# Project:

Firstly, I created 2 storage account

1. airlinedevacc
2. airlineprodacc



Now i created a container name: **united-airlines** and under it i created **landing-zn and processed-data**

Under **landing-zn** I upload 2 csv files

united-airlines  >  📁 landing-zn

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

| | Name | Last modified | Access tier | Blob type | Size | Lease state |
|---|---|---|---|---|---|---|
| ☐ | 📁 [..] | | | | | |
| ☐ | 📄 airports.csv | 11/30/2025, 2:26:52 PM | Hot (Inferred) | Block blob | 15.93 KiB | Available |
| ☐ | 📄 flights.csv | 11/30/2025, 2:27:09 PM | Hot (Inferred) | Block blob | 20.97 MiB | Available |

Showing all 2 items

Only show active objects

## DailyFlightSource (flight.csv)

```
Carrier | OriginAirportID | DestAirportID | DepDelay | ArrDelay
DL      | 11433           | 13303         | -3       | 1
DL      | 14869           | 12478         | 0        | -8
DL      | 14057           | 14869         | -4       | -15
```

## 2. AirlineDimSource (airport.csv)

```
airport_id | city      | state | name
11433      | Hilo      | HI    | Hilo International
13303      | Miami     | FL    | Miami International
14869      | San Diego | CA    | San Diego International
```

Now i created a linked service as to connect data factory with adls so as to fetch dataset files

**Connections**

- 🗄 Linked services
- 🔁 Integration runtimes
- 👁 Microsoft Purview
- 🖋 ADF in Microsoft Fabric

Filter by name          Annotations : **Any**

Showing 1 - 1 of 1 items

| Name | Type | Related |
|---|---|---|
| 🔷 AzureDataLakeStorage1 | Azure Data Lake Storage Gen2 | 3 |

Now I created two Azure Data Factories – one for development and one for production. The development Data Factory is connected to Azure DevOps for source control.

Whenever I make changes in the Dev environment, I publish them, and through a CI/CD pipeline in Azure DevOps, those changes are automatically deployed to the production Data Factory.

This setup ensures proper environment separation, version control, and automated deployment without manual intervention.





## Data Flow Overview

I created a **Mapping Data Flow** in Azure Data Factory to transform and combine airline-related data. The flow consists of:

1. **Source Datasets**
   a. **AirlineDimSource** → Reads airports.csv from Azure Data Lake (landing-zn folder).
   b. **DailyFlightSource** → Reads flights.csv from Azure Data Lake (landing-zn folder).
   c. **ProcessedData** → Stores transformed data in Azure Data Lake.
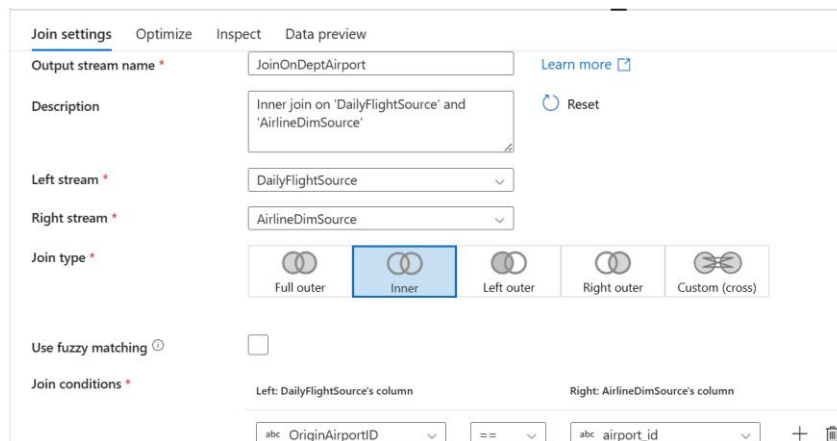
**JoinOnDeptAirport**

Under it I configured an inner join between the flight data and the airport dimension data. The join condition matches the OriginAirportID from the flight dataset with airport_id from the dimension dataset. This helps combine flight records with airport details.



## Select activity

After joining the datasets, I used a Select transformation to choose relevant columns and rename them for clarity. For example, the airport name column was renamed to DepAirportName, and city/state were renamed to DepCity and DepState.

## Join Activity

I added another inner join to match the destination airport ID from the flight dataset with airport_id from the dimension dataset. This step adds destination airport details to the data.

The second join is needed because the flight dataset has two keys: OriginAirportID and DestAirportID. The first join adds departure airport details, and the second join adds arrival airport details.



### Select Activity

I used Select to keep only the required columns and rename them for clarity. For example, city/state/name from the arrival airport were renamed to ArrCity, ArrState, and ArrAirportName.

