

# **Sri Lanka Institute of Information Technology**



## **Data Warehousing and Business Intelligence IT3021 Assignment 1**

### **Submitted by:**

Full name: - H.A.M Senadheera

Registration No.: - IT20125998

Batch: - Y3.S1.DS.WE.04

## Table of Contents

<b><i>Data Set Selection.....</i></b>	<b><i>3</i></b>
<b><i>Introduction to the dataset .....</i></b>	<b><i>3</i></b>
<b><i>Preparation of the data set.....</i></b>	<b><i>4</i></b>
<b><i>Architecture.....</i></b>	<b><i>5</i></b>
<b><i>Class diagram using the sources tables .....</i></b>	<b><i>6</i></b>
<b><i>Snowflake Schema .....</i></b>	<b><i>7</i></b>
<b><i>ETL Process Screenshots .....</i></b>	<b><i>8</i></b>
<b><i>Extraction from Source database and flat files to staging area. ....</i></b>	<b><i>8</i></b>
Procedures for updating DimHost table.....	8
Procedure for updating DimListing table .....	9
Transform and load data to Data Warehousing tables. ....	9
Transform and load slowly changing dimension data to data warehouse. ....	10
Transform and load accumulating timestamp data to data warehouse.....	10

## Data Set Selection

Data Set name: - Boston Airbnb open data

Source: - Kaggle.com

Link to the source: - <https://www.kaggle.com/datasets/airbnb/boston>

## Introduction to the dataset

This is a dataset of Airbnb bookings in the city of Boston in United States of America in the year 2016. The initial dataset is adjusted and arranged to meet the requirements of the assignment and suit the scenario

Content of the dataset	
Hosts	Contains the data of individuals who have listed their properties in Airbnb
Listings	Contains the data of properties listed in the Airbnb
Customers	Contains the data of individuals who have done bookings/reservations for listed properties
CustomerAddress	Contains addresses of the customers.
Bookings	Contains the details of bookings that customers have made to listed properties.

## Preparation of the data set

Dataset was obtained as a single large csv containing listing and hosts.

While preparing the dataset listings and host were broken into two csv files.

Three more csv files were created by generating some sample data named as customer, booking and accumulated complete date csv.

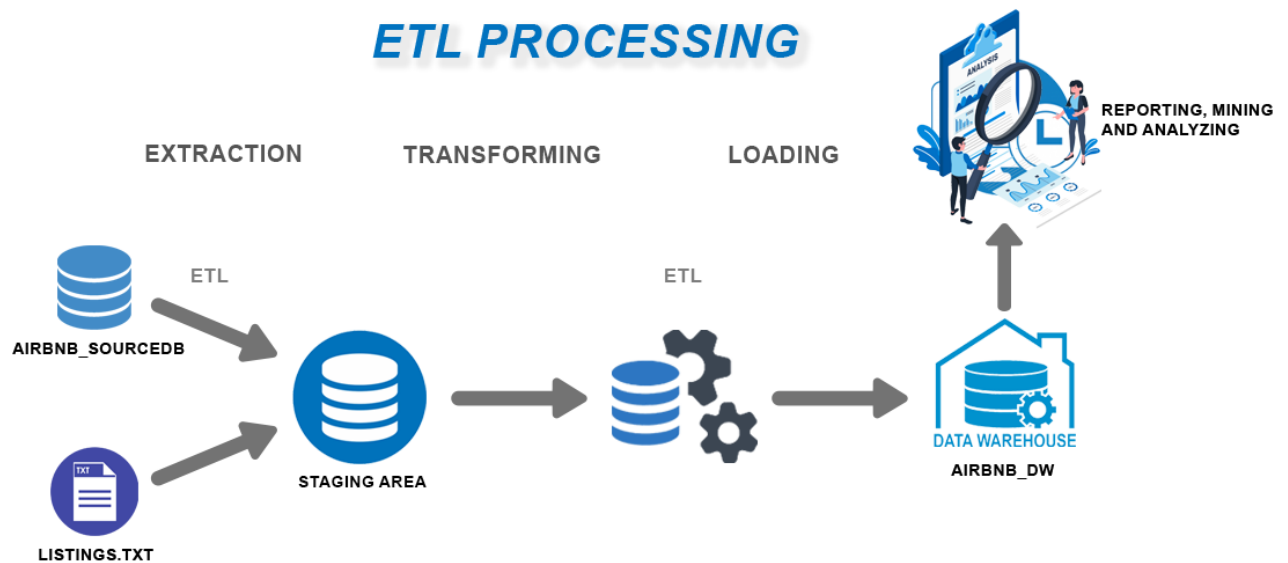
A text file was created by adding customer address information.

