# Chinook OLTP → Snowflake Data Warehouse Integration using Azure Data Factory

## 1. Objective

The goal of this project is to replicate the *Chinook OLTP to Data Warehouse* data integration workflow using **Azure Data Factory (ADF)**, **Azure SQL Database**, **Azure Blob Storage**, and **Snowflake**.

The deliverable demonstrates an end-to-end ETL process that extracts data from the Chinook transactional database hosted in Azure SQL, stages it in Azure Blob Storage as Parquet files, and then loads it into Snowflake's `STAGE` and `DW` schemas using ADF pipelines and data flows.

## 2. Architecture Overview

The architecture follows a layered ELT approach:

1.  **Extract Layer (ADF Pipeline 1)**

    ○   Source: Azure SQL Database (Chinook OLTP)

    ○   Destination: Azure Blob Storage (Parquet format)

2.  **Stage Layer (ADF Pipeline 2)**

    ○   Reads Parquet data from Blob Storage

    ○   Loads into Snowflake `STAGE` schema tables

3.  **Transformation Layer (ADF Data Flow)**

    ○   Performs transformation logic (hashing, deduplication, merge updates)

    ○   Loads transformed data into Snowflake `DW` schema tables

4.  **Data Warehouse Layer**

    ○   Snowflake hosts final `DATE_DIM`, `TIME_DIM`, `CUSTOMER_DIM`, `ARTIST_DIM`, and `SALES_FACT` tables.

    ○   These form a star schema suitable for business analytics.

# 3. Tools and Technologies Used

| Component | Purpose |
|---|---|
| **Azure SQL Database** | Stores source Chinook OLTP tables (Customer, Invoice, Artist, Album) |
| **Azure Data Factory (ADF)** | Orchestrates extraction, staging, and loading via pipelines and data flows |
| **Azure Blob Storage** | Intermediate storage for Parquet files |
| **Azure Key Vault** | Secures Snowflake and SQL credentials |
| **Snowflake** | Target cloud data warehouse hosting STAGE and DW schemas |
| **DBeaver / Snowflake UI** | Used for running SQL DDL, validation, and testing |

# 4. Azure Setup

**Screenshot:** Setup.png

Created the following resources:

- SQL Server: damg7370fall2025.database.windows.net

- SQL Database: DAMG7370FALL2025

- Storage Account: stgchinookdamg

- Azure Data Factory: adfchinookdamg

- Azure Key Vault: AzureKeyVault1

Networking setup:

- Added client IP address to firewall

- Enabled "Allow Azure services to access server"

- Verified connection via SSMS and ADF

# 5. Datasets

**Screenshots:** `Dataset1.png`, `Dataset2.png`

| Dataset Name | Type | Linked Service | Purpose |
|---|---|---|---|
| `sqlserverdb_chinook` | Azure SQL Database | `Sql_db_Chinook` | Source dataset |
| `Parquet_ds` | Parquet (Blob) | `Storage_Chinook` | Staging Parquet output |
| `SnowChinook_Ds` | Snowflake | `snow_chinook` | Target dataset for loading into STAGE schema |

# 6. Pipeline 1 — Extract SQL DB to Parquet

**Screenshot:** `1stpipeline.png`

**Pipeline Name:** `extract_SQLDB_PL`

- Uses a **ForEach** activity looping through the array:
  `["Customer", "Artist", "Album", "Invoice"]`

- Inside loop: Copy Activity (`sql_2_parquet`)

  - Source: Azure SQL Database

  - Sink: Azure Blob Storage (Parquet)

- Managed identity access configured for Blob Storage.

- All activities succeeded with matching row counts.

## 8. Pipeline 2 — Parquet to Snowflake Stage

**Screenshot:** `2nd pipeline.png`

**Pipeline Name:** `Parquet_2_SnowStage_PL`

- Reads Parquet files from Blob container `/stage_data/`

- Writes into Snowflake `STAGE` tables:

    - `STAGE.CUSTOMER`

    - `STAGE.ARTIST`

    - `STAGE.ALBUM`

    - `STAGE.INVOICE`

- Connected using SAS-based Snowflake Linked Service for stage loading.

## 9. Data Flow — Load Customer Dimension

**Screenshot:** `Dataflow.png`

**Data Flow Name:** `DF_Load_Customer_DIM`

**Steps:**

1. **Source:** `STAGE.CUSTOMER`

2. **Derived Column:** Generate `CUSTOMER_HASH` using SHA-256 for change detection.

