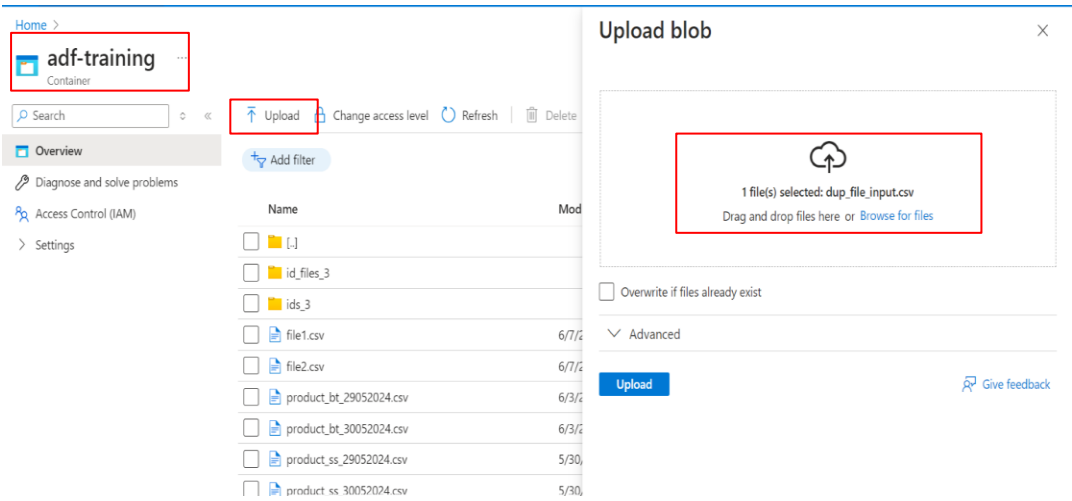# Formal Documentation of Azure Data Factory Pipeline – Training

**Use Case**: Removing Duplicate Data from a File

**Steps:**

## 1. Loading the File into the Container.

- A container named adf-training was previously created in the Azure storage account.
- The files were uploaded into the container, in the input folder (refer to image 1).
- File name: dup_file_input.csv.
- Once the upload was successful, the files were added to the input folder.
- Refer image 2 to review the input data that contains duplicate data.



-Image 1

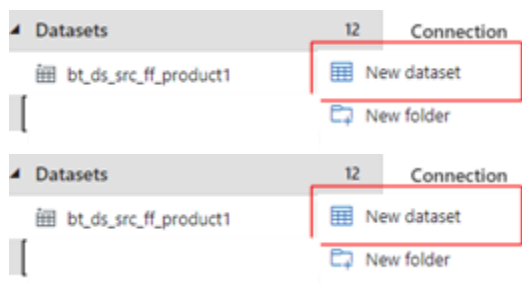Linked service: ADF_Training

Object:    dup_file_input.csv

| 🎛 | id | first_name | last_name | department |
|---|---|---|---|---|
| 1 | 1 | Raj | Challa | IT |
| 2 | 2 | Babita | Tiwari | IT |
| 3 | 3 | Siraj | Shaikh | IT |
| 4 | 1 | Raj | Challa | IT |
| 5 | 4 | Kamran | Khan | IT |
| 6 | 4 | kamran | Khan | IT |