- Listing.csv
- Booking.csv
- Host.csv
- Customer.csv
- CustomerAddress.txt
- AccumulatedCompleteDate.csv

Using the csv files, a database named Airbnb\_SourceDb was created. The sourceDB, csv and text files are as follows.

- Airbnb\_SourceDB (source database)
  - dbo.Host
  - dbo.Listing
  - dbo.Customer
  - dbo.Booking
- CustomerAddress.txt
- AccumulatedCompleteDate.csv

## Architecture



A data source is an initial location where information is obtained for a given scenario. The primary data source in this case is a database and secondary data source is a flat file.

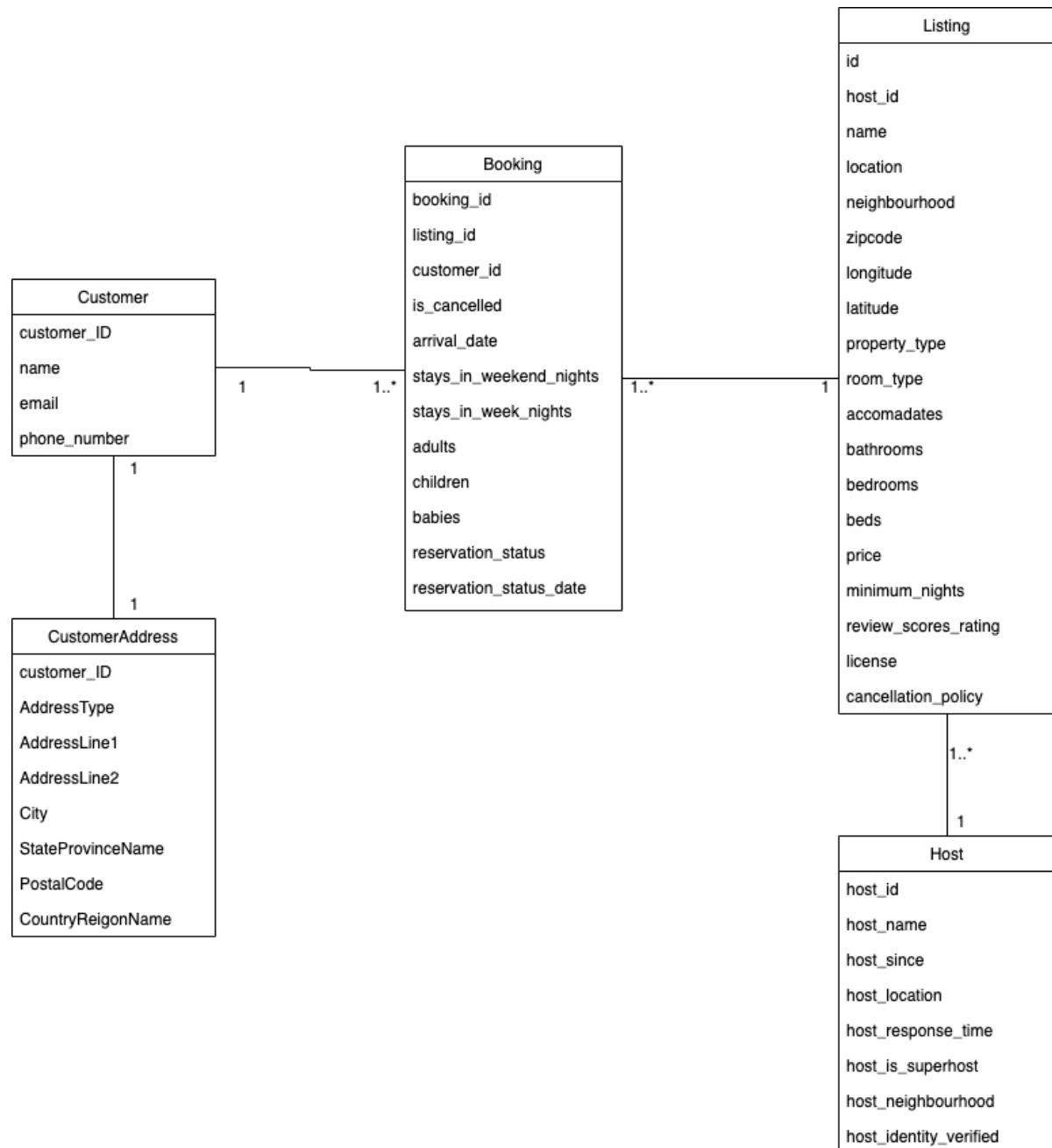
ETL, which stands for extract, transform and load, is a data integration process that extracts and combines data from multiple data sources into a single, consistent data store that is loaded into a Data Warehouse system or other target systems.