3. **Join:** Compare incoming vs existing customers.

4. **Sink:** Upsert into `DW.CUSTOMER_DIM`

    - Merge logic based on hash value.

    - Insert new and update changed records.

# 10. Snowflake Schemas and Tables

**Screenshot:** `Dimtables.png`

Executed SQL scripts in Snowflake:

1. `create_stage_schema.sql`

2. `create_dw_schema.sql`

3. `load_date_dim.sql`

4. `load_time_dim.sql`

5. `merge_artist_dim.sql`

6. `load_sales_fact.sql`

7. `validation_counts.sql`

Final Star Schema Tables:

- **Dimensions:**
  `DATE_DIM`, `TIME_DIM`, `CUSTOMER_DIM`, `ARTIST_DIM`

- **Fact Table:**
  `SALES_FACT` (references DATE_DIM_KEY and CUSTOMER_KEY)

# 11. SQL Script Summary

| Script | Purpose |
|---|---|
| `create_stage_schema.sql` | Defines staging layer tables |
| `create_dw_schema.sql` | Creates DW layer (DIM & FACT) |
| `load_date_dim.sql` | Populates `DATE_DIM` |
| `load_time_dim.sql` | Populates `TIME_DIM` |

`merge_artist_dim.sql`          Incremental load for `ARTIST_DIM`

`load_sales_fact.sql`          Inserts transactional data into `SALES_FACT`

`validation_counts.sql`      Validates record counts across all layers
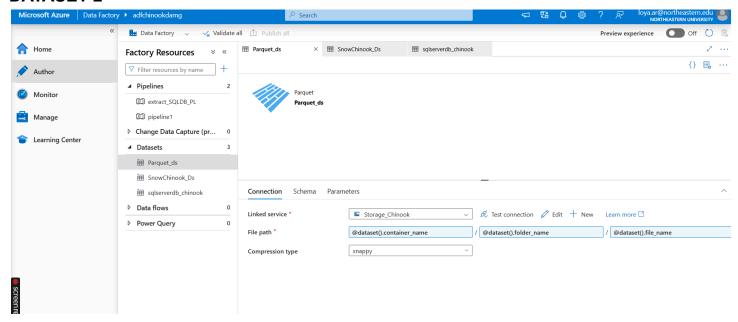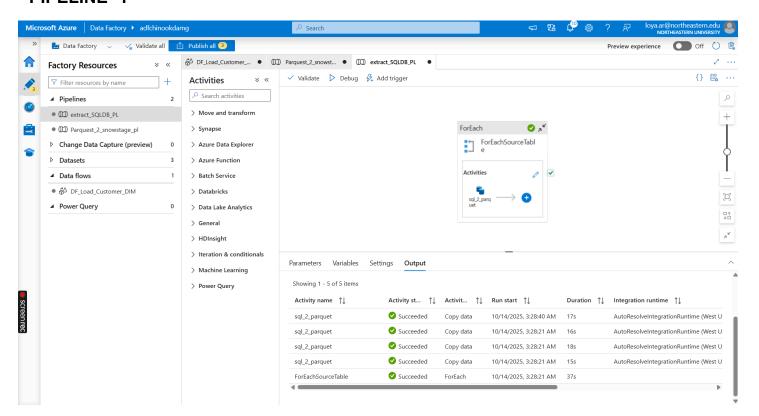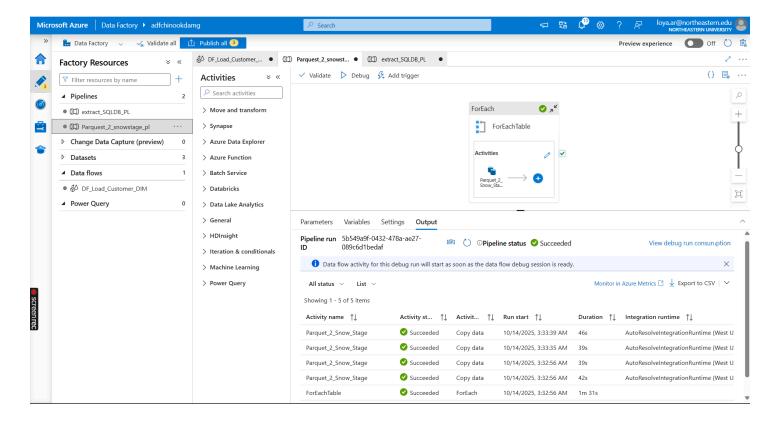
## All Screenshots