- Image 2

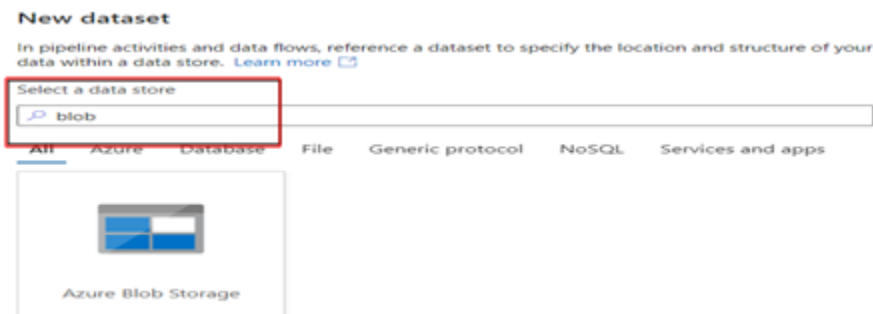## 2. Dataset Creation

### Source Dataset

- A new source dataset was created for the source data files.
- Source dataset name: bt_ds_src_dup_transformation.
- Inside Azure Data Factory, in the Author tab, select the Dataset option and click on "New Dataset" (refer to image 1).
- Choose the Azure Blob Storage option (refer to images 2 and 3). Next, select the Delimited Text file format, which brings you to the properties page where you define the dataset name and path.
- Specify the dataset name, select the linked service, and provide the path of the input file (refer image 4)
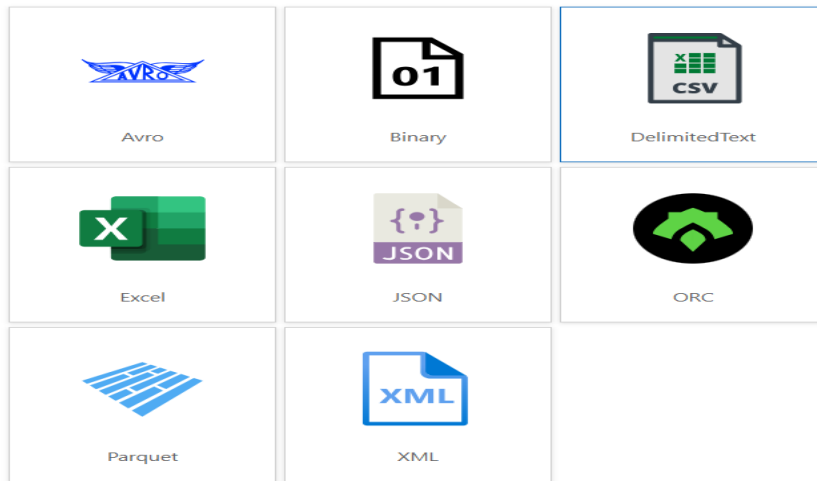- These steps create the source dataset.



- Image 1



- Image 2

**Select format**

Choose the format type of your data

| | | |
|---|---|---|
| Avro | Binary | DelimitedText |
| Excel | JSON | ORC |
| Parquet | XML | |

Continue    Back                                          Cancel   - Image 3

**Set properties**

Name

bt_ds_src_dup_transformation

Linked service *

ADF_Training

File path

adf-training  /  Input  /  dup_file_input.csv

First row as header    ☑

Import schema

◉ From connection/store    ◯ From sample file    ◯ None

> Advanced

- Image 4.

**Target Dataset**

- Follow similar steps for the target dataset. In the Author tab, select the Datasets option, click on "New Dataset," select Azure Blob Storage, and then select the Delimited Text format, which brings you to the properties page where you define the dataset name and path.
- Assigned path: adf-training/output (refer image 1)
- Target dataset name: bt_ds_tgt_dup_transformation.

## Set properties

**Name**

bt_ds_tgt_dup_transformation

**Linked service** *

ADF_Training

**File path**

adf-training / Output / File name

**First row as header** ☑

**Import schema**

◉ From connection/store   ○ From sample file   ○ None

> Advanced

- Image 1

## 3. Creating Dataflow

- A dataflow named bt_df_dup_transformation was created to handle duplicate entries.
- The mandatory step is to enable dataflow debug.

Below are the steps:

A. Source

- In the source settings, the source file containing the duplicate data was added.

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

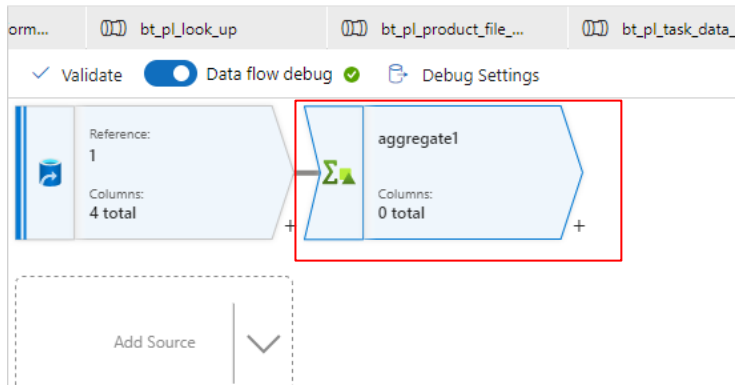| Output stream name * | source1 | Learn more ☐ |
| Description | Import data from bt_ds_src_dup_transformation | ↻ Reset |
| Source type * | Dataset \| Inline | |
| Dataset * | bt_ds_src_dup_transformation | Test connection |

