# Sri Lanka Institute of Information Technology



# Data Warehousing & Business Intelligence

# Assignment 1

# 2022

## M.A.D.G.A. SURIYAWATTA

## IT20135652

## TABLE OF CONTENTS

This data set about investigation data for aviation accidents and incidents from 2002 to 2007. An occurrence associated with the operation of an aircraft, which takes place from the time any person boards the aircraft with the intention of flight until all such persons have disembarked, and in which a) a person is fatally or seriously injured, b) the aircraft sustains significant damage or structural failure, or c) the aircraft goes missing or becomes completely inaccessible. This data set defines an aviation incident as an occurrence, other than an accident, associated with the operation of an aircraft that affects or could affect the safety of operation. Accidents and incidents are investigated by government bodies such as the FAA and NTSB.
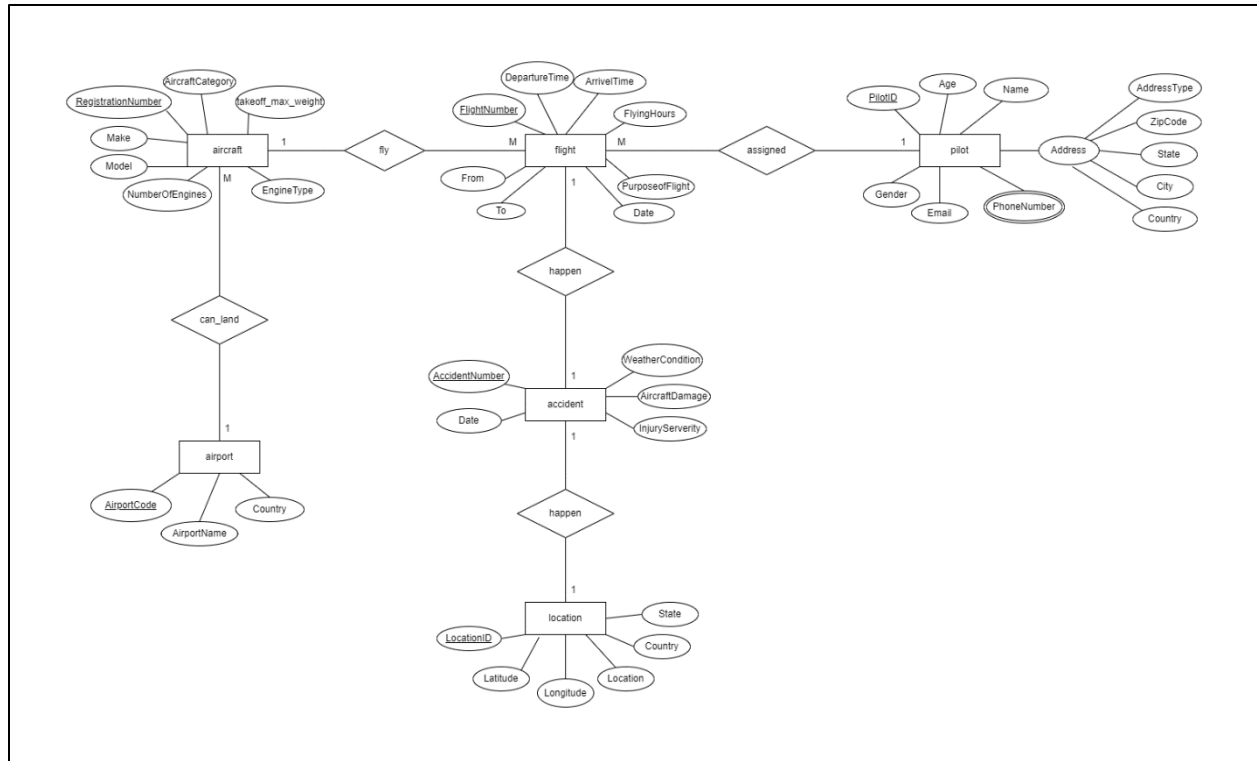
**Link to the selected source data set:**

https://www.kaggle.com/datasets/prathamsharma123/aviation-accidents-and-incidents-ntsb-faa-waas

The original dataset has less tables. I cut the columns of original source tables and put them into different source tables to get more dimensions and a hierarchy, because the assignment document says that we need to enrich the ETL process.

Customized source has seven tables and it include accidents' details, locations' details, aircrafts' details, airports' details, pilot's details, pilots' addresses' details.


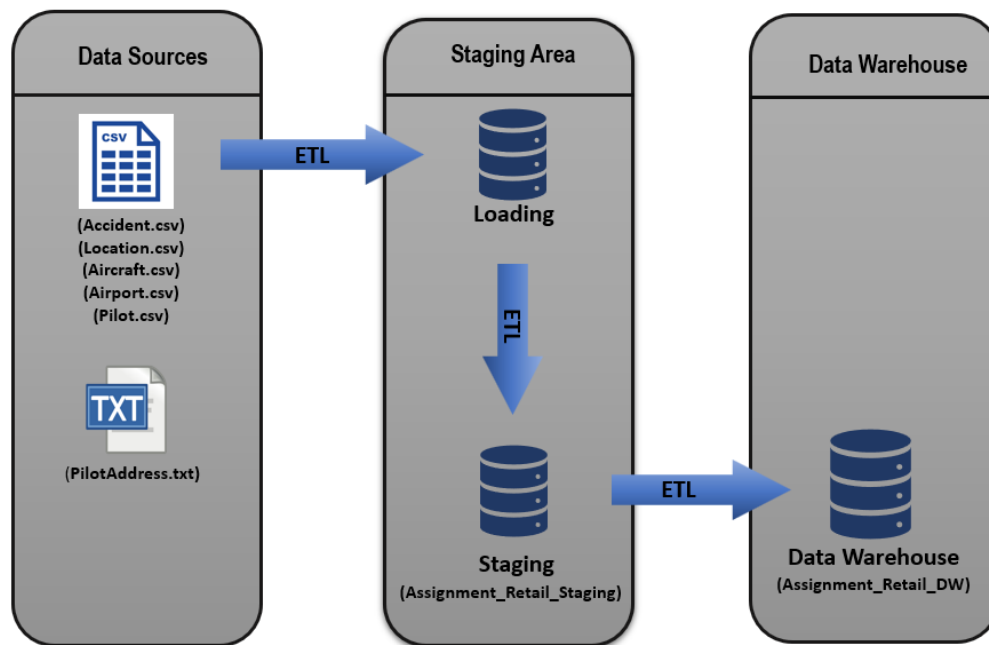## The ER Diagram of the Data Set

## STEP 2: PREPARATION OF DATA SOURCES

There are seven source tables in two formats (.csv & .txt). And they were used to create the following,

| Data Source Name | Data Source Type | Description |
| --- | --- | --- |
| **Accident** | CSV | Details about accidents |
| **Location** | CSV | Details about accidents happened locations |
| **Flight** | CSV | Details about the flights involved the accidents |
| **Aircraft** | CSV | Details about the aircrafts involved the accidents |
| **Airport** | CSV | Details about the airport where the aircrafts land |
| **Pilot** | CSV | Details about the pilots who involved the flights |

| Pilot Addresses | Text | Details abouts pilots' addresses | |
|---|---|---|---|

Above architecture shows the high-level BI solution to the warehouse design.