## Select settings Optimize Inspect Data preview

| | |
|---|---|
| Output stream name * | SelectTransOnArrivalAirport |
| Description | Renaming JoinOnArrivalAirportID to SelectTransOnArrivalAirport with columns 'Carrier, DepDelay, ArrDelay, |
| Incoming stream * | JoinOnArrivalAirportID |
| Options | ☑ Skip duplicate input columns ⓘ |
| | ☑ Skip duplicate output columns ⓘ |

Learn more ↗

↻ Reset

| JoinOnArrivalAirportID's column | Name as |
|---|---|
| abc Carrier | Carrier |
| abc DepDelay | DepDelay |
| abc ArrDelay | ArrDelay |
| abc DepCity    ArrDelay | DepCity |
| abc DepState | DepState |
| abc DepAirportName | DepAirportName |
| abc city | ArrCity |
| abc state | ArrState |
| abc name | ArrAirportName |

## Sink Configuration Details

Finally, I used a Sink transformation to write the processed data back to Azure Data Lake in the processed-data folder. I enabled schema drift, auto mapping, and cleared the folder before writing to ensure clean output.



## Sink Settings Errors Mapping Optimize Inspect Data preview

| | |
|---|---|
| Output stream name * | WriteProcessedData |
| Description | Export data to ProcessedData |
| Incoming stream * | SelectTransOnArrivalAirport |
| Sink type * | ▦ Dataset   ▨ Inline   ▤ Cache |
| Dataset * | ▤ ProcessedData    Test connection ✎ Open + New |
| Skip line count | |
| Options | ☑ Allow schema drift ⓘ |
| | ☐ Validate schema ⓘ |

Learn more ↗

↻ Reset

| Sink | **Settings** | Errors | Mapping | Optimize | Inspect | Data preview |
|------|----------|--------|---------|----------|---------|--------------|

Clear the folder ☑

File name option * [ Default ▾ ]

Quote All ⓘ ☐

Headers ⓘ [ Enter expression...                    ANY ]

Umask ⓘ

|        |      |      |      |
|--------|------|------|------|
| Owner  | ☑ R  | ☑ W  | ☑ X  |
| Group  | ☑ R  | ☑ W  | ☑ X  |
| Others | ☑ R  | ☑ W  | ☑ X  |

| Sink | Settings | Errors | **Mapping** | Optimize | Inspect | Data preview |
|------|----------|--------|---------|----------|---------|--------------|

Options    ☑ Skip duplicate input columns ⓘ

☑ Skip duplicate output columns ⓘ

☑ Auto mapping ⓘ    ↻ Reset    + Add mapping    🗑 Delete    [ Output format ]

# Now I started creating pipeline

▲ Pipelines                           3
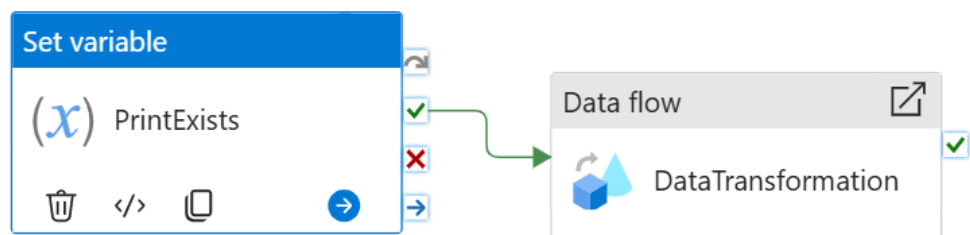
  ◫ AirlinePipeline

## FindDailyFlightFiles

I used Get Metadata to check if the source file exists by retrieving the 'Exists' property. This value drives an If Condition activity, ensuring the pipeline only runs when the file is available

| General | Settings | User properties |
| --- | --- | --- |

Dataset *   📄 DailyFlightSource  ⌄    ✏ Open   + New    Learn more ↗

Field list *   + New  | 🗑 Delete
　　　　　　☐ **Argument**
　　　　　　☐ Exists   ⌄

I used a Set Variable activity to store the result of the Get Metadata check using the expression `@activity('FindDailyFlightFiles').output.exists`. This value drives the If Condition logic, ensuring the pipeline only runs when the file exists.

🔖 AirlinePipeline > 🔀 If Condition1

**Set variable**

$(x)$ PrintExists

🗑  </>  ▢  ➡

**Data flow** ↗

DataTransformation

| General | **Settings** | User properties |
| --- | --- | --- |

Variable type ⓘ      ◉ Pipeline variable   ◯ Pipeline return value

Name *      PrintExists  ⌄   + New

Value      @activity('FindDailyFlightFiles').outp...

I configured the Data Flow activity to run on AutoResolveIntegrationRuntime with a small compute size for cost efficiency. It executes the AirlineDataTransformation flow only when the source file exists, and verbose logging helps in monitoring and troubleshooting.

Set variable

$(x)$  PrintExists

Data flow

DataTransformation

| General | Settings | Parameters | User properties |

Data flow *                    AirlineDataTransformation            ⌄          ✎ Open

Run on (Azure IR) * ⓘ         ✓ AutoResolveIntegrationRuntime      ⌄

Compute size * ⓘ               Small                                ⌄

> Advanced

Logging level * ⓘ             ● Verbose    ○ Basic    ○ None

In the False branch, I set a variable using the same dynamic expression to capture the file existence status. This allows me to log or trigger alerts when the source file is missing, making the pipeline more reliable.