- Image 1

B. Aggregate Function

- Add an aggregate transformation to remove duplicates based on the "id" column (refer image 1)
- In "group by," specified the id column (refer image 2).
- In "aggregates," select three columns: first_name, last_name, department, and use the first expression (refer image 3).
- The first expression in the Aggregate function ensures that for each unique id, the first occurrence of the first_name, last_name, and department columns is selected (refer image 4).
- The duplicate data was successfully removed. I checked the same in the data preview tab (refer image 5).

✓ Validate | ⬤ Data flow debug ✓ | Debug Settings

Reference:
1

Columns:
4 total

aggregate1

Columns:
0 total

Add Source ∨

- Image 1

Aggregate settings | Optimize | Inspect | Data preview ⬤

Output stream name *          aggregate1          Learn more ⬈

Description          Aggregating data by 'id' producing columns 'first_name, last_name, department'          ↻ Reset

Incoming stream *          source1 ∨

Group by | Aggregates

Columns          Name as

abc  id ∨          id          + 🗑

- Image 2

Aggregate settings | Optimize | Inspect | Data preview ⬤

Incoming stream *          source1 ∨

Group by | Aggregates

Grouped by: id

+ Add | 📋 Clone | 🗑 Delete | ⬀ Open expression builder

| ☐ | Column | Expression | | |
|---|---|---|---|---|
| ☐ | first_name ∨ | first(first_name) | abc | + 🗑 |
| ☐ | last_name ∨ | first(last_name) | abc | + 🗑 |
| ☐ | department ∨ | first(department) | abc | + 🗑 |

-Image 3

**Dataflow expression builder**
Σ₄ aggregate1

**Aggregate Columns**

+ Create new ∨

abc  first_name

abc  last_name

abc  department

**Column name** *

first_name

**Expression**

first(first_name)

— Image 4

Aggregate settings    Optimize    Inspect    **Data preview** ●

Number of rows + INSERT 4          ⁕ UPDATE 0          ✕ DELETE 0

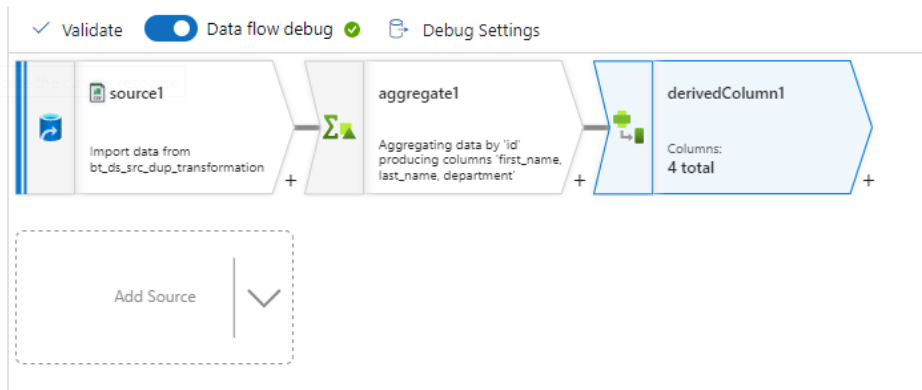↻ Refresh | ∨    Typecast ∨    ⬚ Modify ∨    ⬚ Map drifted    ⬚ Statistics  ✕ Remove    ↓ Export to CSV | ∨
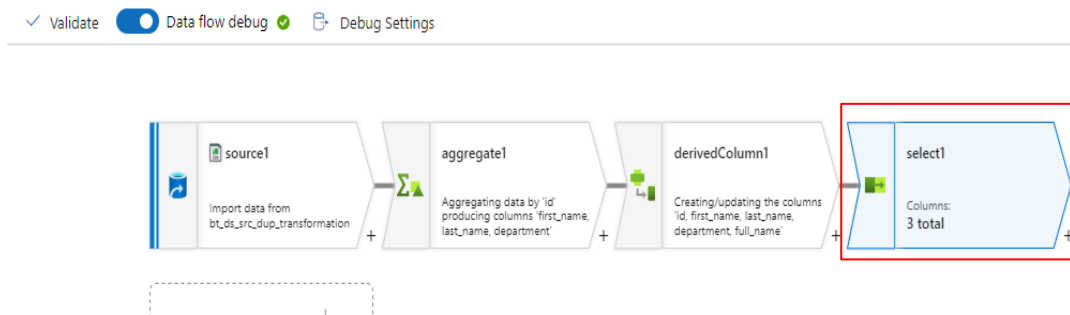
| ↑↓ | id | abc ↑↓ | first_name | abc ↑↓ | last_name | abc ↑↓ | department |
|---|---|---|---|---|---|---|---|
| ✦ | 1 | | Raj | | Challa | | IT |
| ✦ | 2 | | Babita | | Tiwari | | IT |
| ✦ | 3 | | Siraj | | Shaikh | | IT |
| ✦ | 4 | | Kamran | | Khan | | IT |