- Airbnb\_Staging (staging database)
  - dbo.StgHost
  - dbo.StgListing
  - dbo.StgCustomer
  - dbo.StgBooking

The data warehouse (DWH) is a central repository where an organization electronically stores data by extracting it from operational systems and making it available for data analysis and scheduled reporting. In contrast, the process of building a data warehouse entail designing a data model that can quickly generate insights.

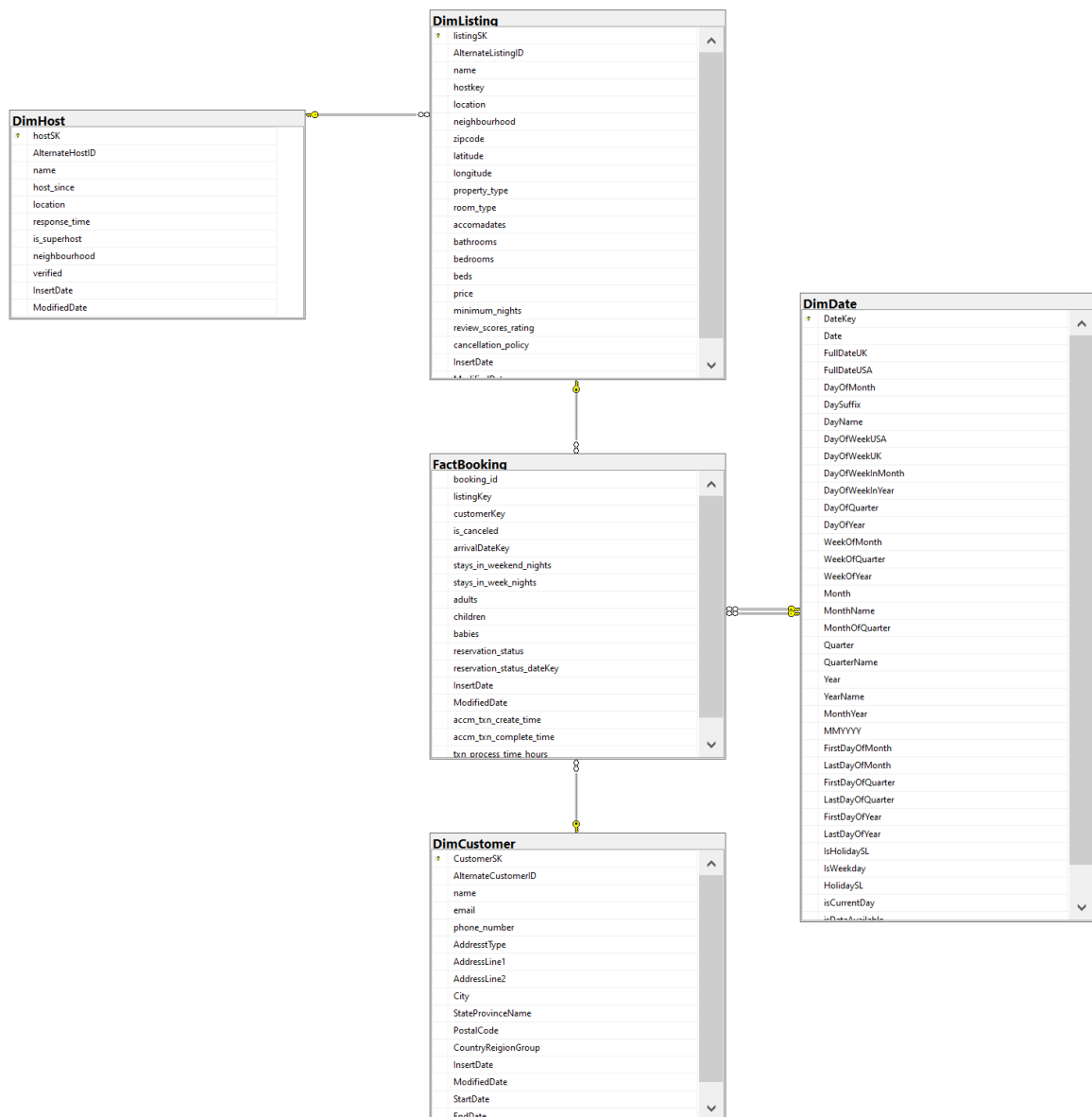
- Airbnb\_DW (data warehousing database)
  - dbo.DimHost
  - dbo.DimListing
  - dbo.DimCustomer
  - dbo.DimDate
  - dbo.FactBooking

