# 01

# Overview of HW 3

HW 3 - Transforming & Staging

# Homework 3

## Transforming & Staging

**01**

SQL Database +
Data Factory

**02**

Transform

Change 'DateA'
to Date

**03**

Staging

Copy Data to
SQL DB Table

# HW3 Final Output

## Most Recently Updated Rows/Observations from the SQL Table Data

# Homework 3

Transforming & Staging

**01** SQL Database

**02** Data Factory

**03** SQL Database

# 0. Deploy hw3 template

**01**

Open Azure CloudShell

**02**

Clone Github Repository
(ALREADY CLONED)

**03**

Create Azure Template

**04**

Deploy Azure Template

*Hint

```
PS /home/dk98> cd './  OMDSMod4  /Homework 3/Transforming and Staging/'
PS /home/dk98/  OMDSMod4  /Homework 3/Transforming and Staging> bash ./formTemplate.sh
Template and parameters created successfully: ./template/template.json and ./template/parameters.json
```

# Deployment Result

| | | |
|---|---|---|
| ☐ 🗄 db47fdf4- | SQL server | East US |
| ☐ 🗄 dba09cc5 | SQL database | East US |

Azure SQL Database                    Azure SQL Server

# 1. Go to Azure SQL Database

| | | |
|---|---|---|
| ☐ 🗄 db47fdf4- [        ] | SQL server | East US |
| ☐ 🗄 dba09cc5 [            ] | SQL database | East US |

Azure SQL Database

Azure SQL Server

slidesmania.com

# 1.1 Access Query Editor in SQL DB

Home > David >

**dba09cc5-** ⬚⬚⬚⬚⬚⬚⬚⬚⬚ | Query editor (preview) ☆ ...
SQL database

🔍 Search

« 👤 Login ＋ New Query ⬆ Open query 🔗 Feedback 📖 Getting started

⬚ Overview

📋 Activity log

🏷 Tags

🔧 Diagnose and solve problems

🔲 Query editor (preview)

**Settings**

⚙ Compute + storage

🔌 Connection strings

📊 Properties

🔒 Locks

**Data management**

🗄 Replicas

🔄 Sync to other databases

Query editor (preview) is a tool to run SQL queries against Azure SQL Database in the Azure portal. It is designed for lightweight querying and object exploration in your database. For more information and troubleshooting, Learn more

