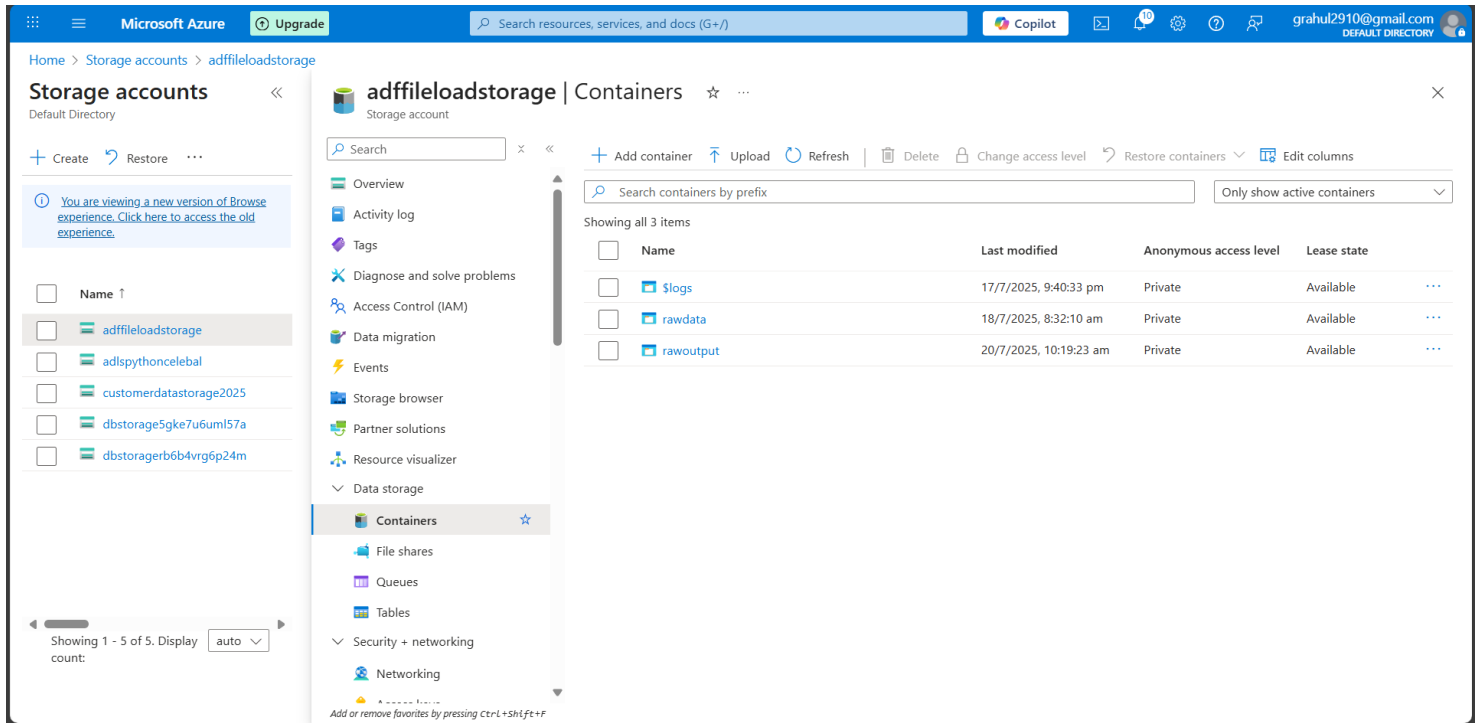# FILE COPY WITH TRANSFORMATION