## Class diagram using the sources tables



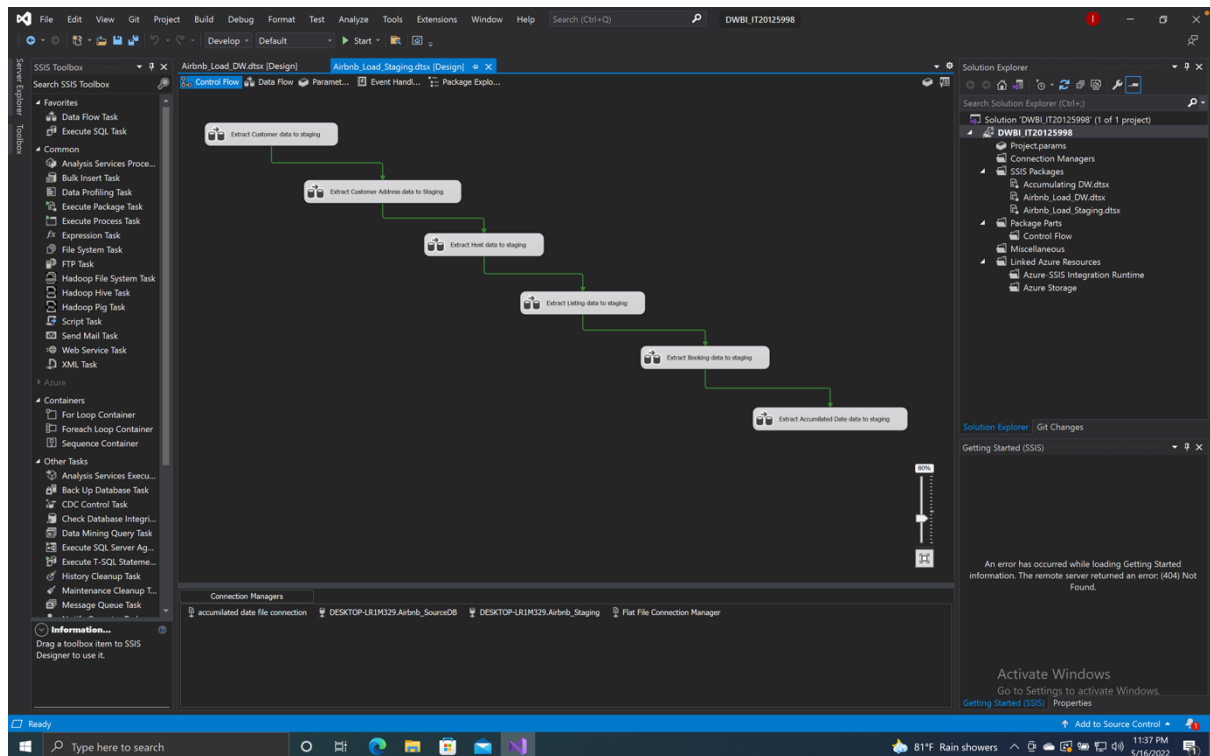
# Snowflake Schema

DimCustomer is a slowly changing dimension