### SETUP



### DATA SET 1

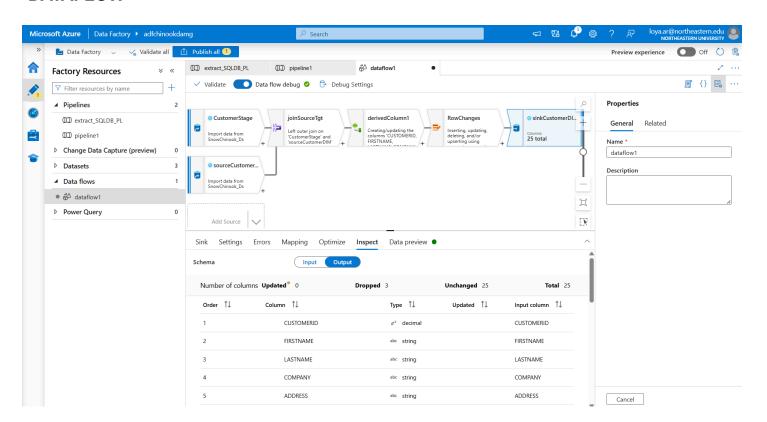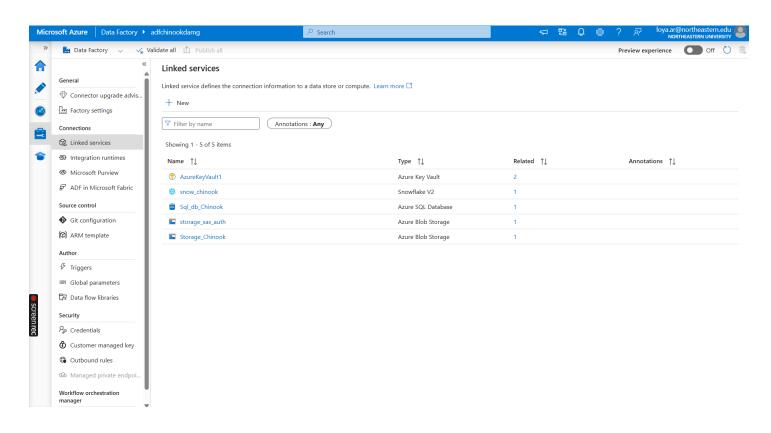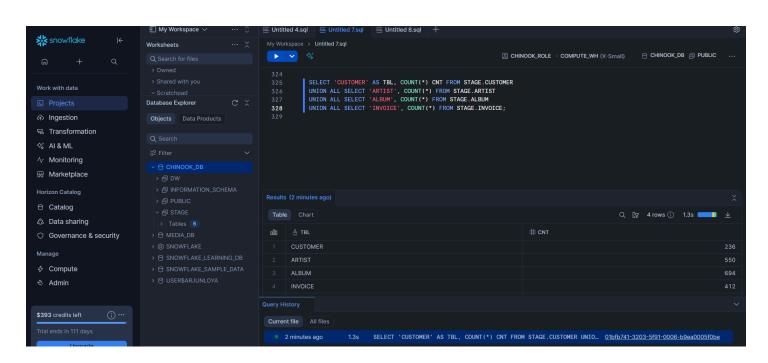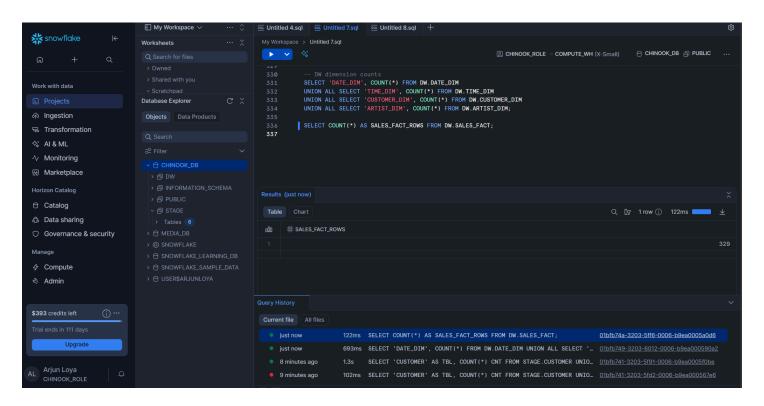# DATASET 2



# PIPELINE- 1



# PIPELINE-2

## DATAFLOW

# LINKED SERVICES



# VALIDATION

# DIMTABLES



# FACT_TABLE_VALIDATION



# THANKYOU