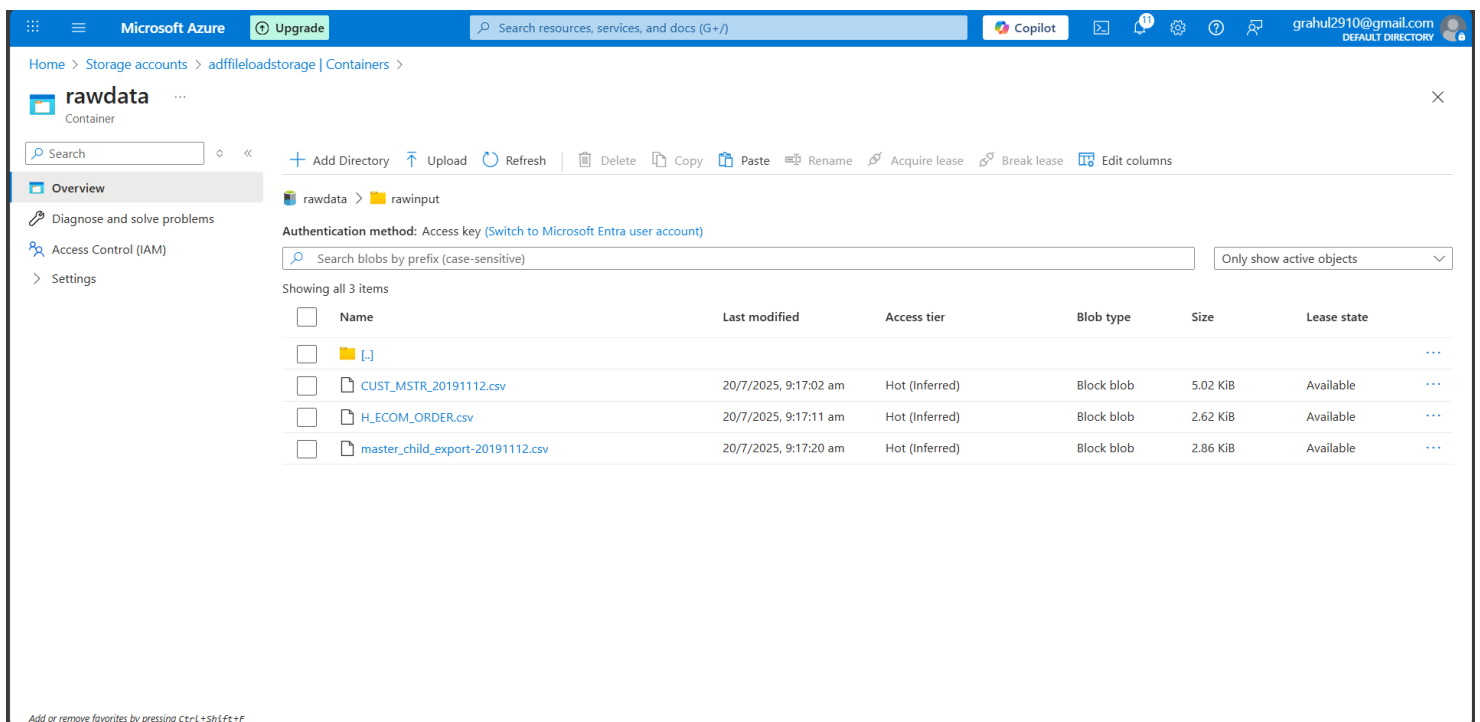
1. Setup Resources