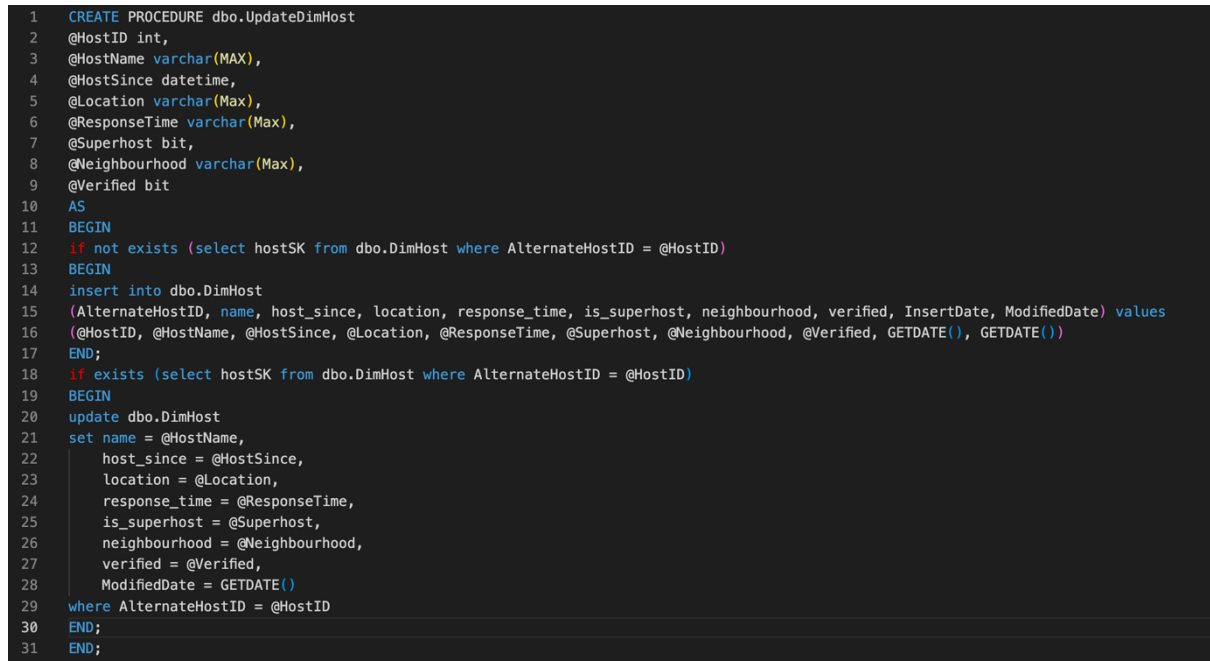


## ETL Process Screenshots

Extraction from Source database and flat files to staging area.