**SQL**

Welcome to SQL Database Query Editor

SQL server authentication                    Microsoft Entra authentication
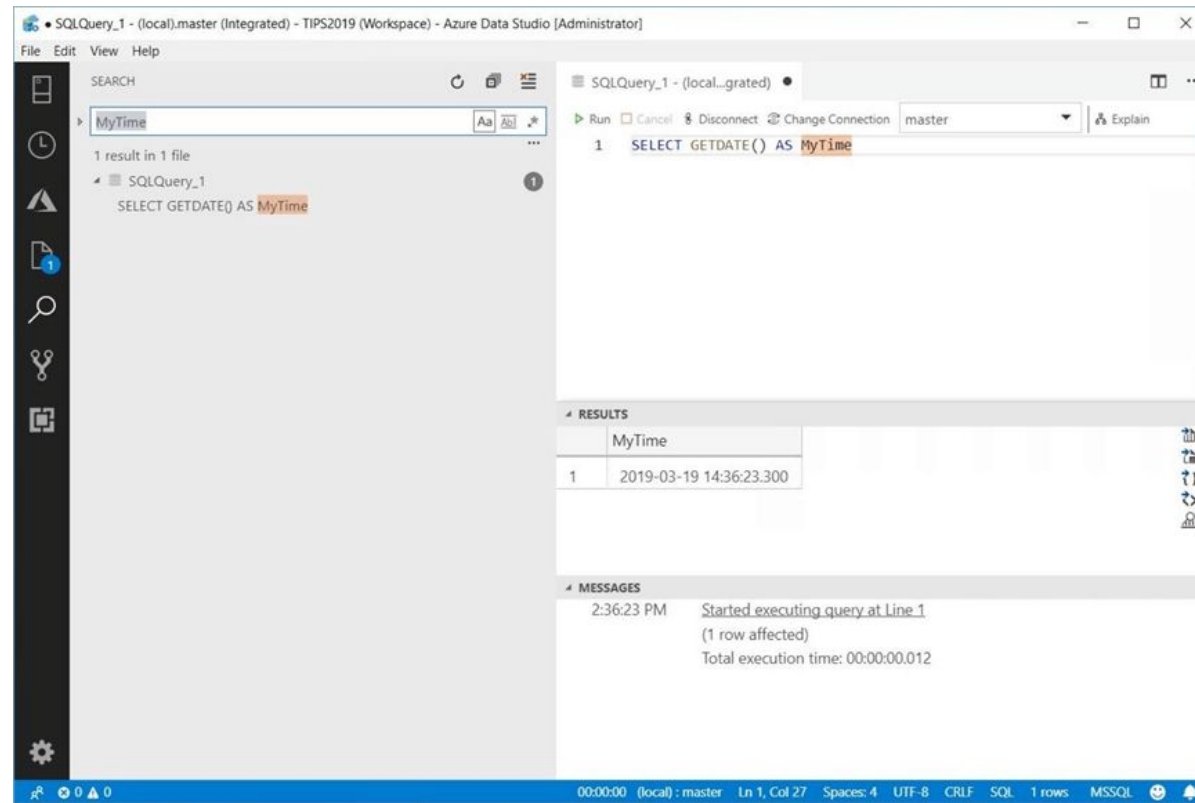
Login *                                      [ Continue as dk98@bu.edu ]
[ omdsmod4admin ]
                                  OR
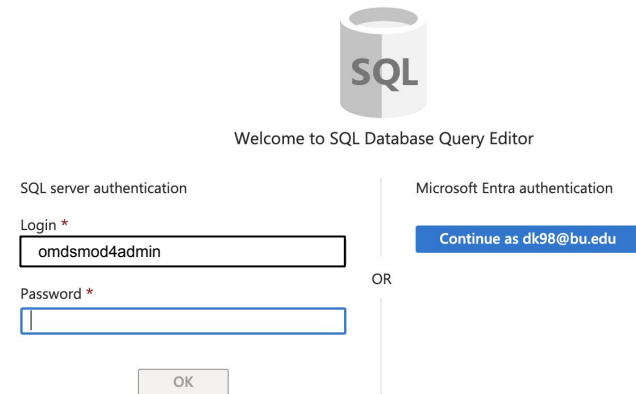Password *
[ omdsmod4password013! ]

[ OK ]

slidesmania.com

# 1.1 Or Azure Data Studio

# 1.2 Authentication

1. Deployed our template:
   a. Login: omdsmod4admin
   b. Password: omdsmod4password013!
2. Created SQL Database Manually:
   a. Own Credentials

**SQL**

Welcome to SQL Database Query Editor

SQL server authentication                    Microsoft Entra authentication

Login *
omdsmod4admin                    Continue as dk98@bu.edu

                                  OR

Password *

            OK

*Firewall Error

If Prompted,

Give Access to your own IP Address

# 1.3 Create Table in SQL Database

Look at [Reference](#) File to create the following query

CREATE TABLE Complaints

(      CMPLID int Primary,

      *Hint: NVARCHAR(max),

      ...