— Image 5

C. Derived Column

- Select the derived column (refer to image 1) and add a new column named full_name to concatenate the first and last names in the file.
- The original file had 4 columns: id, first_name, last_name, department.
- In derived column settings, click on "add column," assign the column name as full_name, and click on "expression" (refer to image 2).
- Wrote a concat expression to create a new column named full_name (refer image 3).
- I then checked my data through data preview and the new column named full_name was added (refer image 4).

- Image 1.

Optimize    Inspect    Data preview ●

| Output stream name * | derivedColumn1 | Learn more ⧉ |
|---|---|---|

Description          Creating/updating the columns 'id,          ↻ Reset
                     first_name, last_name, department,
                     full_name'

| Incoming stream * | aggregate1 ⌄ |
|---|---|

+ Add    ⧉ Clone    🗑 Delete    ⧉ Open expression builder

Columns * ⓘ

| ☐ | Column | Expression |
|---|---|---|
| ☐ | full_name ⌄ | concat(first_name, " ", last_name)  abc  + 🗑 |

- Image 2

## Dataflow expression builder
🔾 derivedColumn1

**Derived Columns**

+ Create new ⌄

    abc    full_name

**Column name** *

full_name

**Expression**

    concat(first_name, " ", last_name)

- Image 3

| | | | | | |
|---|---|---|---|---|---|
| Number of rows ＋ **INSERT** 4 | | ✳ **UPDATE** 0 | ✕ **DELETE** 0 | ✴ **UPSERT** 0 | |

○ Refresh | ∨   Typecast ∨   ⊞ Modify ∨   ⌯ Map drifted   ⊟ Statistics  ✕ Remove   ⭳ Export to CSV | ∨

| | id | first_name | last_name | department | full_name |
|---|---|---|---|---|---|
| ✚ | 1 | Raj | Challa | IT | Raj Challa |
| ✚ | 2 | Babita | Tiwari | IT | Babita Tiwari |
| ✚ | 3 | Siraj | Shaikh | IT | Siraj Shaikh |
| ✚ | 4 | Kamran | Khan | IT | Kamran Khan |

Image 4

## D. Select

- Use the select function to get only the required columns (refer to image 1).
- After the derived function, there were 5 columns, but only id, full_name, and department were needed.
- In select settings, delete the other 2 columns, first_name and last_name, to get the data in the desired format (refer image 2).
- I also changed the position of the full name column from last to middle by just dragging and moving.
- Check the data in the Data Preview tab (refer to image 4).