**Data Sources**

csv' component is used to display Comma Separated files and '.txt' component represents Text files.
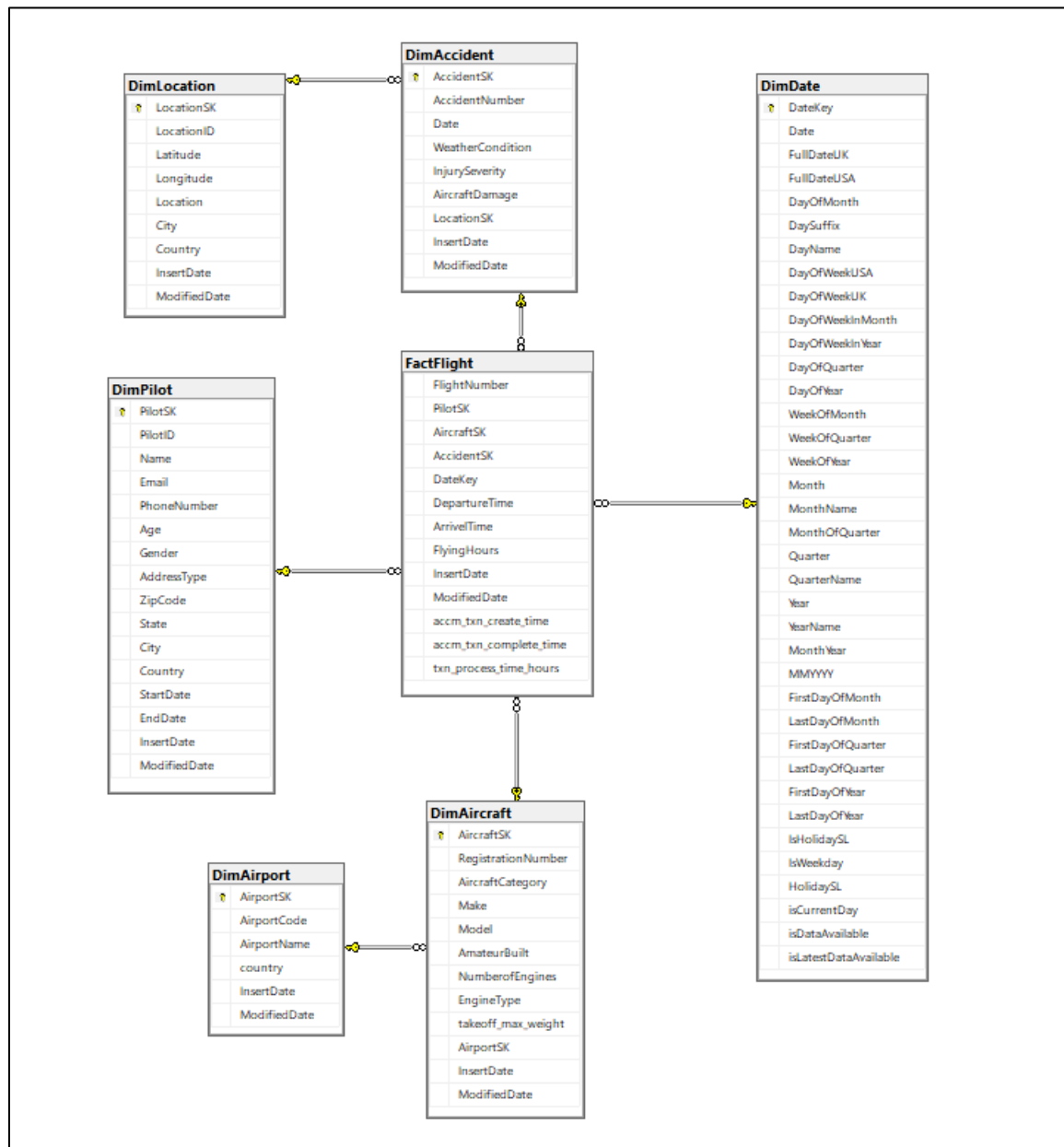
**Staging Area**

This area Loading DB component represents the process of the creating database tables. Accident, Location, Aircraft, Airport, Flight, Pilot and Pilot Addresses text files was imported to the database and was used to create the tables. And these tables were used as the DB source data. Staging DB component represents creating staging level tables through the 'Extract'.

**Data Warehouse**

Data warehouse DB component is used display the cratering dimension tables in the warehouse using 'Transform' and 'Load.'

## STEP 4: DATA WAREHOUSE DESIGN & DEVELOPMENT

Following figure will show how the fact table and dimension tables was combined in a rational manner. For this scenario, **snowflake schema** type was used.

## Dimension Types

**<u>Hierarchical Dimension</u>**

- Date – all the hierarchies in date
- Pilot Addresses – Country -> City -> State -> ZipCode
- Location -  Location -> City -> Country

**<u>Slowly Changing Dimension</u>**

DimPilot is slowly changing dimention. PhoneNumber may be changed in future. Therefore, I get it as slowly changing Attribute.

## 1.Extract

In this step, All the data sources were imported to the staging tables by using the relevant Data connection. Flat file connection was used for text files and csv files. All those tables were imported to the Assignment_Retail_Satging.

- ▪ **Snapshot of SSMS Staging Database**

- **Snapshot of Visual Studio Control Flow of Extract**

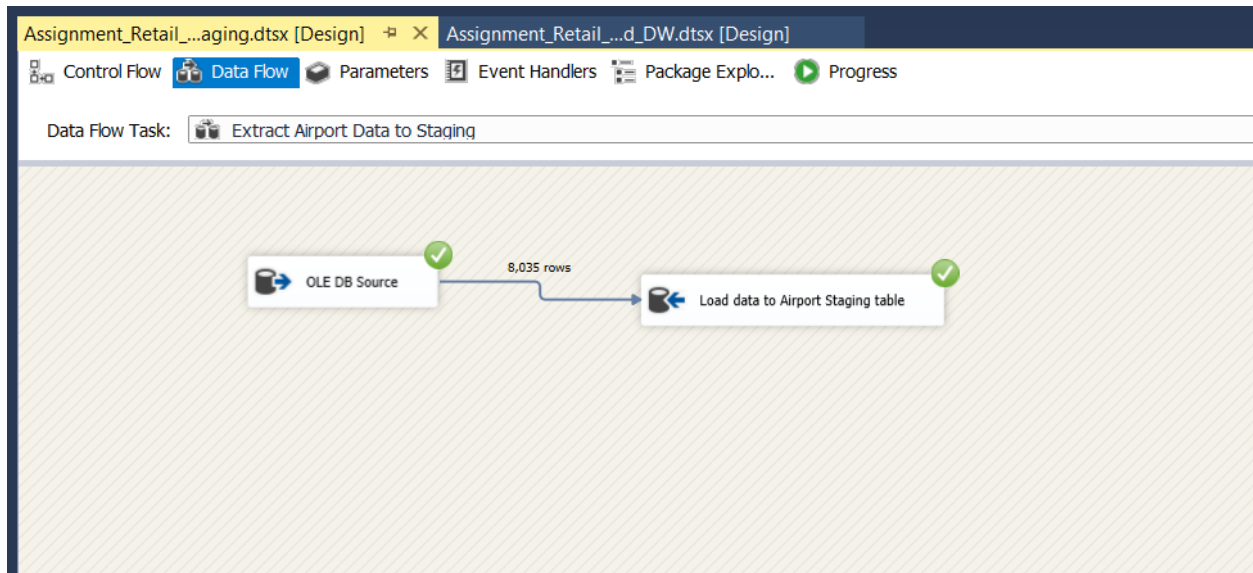Control Flow | Data Flow | Parameters | Event Handlers | Package Explo... | Progress
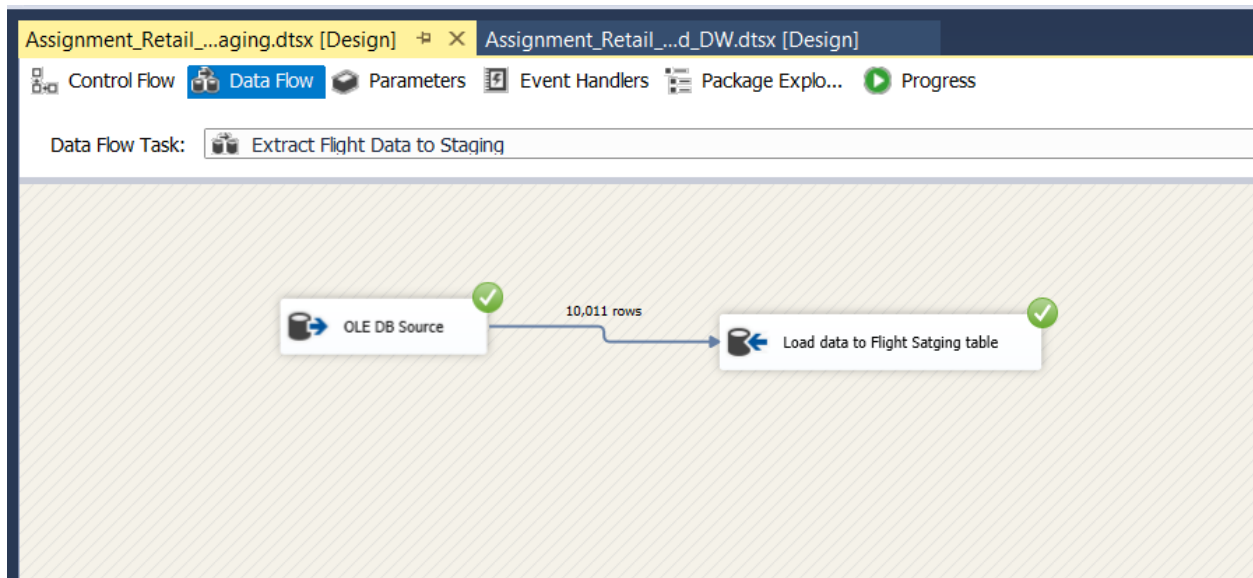
Extract Pilot Data to Staging