      DATEA **date**,

)

```
Query 1 ×

▷ Run    ☐ Cancel query    ↓ Save query

1   CREATE TABLE complaints
2   (
3       CMPLID int primary key,
```

```
Appendix A.  Complaints File Characteristics

Field#  Name        Type/Size   Description
------  ---------   ---------   ------------------------------------
1       CMPLID      CHAR(9)     NHTSA'S INTERNAL UNIQUE SEQUENCE NUMBER.
                                IS AN UPDATEABLE FIELD,THUS DATA FOR A
                                GIVEN RECORD POTENTIALLY COULD CHANGE FROM
                                ONE DATA OUTPUT FILE TO THE NEXT.
2       ODINO       CHAR(9)     NHTSA'S INTERNAL REFERENCE NUMBER.
                                THIS NUMBER MAY BE REPEATED FOR
                                MULTIPLE COMPONENTS.
                                ALSO, IF LDATE IS PRIOR TO DEC 15, 2002,
                                THIS NUMBER MAY BE REPEATED FOR MULTIPLE
                                PRODUCTS OWNED BY THE SAME COMPLAINANT.
3       MFR_NAME    CHAR(40)    MANUFACTURER'S NAME
4       MAKETXT     CHAR(25)    VEHICLE/EQUIPMENT MAKE
5       MODELTXT    CHAR(256)   VEHICLE/EQUIPMENT MODEL
6       YEARTXT     CHAR(4)     MODEL YEAR, 9999 IF UNKNOWN or N/A
7       CRASH       CHAR(1)     WAS VEHICLE INVOLVED IN A CRASH, 'Y' OR 'N'
8       FAILDATE    CHAR(8)     DATE OF INCIDENT (YYYYMMDD)
9       FIRE        CHAR(1)     WAS VEHICLE INVOLVED IN A FIRE 'Y' OR 'N'
10      INJURED     NUMBER(2)   NUMBER OF PERSONS INJURED
11      DEATHS      NUMBER(2)   NUMBER OF FATALITIES
12      COMPDESC    CHAR(128)   SPECIFIC COMPONENT'S DESCRIPTION
13      CITY        CHAR(30)    CONSUMER'S CITY
14      STATE       CHAR(2)     CONSUMER'S STATE CODE
15      VIN         CHAR(11)    VEHICLE'S VIN#
16      DATEA       CHAR(8)     DATE ADDED TO FILE (YYYYMMDD)
17      LDATE       CHAR(8)     DATE COMPLAINT RECEIVED BY NHTSA (YYYYMMDD)
18      MILES       NUMBER(7)   VEHICLE MILEAGE AT FAILURE
19      OCCURENCES  NUMBER(4)   NUMBER OF OCCURRENCES
20      CDESCR      CHAR(2048)  DESCRIPTION OF THE COMPLAINT
21      CMPL_TYPE   CHAR(4)     SOURCE OF COMPLAINT CODE:
```

```sql
CREATE TABLE complaints (
    CMPLID int primary key,
    ODINO nvarchar(max),
    MFR_NAME nvarchar(max),
    MAKETXT nvarchar(max),
    MODELTXT nvarchar(max),
    YEARTXT nvarchar(max),
    CRASH nvarchar(max),
    FAILDATE date,
    FIRE nvarchar(max),
    INJURED int,
    DEATHS int,
    COMPDESC nvarchar(max),
    CITY nvarchar(max),
    STATE nvarchar(max),
    VIN nvarchar(max),
    DATEA date,
    LDATE date,
    MILES int,
    OCCURENCES int,
    CDESCR nvarchar(max),
    CMPL_TYPE nvarchar(max),
    POLICE_RPT_YN nvarchar(max),
    PURCH_DT date,
    ORIG_OWNER_YN nvarchar(max),
    ANTI_BRAKES_YN nvarchar(max),
    CRUISE_CONT_YN nvarchar(max),
    NUM_CYLS int,
    DRIVE_TRAIN nvarchar(max),
    FUEL_SYS nvarchar(max),
    FUEL_TYPE nvarchar(max),
    TRANS_TYPE nvarchar(max),
    VEH_SPEED int,
    DOT nvarchar(max),
    TIRE_SIZE nvarchar(max),
    LOC_OF_TIRE nvarchar(max),
    TIRE_FAIL_TYPE nvarchar(max),
    ORIG_EQUIP_YN nvarchar(max),
    MANUF_DT date,
    SEAT_TYPE nvarchar(max),
    RESTRAINT_TYPE nvarchar(max),
    DEALER_NAME nvarchar(max),
    DEALER_TEL nvarchar(max),
    DEALER_CITY nvarchar(max),
    DEALER_STATE nvarchar(max),
    DEALER_ZIP nvarchar(max),
    PROD_TYPE nvarchar(max),
    REPAIRED_YN nvarchar(max),
    MEDICAL_ATTN nvarchar(max),
    VEHICLES_TOWED_YN nvarchar(max)
)
```

