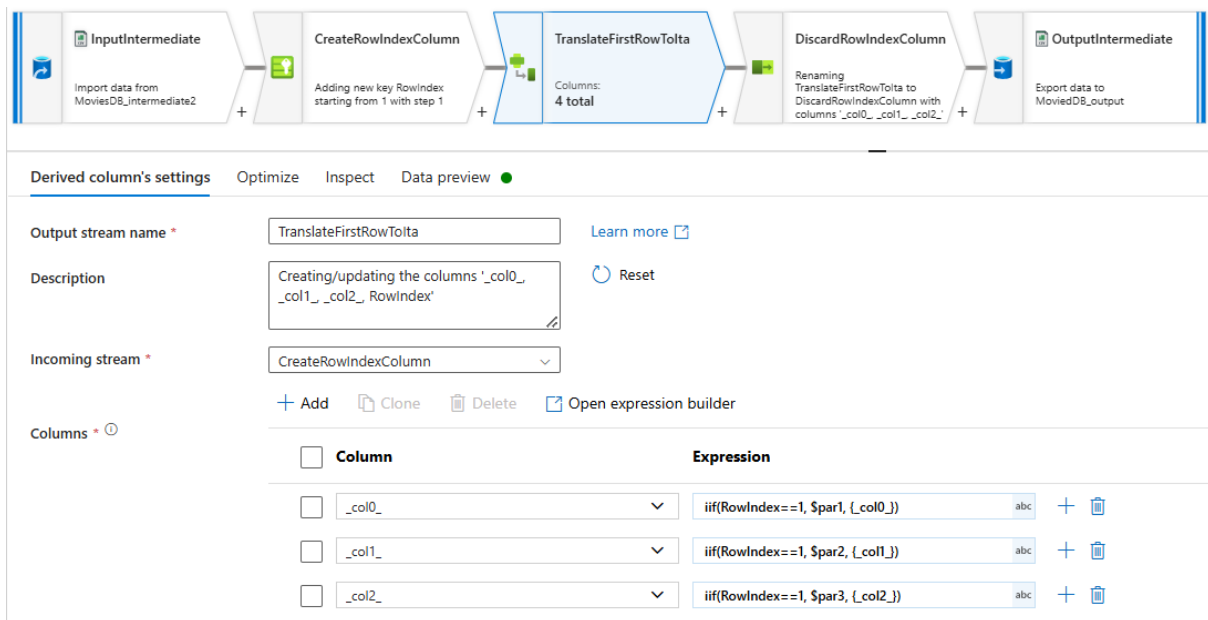
Fig. 17 – Replacement of the dummy columns with the values of the Data Flow parameters, i.e. the pipeline variables.

As there is currently no activity to dynamically change the column names, a trick was found to create dummy columns using an additional dataflow – 'Create DummyColumns' (Fig. 18). In this way, the column names are values that are inserted into a row of the dataset, making them replaceable with the parameters created via a 'Derived Column'.
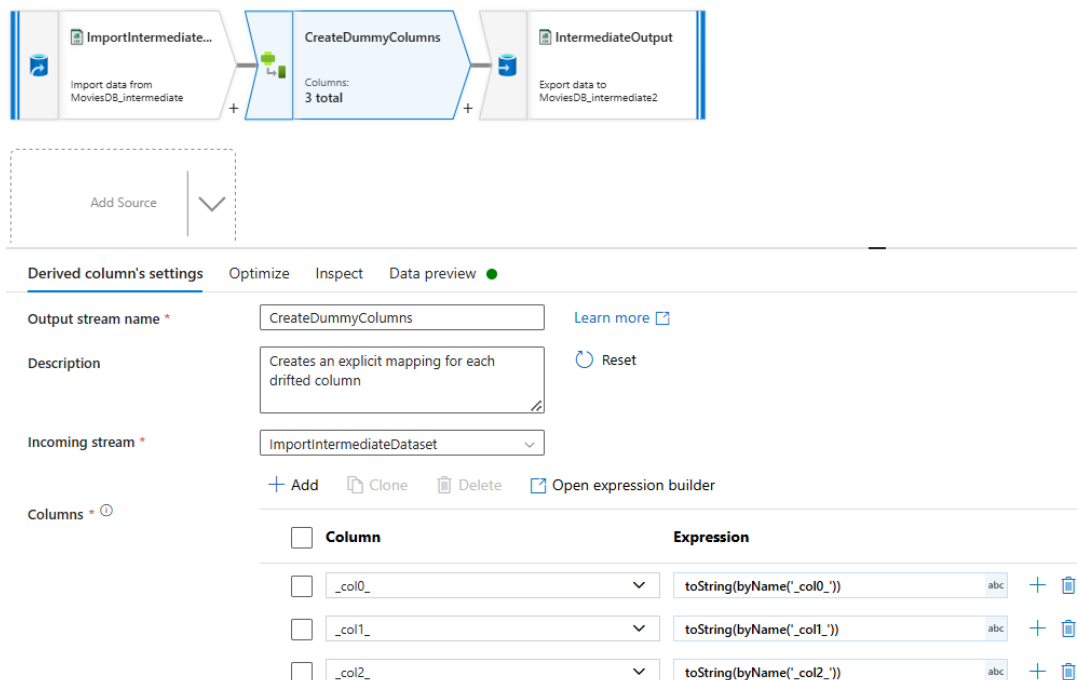


Fig. 18 – Use of a Data Flow to create dummy columns for the purposes described..

In parallel, a Get Metadata activity was used to extract information from the original dataset.

A Shared Access Signature (SAS) web activity was used to obtain custom metadata (Fig. 19), using the SAS URL blob (Fig. 20).

Fig. 19 – Custom Metadata.



Fig. 20 – URL to be inserted in the web activity to allow temporary access to the dataset.

All metadata, both from the original dataset and the final dataset, were stored in the final output dataset via a 'Copy data' activity (Fig.21 and Fig. 22).



Fig. 21 – Metadata saved in the final output file.

# MoviesDB_output.csv
Blob

Save  ✕ Discard  ⬇ Download  ↻ Refresh  🗑 Delete  ⇄ Change tier  Acquire lease  Break lea

| | |
|---|---|
| CONTENT-DISPOSITION | |
| LEASE STATUS | Unlocked |
| LEASE STATE | Available |
| LEASE DURATION | - |
| COPY STATUS | - |
| COPY COMPLETION TIME | - |

**Undelete**

Metadata

| Key | Value | |
|---|---|---|
| OriginalFileColumnNumber | 5 | 🗑 |
| OriginalFileSize | 450221 | 🗑 |
| OriginalFileLastModified | 2024-11-06T10:59:47Z | 🗑 |
| OriginalFileName | moviesDB.csv | 🗑 |
| OriginalFileCustomMetadata_Author | Au-Thor | 🗑 |
| OriginalFileCustomMetadata_Format | CSV | 🗑 |
| OutputFileColumnNumber | 3 | 🗑 |
| OutputFileSize | 71812 | 🗑 |
| OutputFileLastModified | 2024-11-07T18:16:39Z | 🗑 |
| | | |

---

Upload  🔒 Change access level  ⋯

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

Search blobs by prefix (case-...)

⬤ Show deleted blobs

➕ Add filter

**Name**

☐ 📄 MoviesDB_output.csv  ⋯

*Fig. 22 – Output File Metadata.*