Extract Aircraft Data to Staging
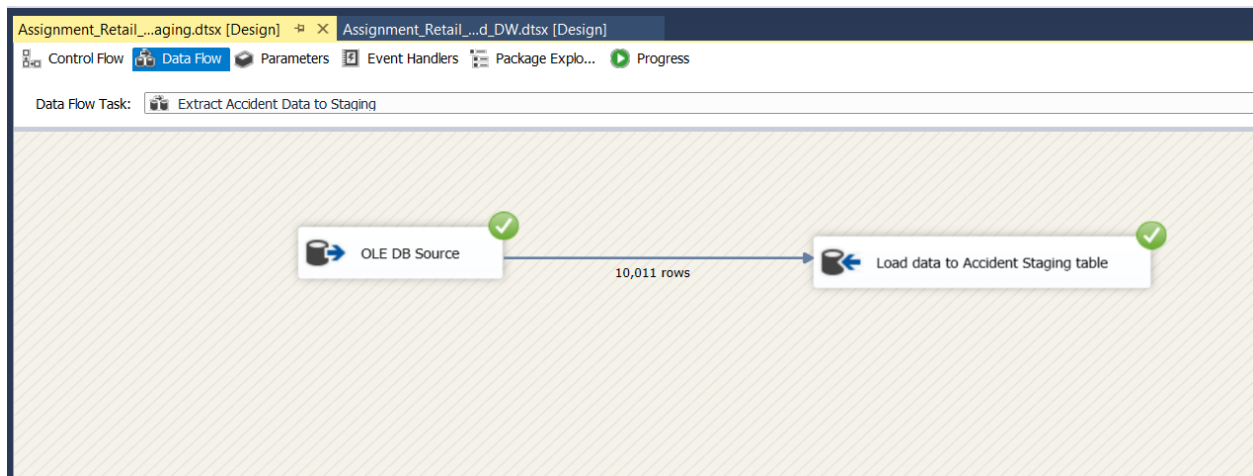
Extract Airport Data to Staging
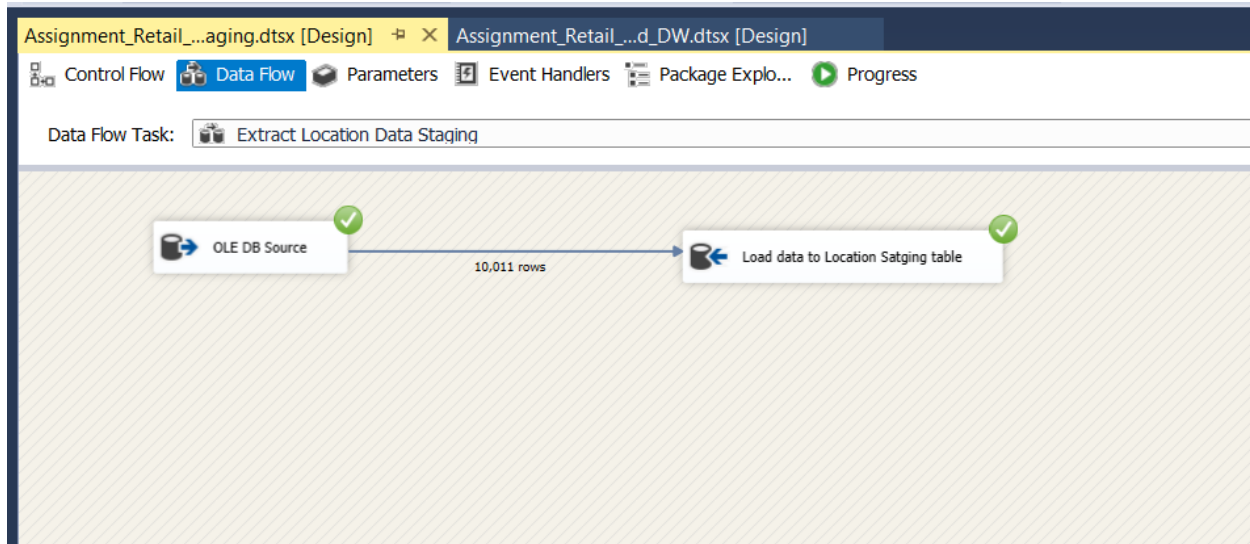
Extract Flight Data to Staging
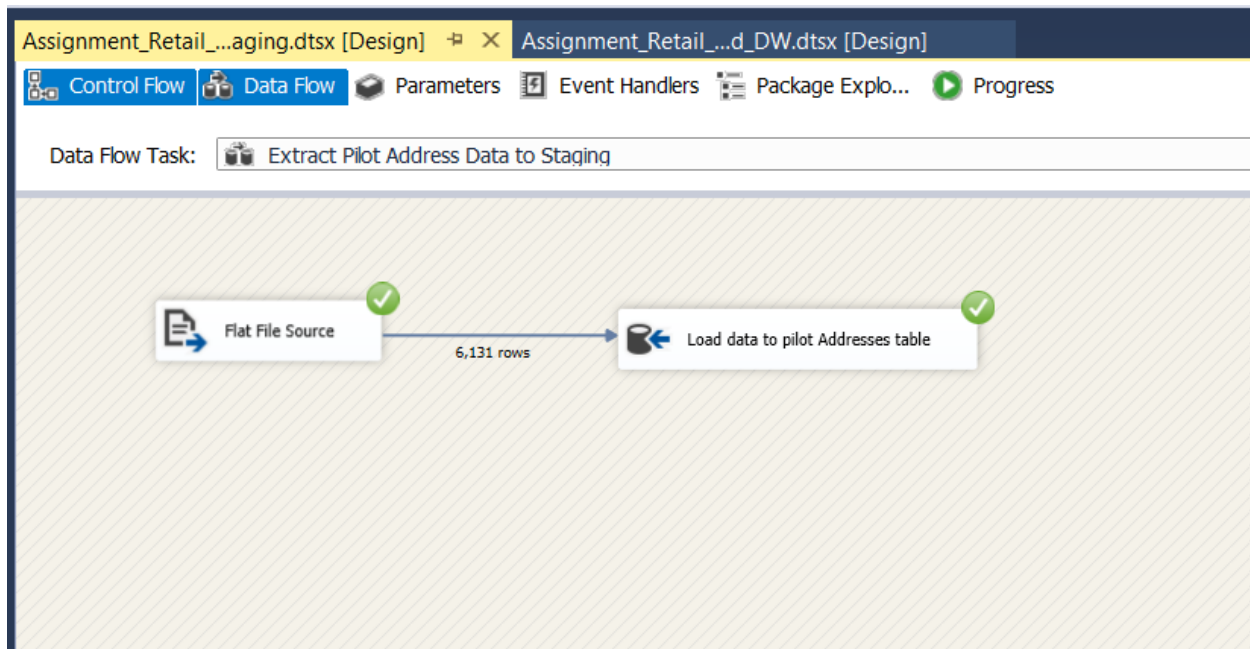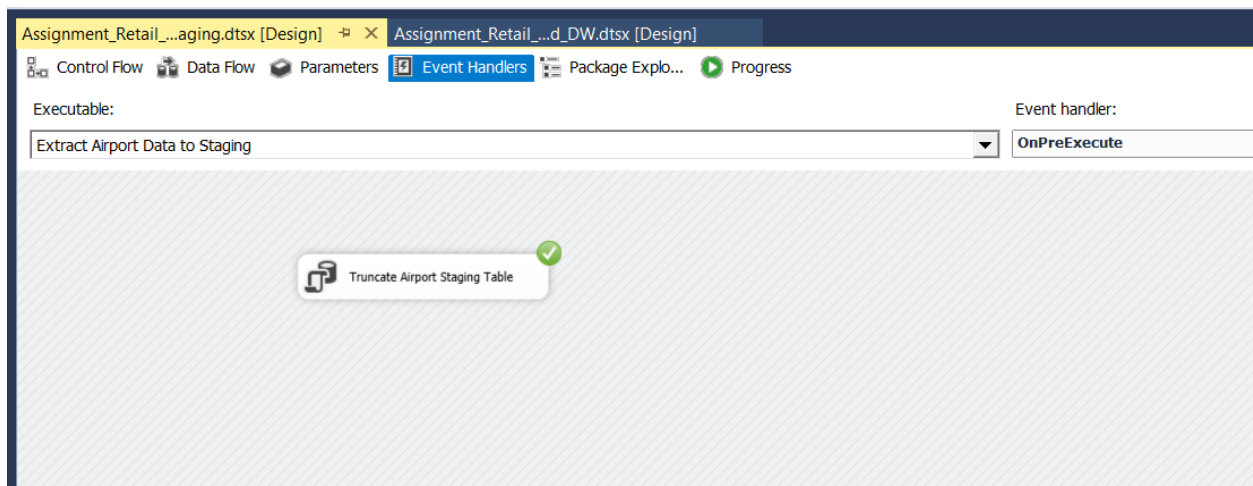
Extract Accident Data to Staging

Extract Location Data Staging

Extract Pilot Address Data to Staging

Extract Complete Time Data to Staging

Execute Assignment_Retail_Load_DW

- **Snapshots of several data types of Data Flows**

- Extract Pilot Data to staging

Control Flow | Data Flow | Parameters | Event Handlers | Package Explo... | Progress