- Created Azure Data Lake Gen2 container and uploaded all types of test files.



- Created Azure SQL Database with the three required tables:
    - CUST_MSTR(Date, …)
    - master_child(Date, DateKey, …)
    - H_ECOM_Orders(…)

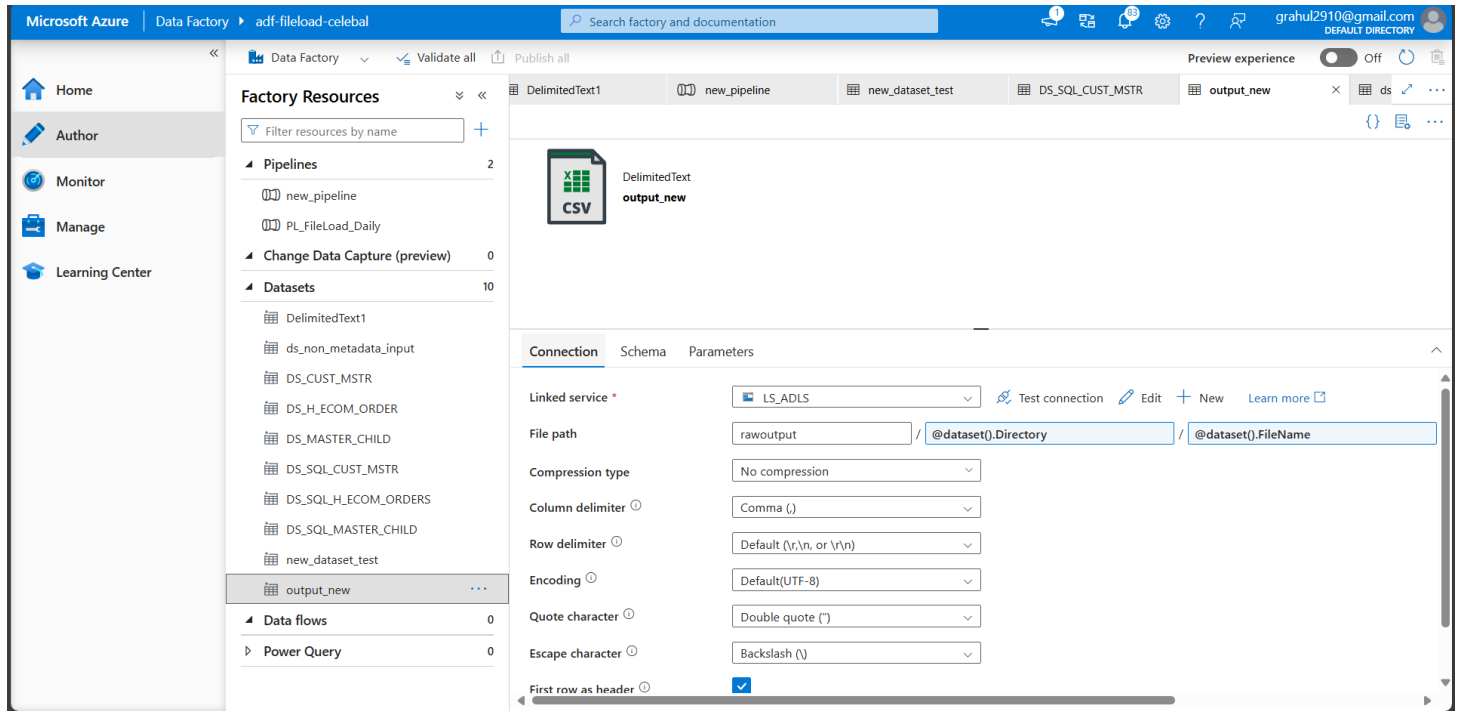| | A | B | C |
|---|---|---|---|
| 1 | CustomerID | Name | Location |
| 2 | 1 | Samaira Kar | Mirzapur |
| 3 | 2 | Ira Agate | Purnia |
| 4 | 3 | Lakshit Kaur | Khammam |
| 5 | 4 | Rasha Rastogi | Nagaon |
| 6 | 5 | Saira Bhatia | Mirzapur |
| 7 | 6 | Kavya Sunder | Ghaziabad |
| 8 | 7 | Kartik Butala | Bangalore |
| 9 | 8 | Nirvaan Kibe | Bhusawal |
| 10 | 9 | Advika Singhal | Thanjavur |
| 11 | 10 | Advika Baria | Kurnool |
| 12 | 11 | Sana Barad | Amroha |
| 13 | 12 | Shray Sathe | Kolhapur |
| 14 | 13 | Onkar Swaminathan | Ghaziabad |
| 15 | 14 | Dhruv Goyal | Surendranagar Dudhrej |
| 16 | 15 | Nayantara Dey | Varanasi |
| 17 | 16 | Abram Karan | Phagwara |
| 18 | 17 | Nehmat Apte | Aurangabad |
| 19 | 18 | Ahana  Toor | Bilaspur |
| 20 | 19 | Nirvaan Mangat | Malegaon |
| 21 | 20 | Tanya Rout | Bokaro |
| 22 | 21 | Lakshay Sangha | Naihati |
| 23 | 22 | Krish DAlia | Allahabad |
| 24 | 23 | Mannat Kashyap | Kulti |
| 25 | 24 | Anahi Thaker | Khandwa |
| 26 | 25 | Adah Kaul | Raurkela Industrial Township |
| 27 | 26 | Nitya Sachar | Sultan Pur Majra |
| 28 | 27 | Sara Guha | Saharsa |

## 2. Created ADF Pipeline

- Designed a single pipeline with:

  - A Get Metadata activity to list all files from the container.