## Procedures for updating DimHost table

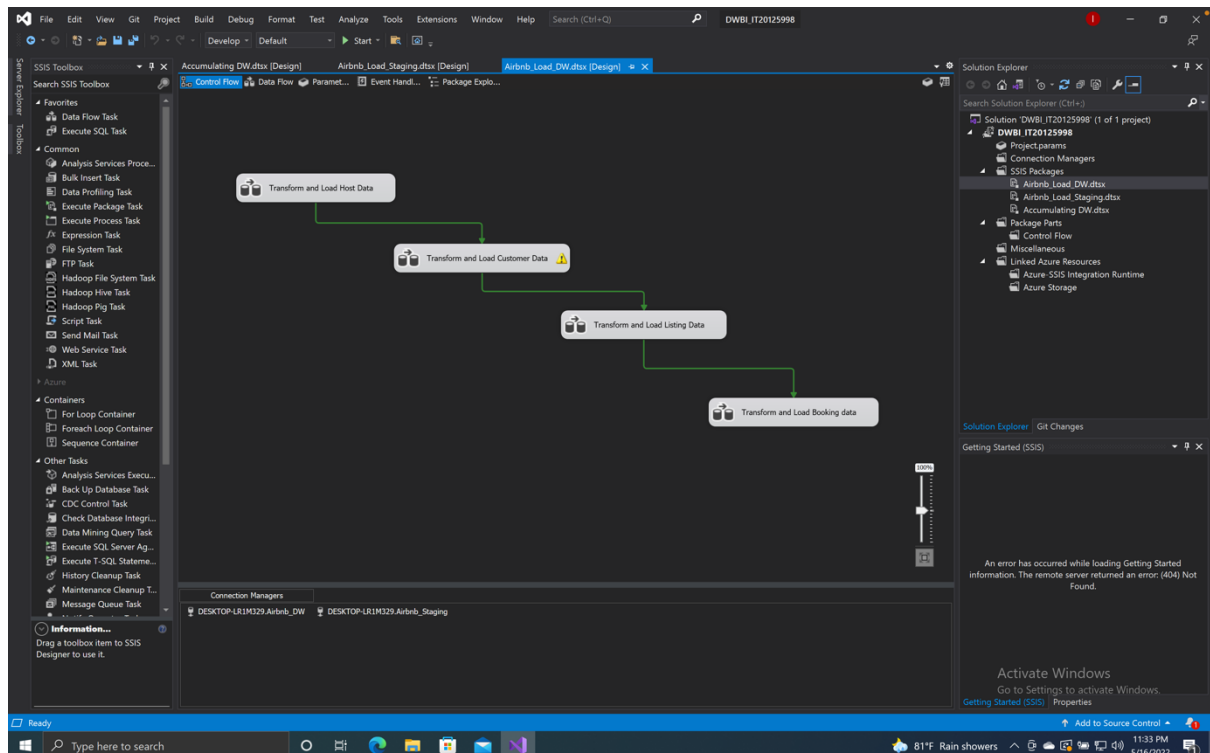




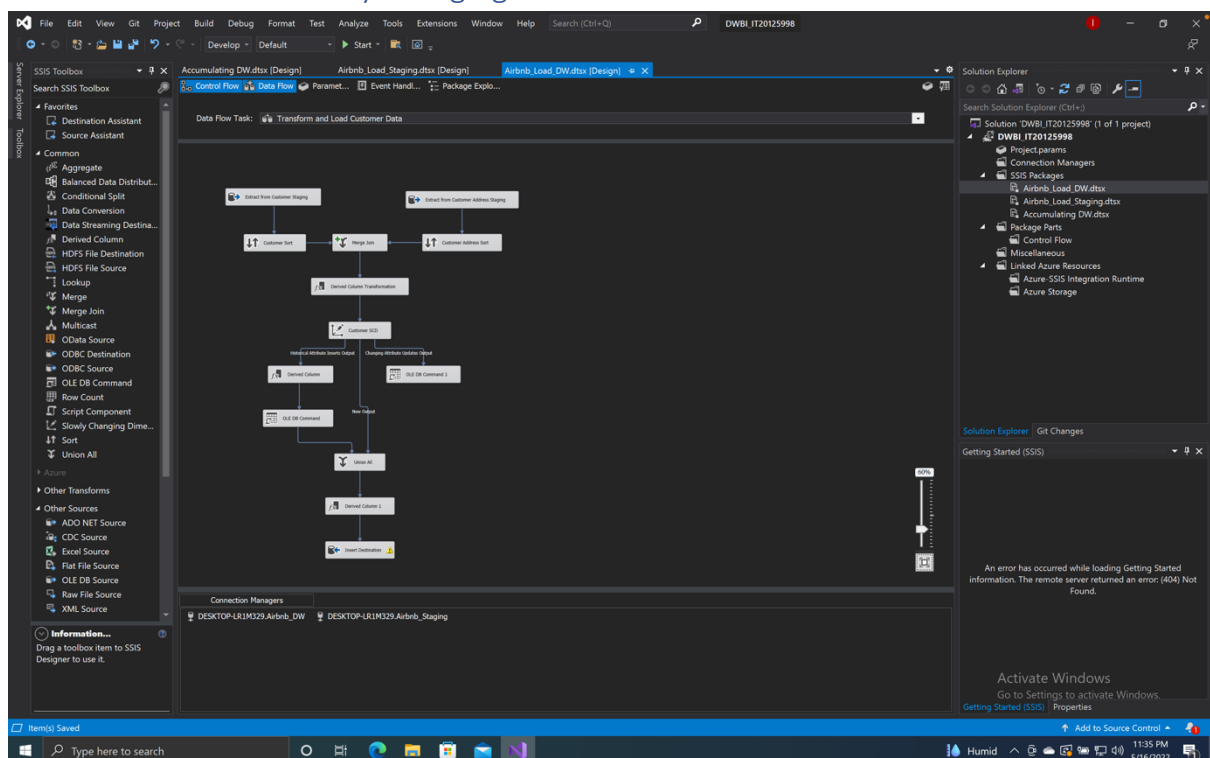
## Procedure for updating DimListing table

```
1 CREATE PROCEDURE dbo.UpdateDimListing
2 @ListingID int,
3 @Name varchar(100),
4 @HostKey int,
5 @Location varchar(100),
6 @Neighbourhood varchar(50),
7 @Zipcode varchar(50),
8 @Latitude varchar(15),
9 @Longitude varchar(50),
10 @PropertyType varchar(50),
11 @RoomType varchar(50),
12 @Accommodates int,
13 @Bathrooms float,
14 @Bedrooms int,
15 @Beds int,
16 @Price money,
17 @MinimumNights int,
18 @Review int,
19 @Cancellation varchar(50)
20 AS
21 BEGIN
22 if not exists (select listingSK from dbo.DimListing where AlternateListingID = @ListingID)
23 BEGIN
24 insert into dbo.DimListing (AlternateListingID, [name], hostkey, [location], neighbourhood, zipcode, latitude, longitude, property_type, room_type,
25 accomadates, bathrooms, bedrooms, beds, price, minimum_nights, review_scores_rating, cancellation_policy, InsertDate, ModifiedDate)