Data Flow Task: Extract Pilot Data to Staging

OLE DB Source → 6,132 rows → Load data to Pilot Staging table

- Extract Aircraft Data to Staging



- Extract Airport Data to staging
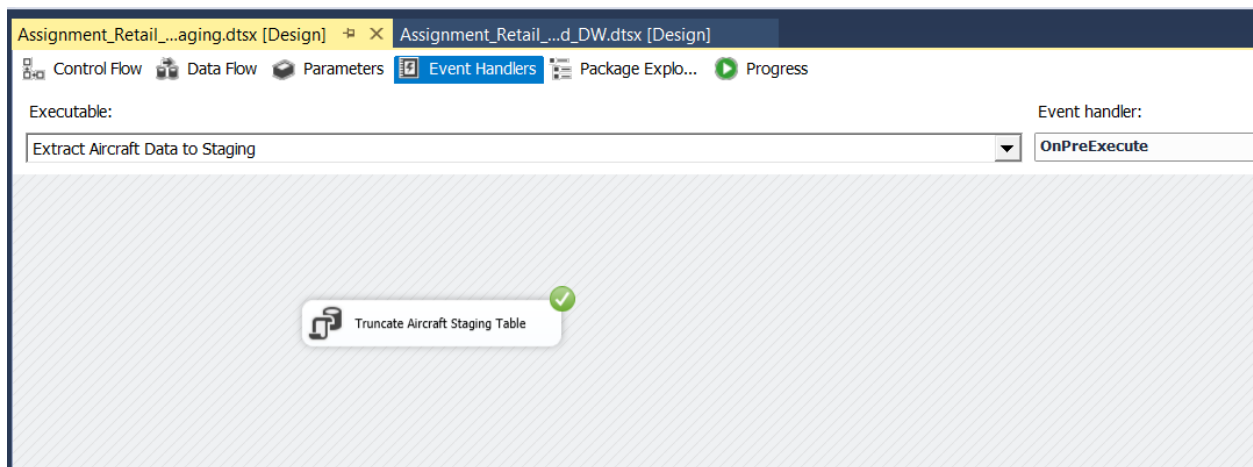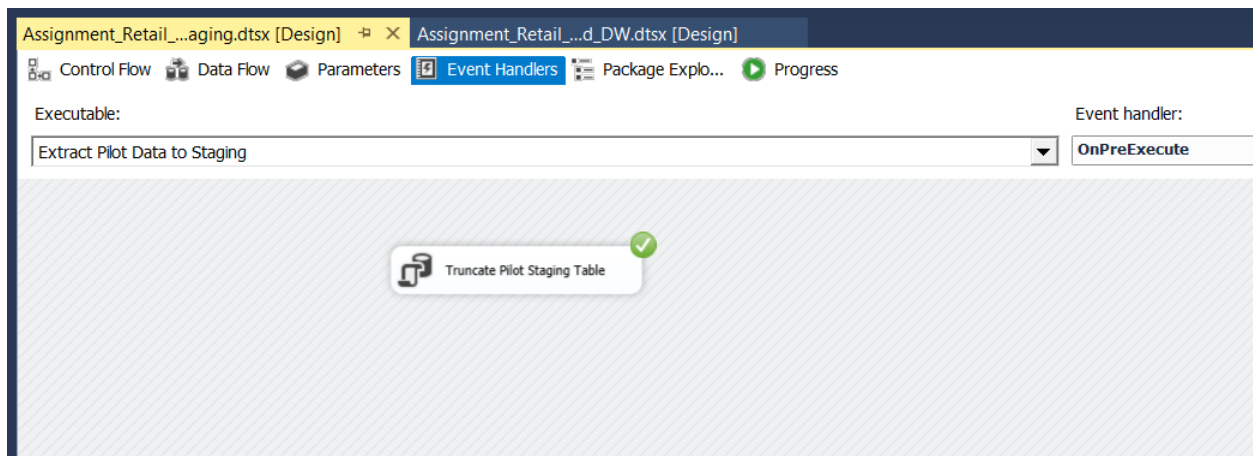
- Extract Flight Data to Staging