The decision to utilize nvarchar(max) instead of adhering to the reference file is due to the presence of human error in the data.

slidesmania.com

# 1.4 Check your 'complaints' Table

# 2. Data Factory

**01**

## COPY DATA 1

Transform + Stage

# Copy Data 1 – Source

Azure Blob Storage

(Output from HW 2)

complaints.txt

1.2GB

DelimitedText

**Source**

Azure Blob Storage

Delimited Text

complaints.txt

*Import Schema - 49 Columns

slidesmania.com

# Delimited Text? Delimiter?

# Copy Data 1 - Sink

Azure SQL Database

Complaints Table

*Authentication:
Back to Slide 'Authentication'

**Sink**

SQL Database

Table dbo.complaints

*Import Schema - 49 Columns
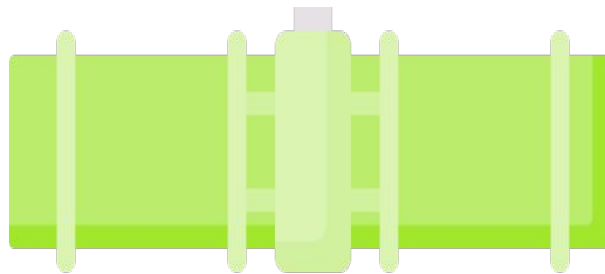
# Copy Data Overview

## Source

(Output from HW 2)

Azure Blob Storage

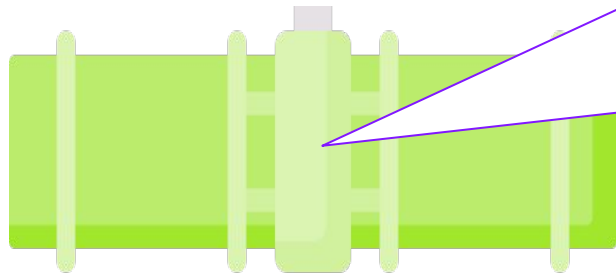Delimited Text

.txt

## 01

**COPY DATA 1**

Transform + Stage

## Sink

SQL Database

Table dbo.Complaints

# Transform

**01**

## COPY DATA 1

Transform

**<Mapping>**

**'DateA' Column**

From 'String' to 'Date' Type

**\*Type Conversion Setting:**

- Datetime Format: yyyyMMdd

\* From <u>Reference</u> File

| 16 | DATEA | CHAR(8) | DATE ADDED TO FILE (YYYYMMDD) |

# IMPORTANT!

## Create Table

### CMPLID INT PRIMARY

- Optional: Add '**NOT NULL**'
- Do not add **IDENTITY(1,1)**

## Copy Data 1 'Settings'

\*Navigate to **'Fault Tolerance'**

- ✅ **Skip Incompatible Rows**
- ☒ Logging
- This is to skip rows affected by **human errors**

## Running Pipeline

- Pipeline runs on **Azure Cloud**
  - Turning off your computer will not affect the run
- Approx. **2.0 hours** for completion
- **Monitor** to ensure the data is being read & written properly

# *RUNNING A PIPELINE

- **Publish & Add Trigger** >> Cheaper than 'Debug' Action

- Monitor pipeline runs through the **'monitor'** section

  - Go to details (eyeglasses icon) for details on data read & written)
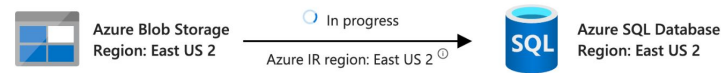
## Left Panel

Details  ↻ Refresh

🚀 **Performance tuning tips:**
Writing to sink is slow, as 11 rows are found incompatible with source format and skipped. To achieve better performance, you are suggested to make sure that source and sink have compatible format. Refer to this document .

Learn more on copy performance details from here.