✓ Validate   ⬤ Data flow debug ✓   ⬚ Debug Settings

| source1 | aggregate1 | derivedColumn1 | select1 |
|---|---|---|---|
| Import data from bt_ds_src_dup_transformation | Aggregating data by 'id' producing columns 'first_name, last_name, department' | Creating/updating the columns 'id, first_name, last_name, department, full_name' | Columns: 3 total |

- Image 1

| | | |
|---|---|---|
| Incoming stream * | derivedColumn1 ⌄ | |
| Options | ☑ Skip duplicate input columns ⓘ | |
| | ☑ Skip duplicate output columns ⓘ | |
| Input columns * | ☐ Auto mapping ⓘ   ○ Reset   ＋ Add mapping   🗑 Delete | |

| ☐ | derivedColumn1's column | ▽ | Name as | | |
|---|---|---|---|---|---|
| ☐ | abc id ⌄ | → | id | ＋ | 🗑 |
| ☐ | abc full_name ⌄ | → | full_name | ＋ | 🗑 |
| ☐ | abc department ⌄ | → | department | ＋ | 🗑 |

- Image 2

- Image 3

E. Sink

- In sink, select the target dataset to copy the cleaned file to the output folder (refer to image 1, 2).
- I also specified the name of my output file (refer image 3).



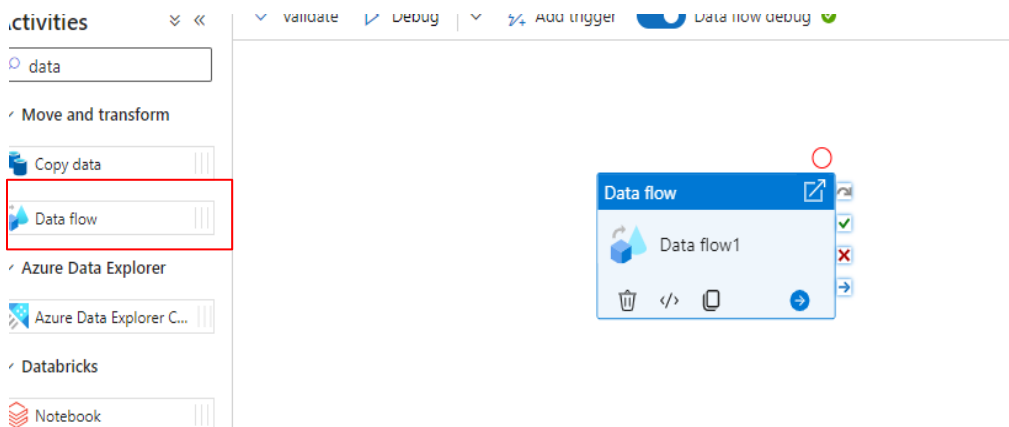- Image 1



Image 2



- Image 3

## 4. Pipeline Creation

- Create a pipeline named bt_pl_dup_transformation to execute the dataflow activities (refer to image 1).



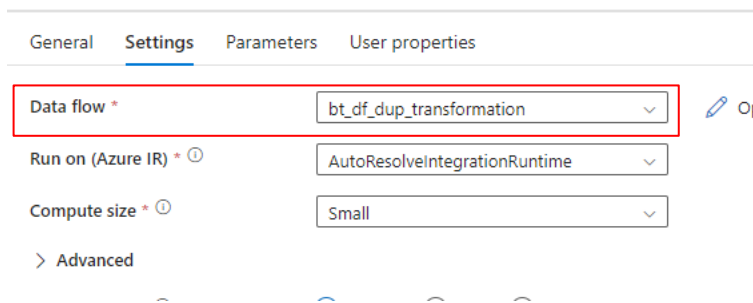- Image 1

## 5. Dataflow Activity

- From the activities, drag and drop the dataflow activity (refer to image 1).
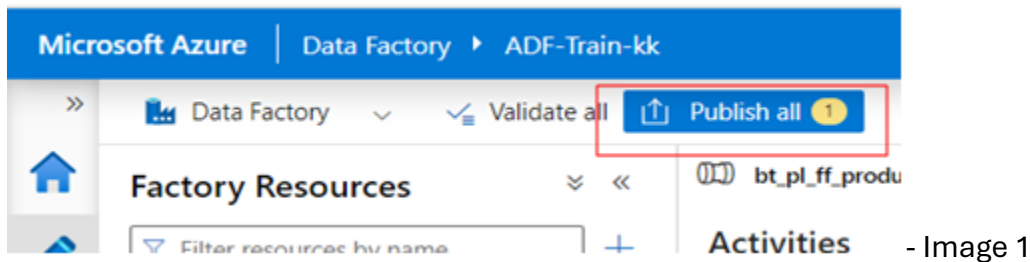- In the settings, selected my dataflow. (refer to image 2).
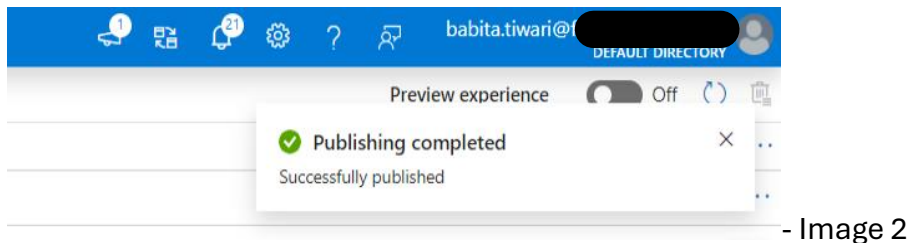