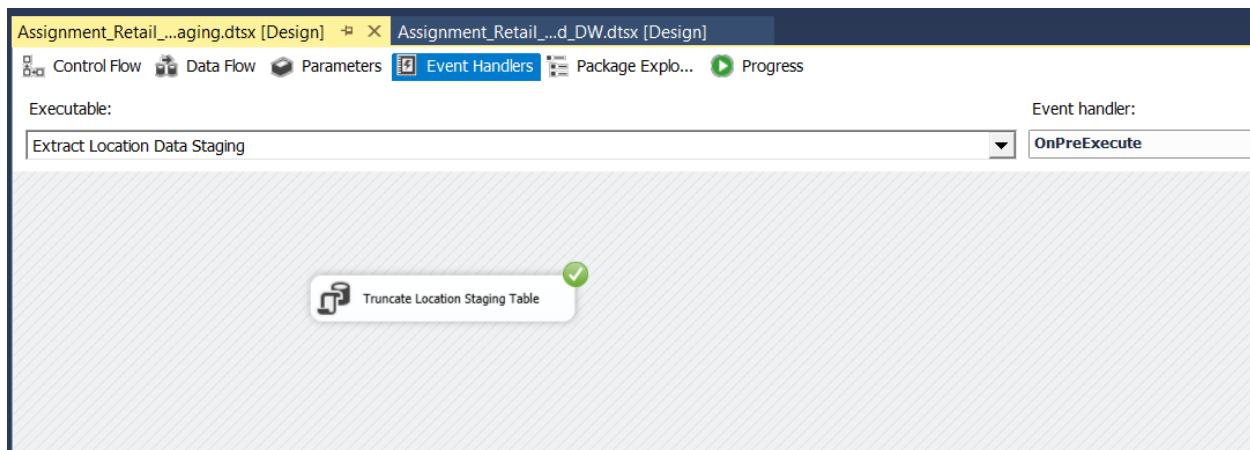- Extract Accident Data to Staging

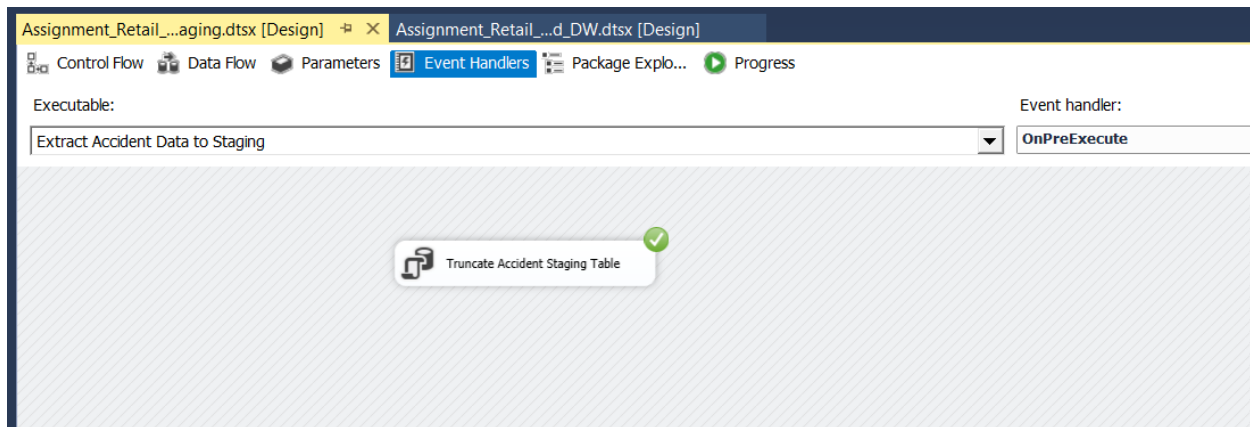- Extract Location Data to Staging
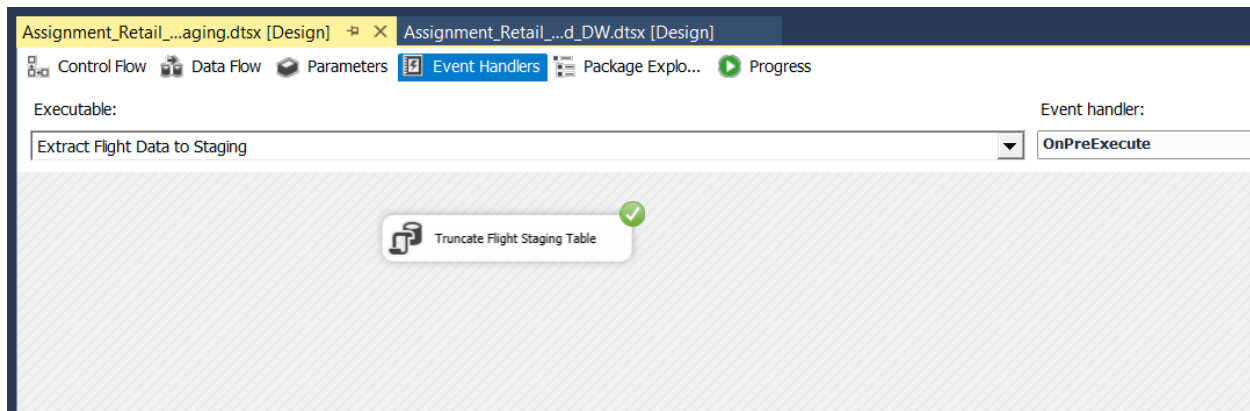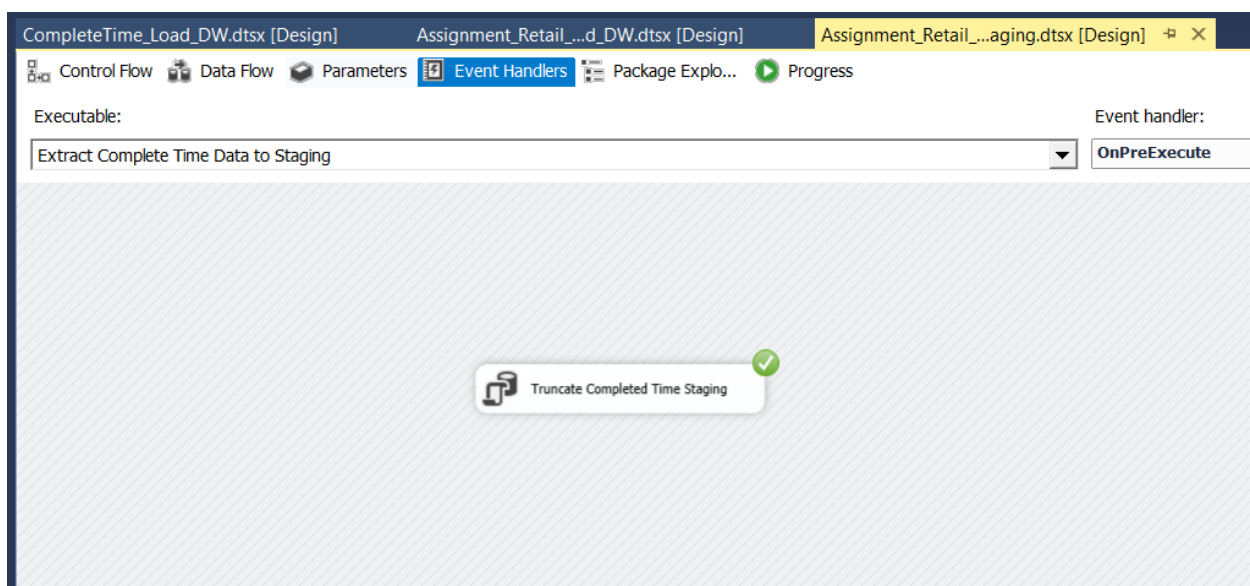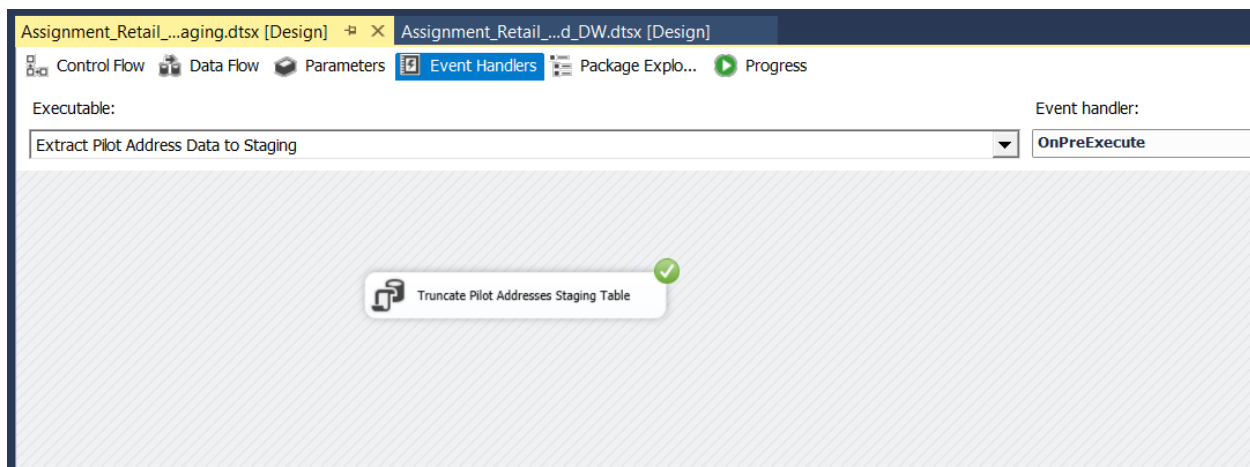