📖 AirlinePipeline > ⚙ If Condition1 - False activities

$(x)$  PrintDoesNotExist

| General | Settings | User properties |

Variable type ⓘ          ● Pipeline variable    ○ Pipeline return value

Name *                    FalseParam                          ⌄         + New

Value                     @activity('FindDailyFlightFiles').outp...

## Pipeline expression builder

Add dynamic content below using any combination of expressions, func

```
@activity('FindDailyFlightFiles').output.exists
```

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

| 🔍 Search blobs by prefix (case-sensitive) | | Only show active objects |

Showing all 5 items

| | Name | Last modified | Access tier | Blob type | Size | Lease state |
|---|---|---|---|---|---|---|
| ☐ | 📁 [..] | | | | | |
| ☐ | 🗋 _SUCCESS | 11/30/2025, 5:57:16 PM | Hot (Inferred) | Block blob | 0 | Available |
| ☐ | 🗋 part-00000-ffdc7... | 11/30/2025, 5:57:16 PM | Hot (Inferred) | Block blob | 26.43 MiB | Available |
| ☐ | 🗋 part-00001-ffdc7... | 11/30/2025, 5:57:16 PM | Hot (Inferred) | Block blob | 26.42 MiB | Available |
| ☐ | 🗋 part-00002-ffdc7... | 11/30/2025, 5:57:16 PM | Hot (Inferred) | Block blob | 25.99 MiB | Available |
| ☐ | 🗋 part-00003-ffdc7... | 11/30/2025, 5:57:14 PM | Hot (Inferred) | Block blob | 9.55 MiB | Available |

# 1. Connected Dev Data Factory to Azure DevOps

I linked my Development Data Factory to an Azure DevOps Git repository. This allows all pipelines, datasets, and data flows to be version-controlled. Any changes I make in ADF are first saved in the collaboration branch before publishing.
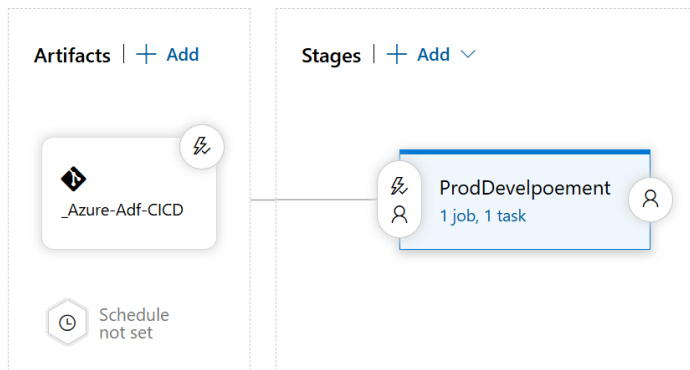


I connected my Development Data Factory to an Azure DevOps Git repository for source control. The collaboration branch is `main`, where I make changes, and the

publish branch is `adf_publish`, which stores ARM templates generated when I publish. This setup ensures version control and enables CI/CD automation.

When I publish changes from Dev Data Factory, it automatically generates ARM templates and commits them to the `adf_publish` branch in Azure DevOps. These templates represent the entire ADF configuration, including pipelines, datasets, and linked services. This branch acts as the source for the CI/CD pipeline.



Now I created a Release Pipeline in Azure DevOps that takes the ARM templates from the **artifact** generated during the build stage and deploys them to the Production Data Factory. This **stage** runs automatically after publishing changes in Dev Data Factory, ensuring a smooth and consistent deployment process.

## Under Artifact

After adding these settings in the artifact section, the pipeline knows which repository and branch to use. This is important because without linking the artifact, the release pipeline cannot fetch the latest ARM templates. Once this is set, every time I publish from the development Data Factory, the pipeline picks up the new changes and deploys them to production.

After linking the artifact and enabling the continuous deployment trigger, the pipeline automatically creates a release whenever there is a new commit in the `adf_publish` branch. This means every time I publish changes from the development Data Factory, the updated ARM templates are pushed to `adf_publish`, and the pipeline deploys those changes to production without manual intervention.



After configuring the artifact and enabling the trigger, I created a stage called `ProdDevelopement` for production deployment. I am the stage owner, which means I control approvals and deployment settings. Inside this stage, I added an agent job that runs on my selected agent pool. The main task in this job is an ARM Template Deployment, which deploys the Data Factory resources to the production environment using the templates from the `adf_publish` branch.

These jobs show that my self-hosted agent executed the release pipeline tasks for production deployment. Each job corresponds to a release triggered by changes in the `adf_publish` branch. The agent picked up the job instantly and completed the ARM template deployment task successfully.



Whatever I was having in dev data factory is also now in prod data factory