26 values(@ListingID, @Name, @HostKey, @Location, @Neighbourhood, @Zipcode, @Latitude, @Longitude, @PropertyType, @RoomType,
27 @Accommodates, @Bathrooms, @Bedrooms, @Beds, @Price, @MinimumNights, @Review, @Cancellation, GETDATE(), GETDATE())
28 END;
29 if exists (select listingSK from dbo.DimListing where AlternateListingID = @ListingID)
30 BEGIN
31 update dbo.DimListing
32 set [name] = @Name,
33 hostkey = @HostKey,
34 [location] = @Location,
35 neighbourhood = @Neighbourhood,
36 zipcode = @Zipcode,
37 latitude = @Latitude,
38 longitude = @Longitude,
39 property_type = @PropertyType,
40 room_type = @RoomType,
41 accomadates = @Accommodates,
42 bathrooms = @Bathrooms,
43 bedrooms = @Bedrooms,
44 beds = @Beds,
45 price = @Price,
46 minimum_nights = @MinimumNights,
47 review_scores_rating = @Review,
48 cancellation_policy = @Cancellation,
49 ModifiedDate = GETDATE()
50 where AlternateListingID = @ListingID
51 END;
52 END;
```

## Transform and load data to Data Warehousing tables.



Transform and load slowly changing dimension data to data warehouse.



Transform and load accumulating timestamp data to data warehouse.