- Extract Pilot Addresses Data to Staging

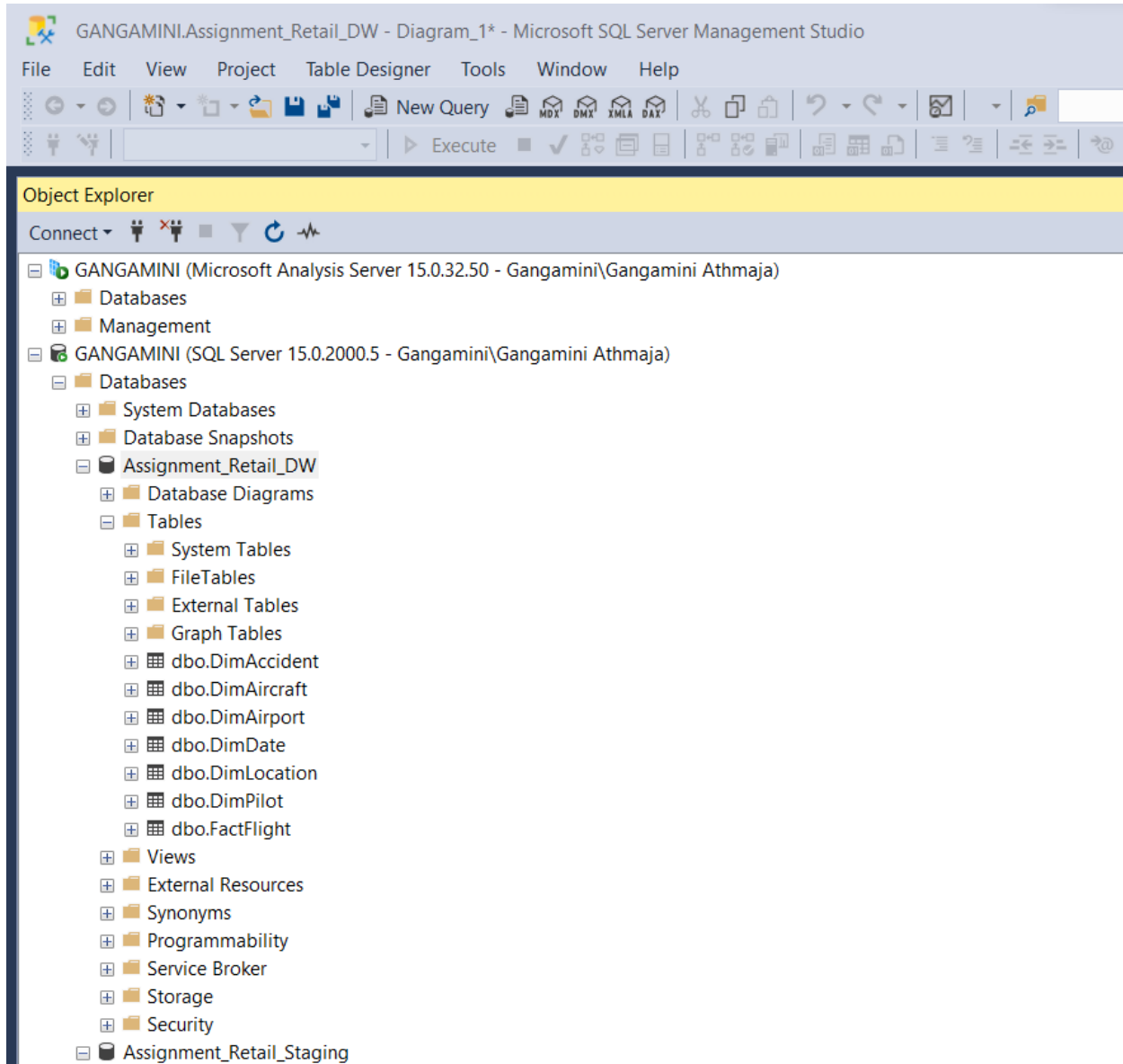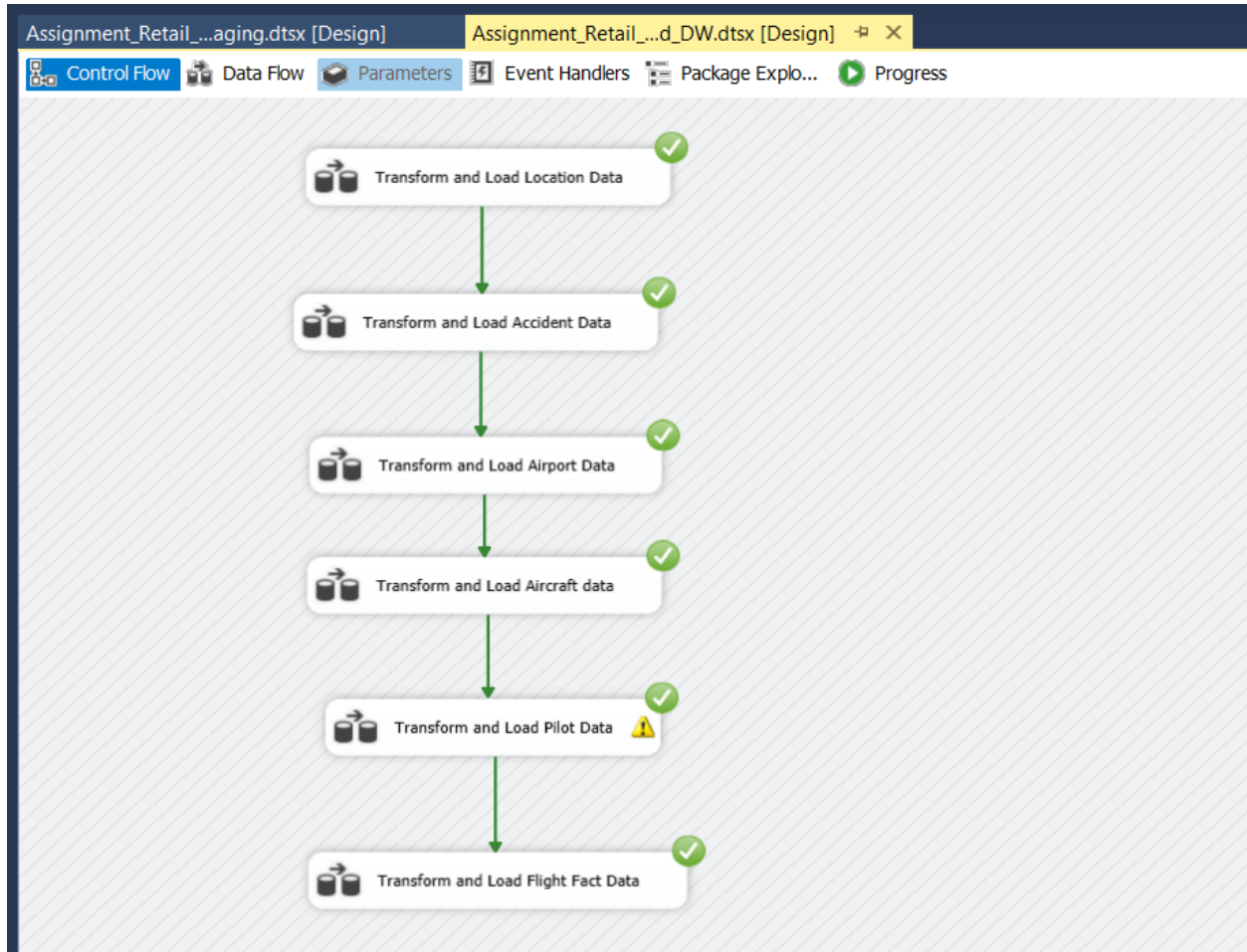## 2. Event Handling (Truncate Staging Data)

## 3.Transform & Load

In this step, both the 'Transform' and 'Load' are done. Firstly, The Dimension tables in the Datawarehouse DB data were created. Then, using the relevant components, data from the staging tables was loaded into the warehouse tables, Assignment_Retail_DW, which contains the below tables,

- **Snapshot of SQL server Data warehouse Database**

- **Snapshot of Visual Studio Control Flow of Extraction**



## Stored Procedures

- Location

```sql
USE [Assignment_Retail_DW]
GO
/****** Object:  StoredProcedure [dbo].[UpdateDimLocation]    Script Date: 5/15/2022 1:48:40 PM ******/
SET ANSI_NULLS ON
GO
SET QUOTED_IDENTIFIER ON
GO
ALTER PROCEDURE [dbo].[UpdateDimLocation]
@LocationID varchar(50),
@Latitude varchar(20),
@Longitude varchar(20),
@Location varchar(50),
@City varchar(50),
@Country varchar(50)

AS
BEGIN
if not exists (select LocationSK
from dbo.DimLocation
where LocationID = @LocationID)
BEGIN
insert into dbo.DimLocation
(LocationID, Latitude, Longitude, Location, City, Country, InsertDate, ModifiedDate)
values
(@LocationID, @Latitude, @Longitude, @Location, @City, @Country, GETDATE(), GETDATE())
END;
if exists (select LocationSK
from dbo.DimLocation
where LocationID = @LocationID)
BEGIN
update dbo.DimLocation
set Latitude = @Latitude,
Longitude = @Longitude,
Location = @Location,
City = @City,
Country = @Country,
ModifiedDate = GETDATE()
where LocationID = @LocationID
END;
END;
```

- Accident

```sql
USE [Assignment_Retail_DW]
GO
/****** Object:  StoredProcedure [dbo].[UpdateDimAccident]    Script Date: 5/15/2022 1:42:51 PM ******/
SET ANSI_NULLS ON
GO
SET QUOTED_IDENTIFIER ON
GO
ALTER PROCEDURE [dbo].[UpdateDimAccident]
@AccidentNumber varchar(30),
@Date datetime,
@WeatherCondition varchar(50),
@InjurySeverity varchar(50),
@AircraftDamage varchar(50),
@LocationSK int

AS
BEGIN
if not exists (select AccidentSK
from dbo.DimAccident
where AccidentNumber = @AccidentNumber)
BEGIN
insert into dbo.DimAccident
(AccidentNumber, Date, WeatherCondition, InjurySeverity, AircraftDamage, LocationSK, InsertDate,  ModifiedDate)
values
(@AccidentNumber, @Date, @WeatherCondition, @InjurySeverity, @AircraftDamage, @LocationSK, GETDATE(), GETDATE())
END;
if exists (select AccidentSK
from dbo.DimAccident
where AccidentNumber = @AccidentNumber)
BEGIN
update dbo.DimAccident
set Date = @Date,
WeatherCondition = @WeatherCondition,
InjurySeverity = @InjurySeverity,
AircraftDamage = @AircraftDamage,
LocationSK = @LocationSK,
ModifiedDate = GETDATE()
where AccidentNumber = @AccidentNumber
END;
END;
```

- Airport

```
USE [Assignment_Retail_DW]
GO
/****** Object:  StoredProcedure [dbo].[UpdateDimAirport]    Script Date: 5/15/2022 1:47:42 PM ******/
SET ANSI_NULLS ON
GO
SET QUOTED_IDENTIFIER ON
GO
ALTER PROCEDURE [dbo].[UpdateDimAirport]
@AirportCode varchar(20),
@AirportName varchar(70),
@country varchar(50)

AS
BEGIN
if not exists (select AirportSK
from dbo.DimAirport
where AirportCode = @AirportCode)
BEGIN
insert into dbo.DimAirport
(AirportCode, AirportName, country, InsertDate, ModifiedDate)
values
(@AirportCode, @AirportName, @country, GETDATE(), GETDATE())
END;
if exists (select AirportSK
from dbo.DimAirport
where AirportCode = @AirportCode)
BEGIN
update dbo.DimAirport
set AirportName = @AirportName,
country = @country,
ModifiedDate = GETDATE()
where AirportCode = @AirportCode
END;
END;
```