- Image 1



-Image 2

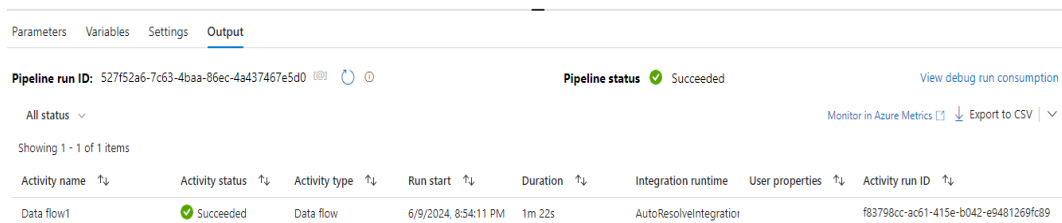## 6. Publishing and Executing the Pipeline

- Save/publish the activities, and successfully publish and execute them (refer to images 1, 2, 3).
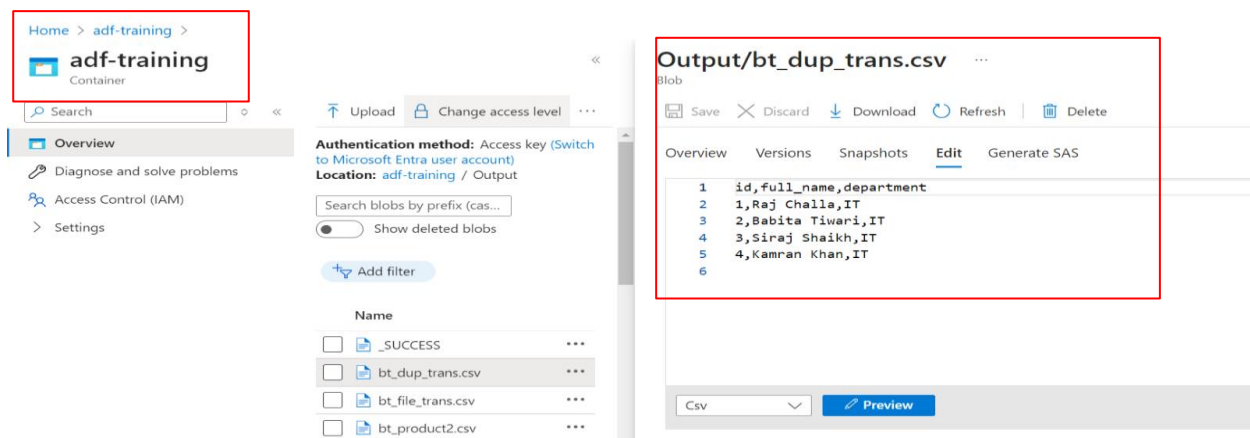


- Image 1



- Image 2



- Image 3

## 7. Verifying the Output

- Visited the output folder in the adf-training container and checked the output.
- The file was generated correctly, and duplicate data was removed from the file (refer image 1)



- Image 1

**Summary:**

In this use case, duplicate data was removed from a file using Azure Data Factory (ADF). The process included loading the file into a container, creating source and target datasets, setting up a dataflow for duplicate handling, and configuring aggregate and derived column transformations. A pipeline was created to execute these activities, and the results were verified in the output folder. The file was processed successfully, with duplicates removed and desired columns retained.