  - A ForEach activity to iterate through each file.

Make the datasets and make them parameterized, so that they can get values
directly from the file name



3. Used Set Variable Activities

- Extracted fileName using @item().name.

- Created two variables:

  o   file_date → Format YYYY-MM-DD

  o   CUST_date_format → Format YYYYMMDD

- Used expressions like:

substring(replace(item().name, '.csv', ''), length(...) - 8, 8)

to extract the date from the filename.

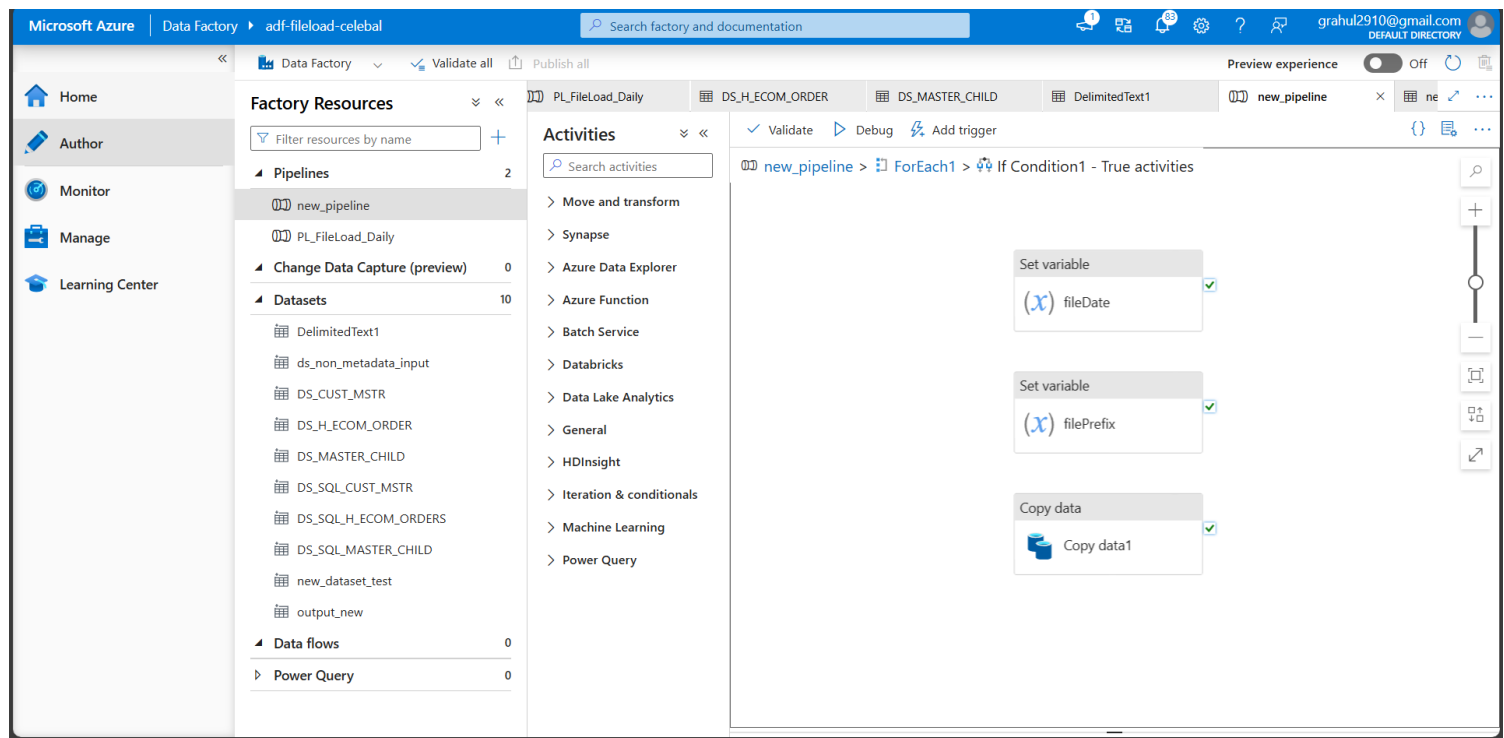- Transformed the extracted dateKey into proper date format using:

@concat(substring(variables(CUST_date_format '), 0, 4), '-', substring(variables(' CUST_date_format '), 4, 2), '-', substring(variables(' CUST_date_format '), 6, 2))

4. Used If Conditions

- Inside the ForEach, placed 3 If Conditions:
    - If filename starts with CUST_MSTR
    - If filename starts with master_child_export
    - If filename starts with H_ECOM_ORDER

## 5. Defined Data Flows

- For CUST_MSTR:

  o Used Data Flow to add one derived column: Date from file_date

  o Sink: Load to CUST_MSTR table (with truncate behavior)

- For master_child_export:

  o Used Data Flow to add two derived columns:

    ▪ Date from file_date

    ▪ DateKey from CUST_date_format

  o Sink: Load to master_child table

- For H_ECOM_ORDER:

  o Direct copy activity (or simple data flow) without any transformations.

  o Sink: Load to H_ECOM_Orders table

## 6. Truncate Load Logic

- Enabled truncate option in each sink to ensure tables are cleared before inserting new data.

## 7. Debug and Testing

- The pipeline should be debugged in ADF to:

  o Confirm file iteration logic

  o Validate column additions

  o Ensure successful data writes

o    Monitor any errors or data issues



## 8. Output Folder Structure

- After the pipeline runs, the processed files are stored in the rawdata/rawoutput/ directory of the data lake.



- New folders are automatically created for each file type inside rawoutput.

- For example, if the file is master_child_export-20191112.csv, the file will be stored in:

rawdata/rawoutput/master_child/

- This folder naming is based on the first part of the file name (before the date or extension), such as:

  o CUST_MSTR → folder cust_mstr/

  o master_child_export → folder master_child/

  o H_ECOM_ORDER → folder h_ecom_order/

- CUST_MSTR Files
  These files will have one additional column:

    A Date column will be added based on the date extracted from the file name (in the format YYYY-MM-DD).

    All other columns will remain unchanged.

    There will be only one new column added in this case.

- master_child_export Files
  These files will have two new columns added:

    A Date column (format: YYYY-MM-DD)

    A DateKey column (format: YYYYMMDD)

    These columns will be appended at the end of each row during transformation.

- H_ECOM_ORDER Files
  These files are loaded as-is, with no transformation or additional columns.

    The output will exactly match the structure of the input file.

    This is because these files do not contain any embedded date in the filename.

---

# master_child_export-20191112.csv
Blob

💾 Save    ✕ Discard    ↓ Download    ⟳ Refresh    |    🗑 Delete

Overview        Versions        **Edit**        Generate SAS

| MasterID | ChildID | Name | filenameDate | filenameDateKey |
|----------|---------|-------|--------------|-----------------|
| 10 | 1001 | ItemA | 2019-11-12 | 20191112 |
| 11 | 1002 | ItemB | 2019-11-12 | 20191112 |
| 12 | 1003 | ItemC | 2019-11-12 | 20191112 |
| 13 | 1004 | ItemD | 2019-11-12 | 20191112 |
| 14 | 1005 | ItemE | 2019-11-12 | 20191112 |
| 15 | 1006 | ItemF | 2019-11-12 | 20191112 |
| 16 | 1007 | ItemG | 2019-11-12 | 20191112 |
| 17 | 1008 | ItemH | 2019-11-12 | 20191112 |
| 18 | 1009 | ItemI | 2019-11-12 | 20191112 |
| 19 | 1010 | ItemJ | 2019-11-12 | 20191112 |
| 20 | 1011 | ItemK | 2019-11-12 | 20191112 |
| 21 | 1012 | ItemL | 2019-11-12 | 20191112 |
| 22 | 1013 | ItemM | 2019-11-12 | 20191112 |
| 23 | 1014 | ItemN | 2019-11-12 | 20191112 |
| 24 | 1015 | ItemO | 2019-11-12 | 20191112 |
| 25 | 1016 | ItemP | 2019-11-12 | 20191112 |
| 26 | 1017 | ItemQ | 2019-11-12 | 20191112 |