- Aircraft



```sql
SQLQuery3.sql - GA...mini Athmaja (67))*    SQLQuery2.sql - GA...mini Athmaja (64))    SQLQuery1.sql - GA...mini Athmaja (54))    GANGAMINI.Assign...il_DW - Diagram_1*
ALTER PROCEDURE [dbo].[UpdateDimAircraft]
    @RegistrationNumber varchar(50),
    @AircraftCategory varchar(50),
    @Make varchar(50),
    @Model varchar(50),
    @AmateurBuilt varchar(50),
    @NumberofEngines varchar(50),
    @EngineType varchar(50),
    @takeoff_max_weight varchar(50),
    @AirportSK int

    AS
BEGIN
if not exists (select AircraftSK
    from dbo.DimAircraft
    where RegistrationNumber = @RegistrationNumber)
BEGIN
insert into dbo.DimAircraft
    (RegistrationNumber, AircraftCategory, Make, Model, AmateurBuilt, NumberofEngines, EngineType, takeoff_max_weight, AirportSK, InsertDate, ModifiedDate)
    values
    (@RegistrationNumber, @AircraftCategory, @Make, @Model, @AmateurBuilt, @NumberofEngines, @EngineType, @takeoff_max_weight, @AirportSK, GETDATE(), GETDATE())
END;
if exists (select AircraftSK
    from dbo.DimAircraft
    where RegistrationNumber = @RegistrationNumber)
BEGIN
update dbo.DimAircraft
    set AircraftCategory = @AircraftCategory,
    Make = @Make,
    Model = @Model,
    AmateurBuilt = @AmateurBuilt,
    NumberofEngines = @NumberofEngines,
    EngineType = @EngineType,
    takeoff_max_weight = @takeoff_max_weight,
    AirportSK = @AirportSK,
    ModifiedDate = GETDATE()
    where RegistrationNumber = @RegistrationNumber
END;
END;
```
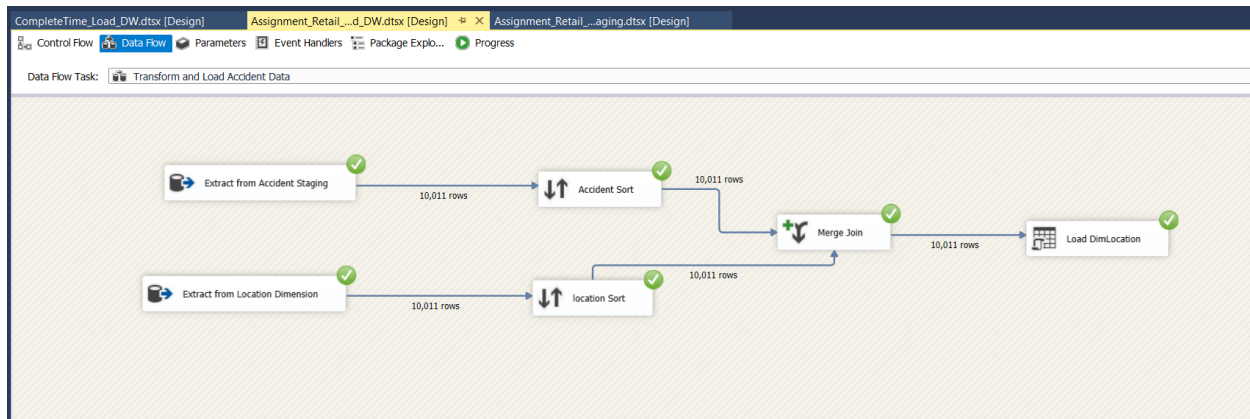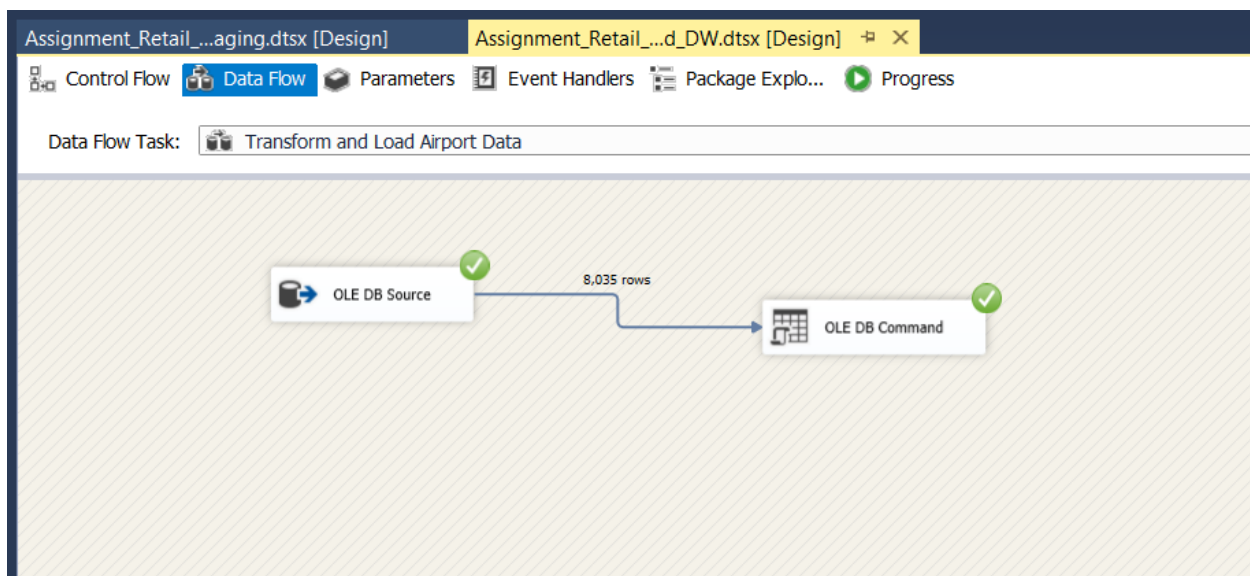
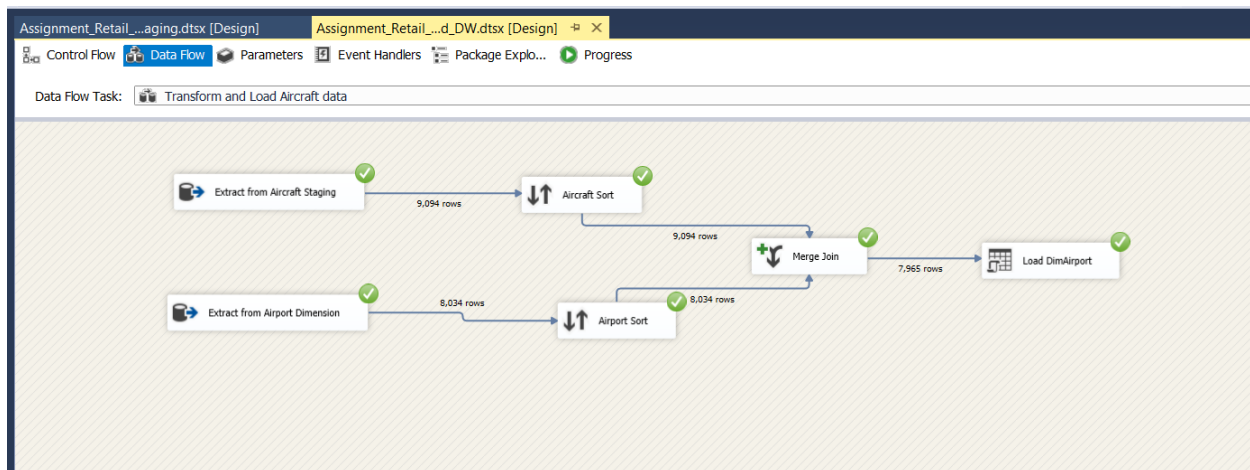- **Transform and Load Location Data**
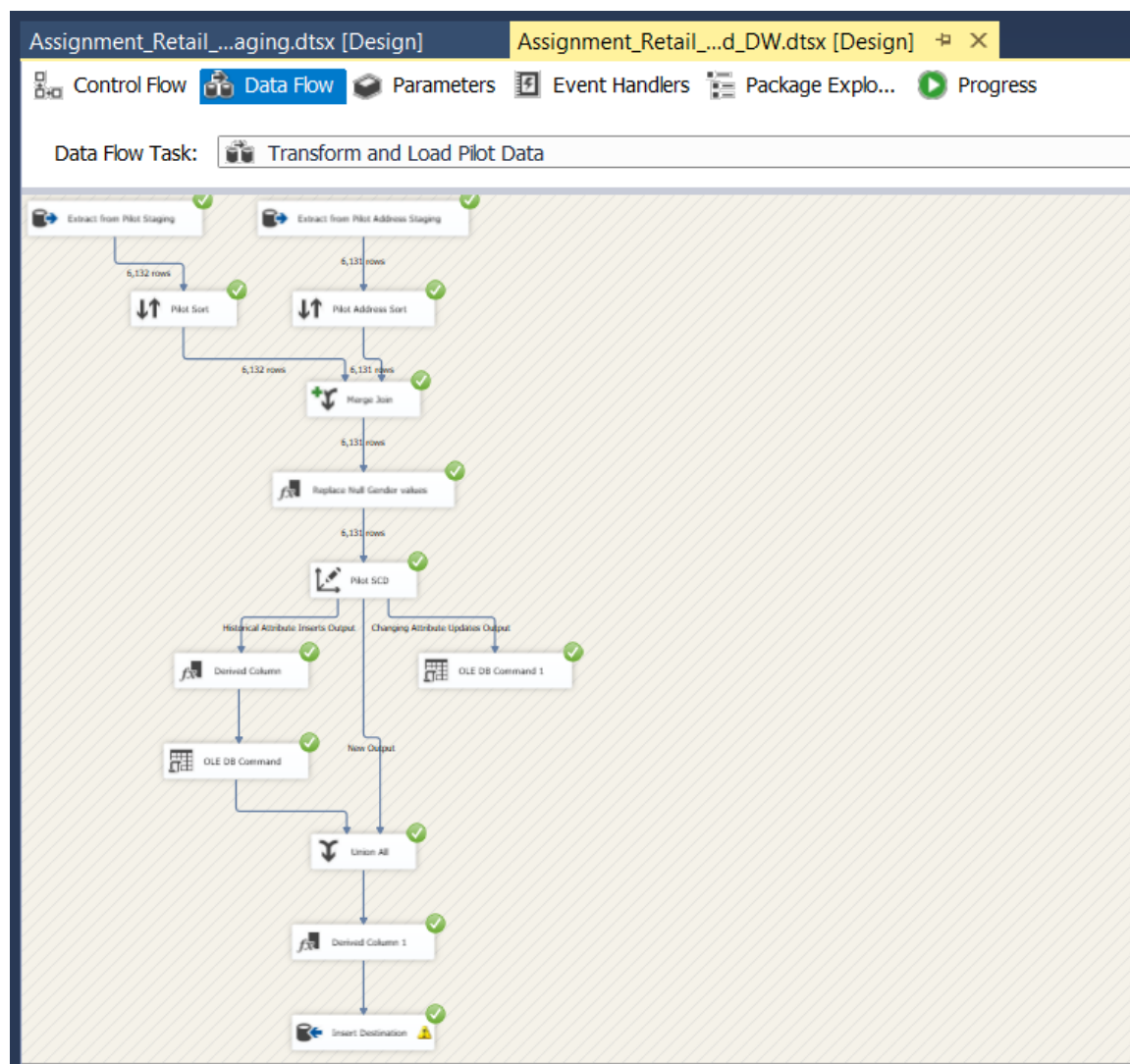
- **Transform and Load Accident Data**