Activity run id: 88ce5804-420e-4ea2-87a9-44248f283a87

**Azure Blob Storage** Region: East US 2  — In progress →  Azure IR region: East US 2  **Azure SQL Database** Region: East US 2

| | | | | |
|---|---|---|---|---|
| Data read: | 1.347 GB | | Data written: | 897.177 MB |
| Files read: | 1 | | Rows written: | 829,530 |
| Rows read: | 1,960,164 | | Rows skipped: | 12 |
| Peak connections: | 10 | | Peak connections: | 2 |

Copy duration  00:47:54
Throughput:  470.162 KB/s

⌄ Azure Blob Storage → Azure SQL Database

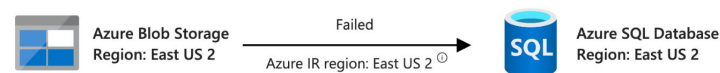| | Working duration | Total duration |
|---|---|---|
| Start time | 2024-02-08T19:48:15.0101701Z | |
| Used DIUs | 4 | |
| Used parallel copies | 1 | |
| ⌄ Duration | 00:47:54 | |
| **Details** | | |
| ✅ Queue | | 00:00:07 |
| ⊙ Transfer | | 00:47:45 |
| Listing source | 00:00:00 | |
| Reading from source | 00:00:09 | |
| Writing to sink | 00:47:34 | |

⌄ Data consistency verification  Verified
Verification result  100% of data copied was verified to be consistent

## Right Panel

Details  ↻ Refresh

🚀 **Performance tuning tips:**
Writing to sink is slow, as 29 rows are found incompatible with source format and skipped. To achieve better performance, you are suggested to make sure that source and sink have compatible format. Refer to this document .

Learn more on copy performance details from here.

Activity run id: 88ce5804-420e-4ea2-87a9-44248f283a87

**Azure Blob Storage** Region: East US 2  — Failed →  Azure IR region: East US 2  **Azure SQL Database** Region: East US 2

| | | | | |
|---|---|---|---|---|
| Data read: | 1.347 GB | | Data written: | 1.71 GB |
| Files read: | 1 | | Rows written: | 1,393,441 |
| Rows read: | 1,960,164 | | Rows skipped: | 30 |
| Peak connections: | 10 | | Peak connections: | 2 |

Copy duration  01:33:55
Throughput:  239.427 KB/s

⌄ Azure Blob Storage → Azure SQL Database

| | Working duration | Total duration |
|---|---|---|
| Start time | 2/8/2024, 2:48:15 PM | |
| Used DIUs | 4 | |
| Used parallel copies | 1 | |
| ⌄ Duration | 01:33:55 | |
| **Details** | | |
| ✅ Queue | | 00:00:07 |
| ⊙ Transfer | | 01:33:46 |
| Listing source | 00:00:00 | |
| Reading from source | 00:00:09 | |
| Writing to sink | 01:33:42 | |

⌄ Data consistency verification  Verified
Verification result  100% of data copied was verified to be consistent

# 3. Go back to Azure SQL Database

| | | |
|---|---|---|
| ☐ 🗄 db47fdf4-[_____] | SQL server | East US |
| ☐ 🗄 dba09cc5-[_____] | SQL database | East US |

Azure SQL Database

Azure SQL Server

# SQL Queries

## Final Query:

SELECT *

FROM Complaints

WHERE DATEA = CONVERT(Date, GETDATE()-1)

- ^^ Get Data From Yesterday?
- But... No Data From Yesterday?
- **Change the 1** in GETDATE()-1 to find the most recent date in your data

## Find Today's Date

## (Azure Time Zone)

SELECT GETDATE()

- Azure's Time Zone is not EST

## Find Dates in Data:

SELECT DISTINCT CONVERT(VARCHAR, DATEA)

FROM Complaints

ORDER BY CONVERT(VARCHAR, DATEA) DESC

- List the dates in 'DateA' column
- Descending order (Most recent on Top)

# HW3 Screenshot
## Transforming & Staging

Best of Luck!