- **Transform and Load Airport Data**
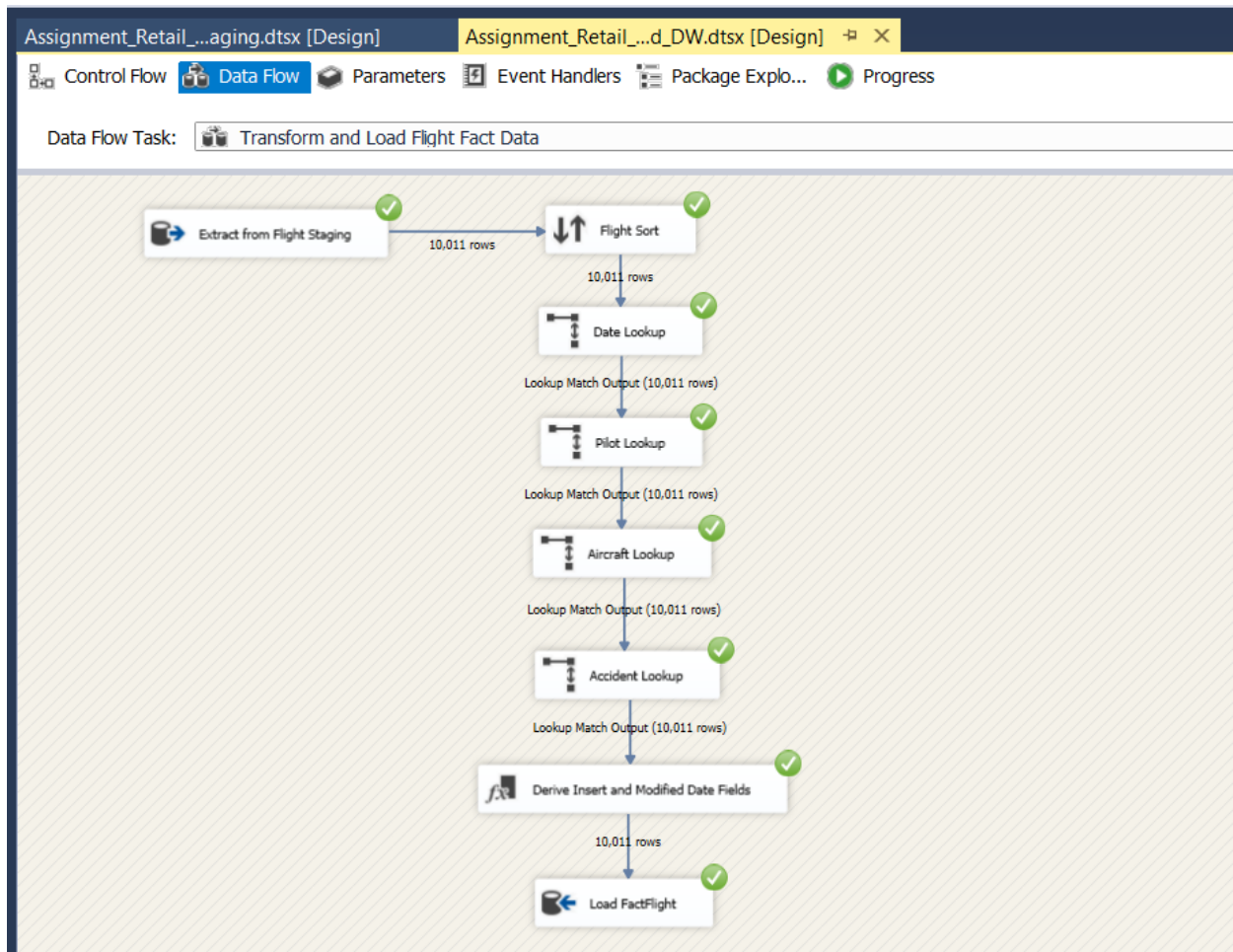


- **Transform and Load Aircraft Data**
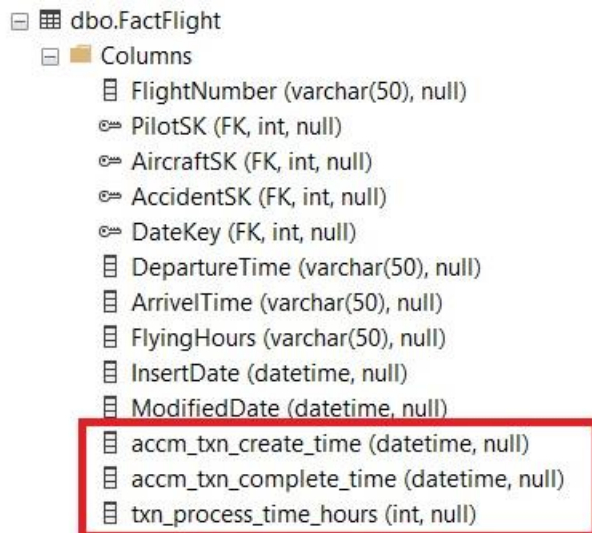
- **Transform and Load Pilot Data**
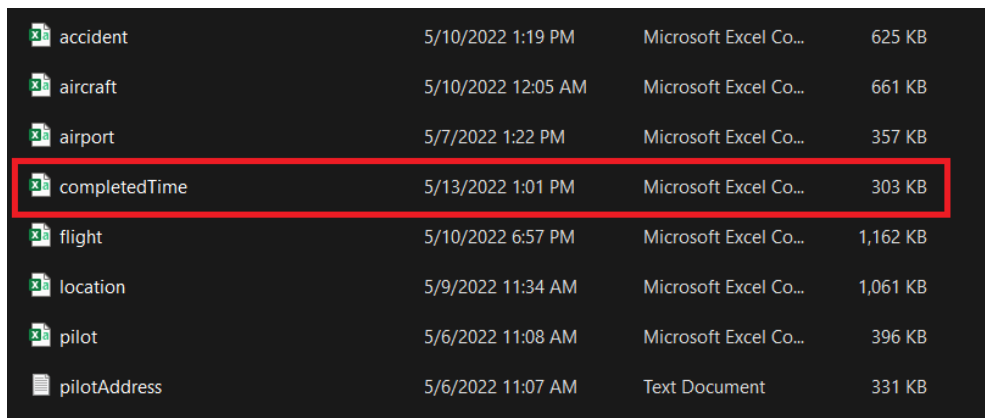
- **Transform and Load Flight Fact data**

- **Extending Fact Table with Additional Columns**



- **Prepare separate data set for complete time**



- **Update Complete Time and Process Time in Fact Table**

- Control flow

- Data Flows