# Sri Lanka Institute of Information Technology

**Assignment 1**

Data Warehouse & Business Intelligence

2022

Submitted by:

Wickramaarachchi W.A.K.M

IT20073428

# Table of Contents

# 1. Data Selection and Preparation

The chosen data source is a collection of transactional data. A link to the source data set is provided below.

Dataset: - [AmExpert 2019](#)

The data set derived from the source was modified as needed. This dataset is about orders made at the store. The dataset contains information about around 1000k orders placed at several marketplaces along with discounts.

The dataset contains CSV files with information about campaign data, customer data, and product data. Modifications were made to the data set derived from the source as needed. This data set consists of combinations of transactions and promotional initiatives.

- campaign_data.csv – campaign information for each campaign
- customer_demographics_data.csv – customer information for customers
- customer_transaction_data.csv – transaction data for all customers for the duration of campaigns
- item_data.csv – item information for each item

All the data sources are provided as CSV files by the Kaggle. Therefore, in the preparation of data sources, some CSV files were imported to the source database and added columns and separated some data to make other data files.

The final types of data sources are mentioned below:

- **SQL Database**
  - Campaign
  - Customer
  - CustomerAddress
  - Transaction

- **CSV files**
  - item_category.csv
  - item_data.csv

# 2. Description of the data set

| Name: | Customer | |
|---|---|---|
| **Source Type:** | SQL Database | |

| Column Name | Data Type | Description |
|---|---|---|
| customer_id | int | Unique identifier of a customer |
| first_name | nvarchar(50) | First name of the customer |
| last_name | nvarchar(50) | Last name of the customer |
| title | nvarchar(10) | Title of the customer |
| gender | nvarchar(3) | Customer's gender |
| email | nvarchar(50) | Email address of the customer |
| phone | nvarchar(20) | Phone number of the customer |
| age_range | nvarchar(50) | Customer age range |
| marital_status | nvarchar(20) | Marital status of the customer |
| rented | bit | Rented – 1 / not – 0 |
| family_size | nvarchar(10) | Size of the customer's family |
| no_of_children | nvarchar(10) | No. of children in the family |
| income_bracket | int | Customer income bracket |

| Name: | CustomerAddress | |
|---|---|---|

| Source Type: | SQL Database | |
|---|---|---|

| Column Name | Data Type | Description |
|---|---|---|
| customer_id | int | Customer's unique Id |
| city | nvarchar(30) | Customer's city |
| country | nvarchar(50) | Customer's country |
| country_code | nvarchar(5) | Customer's country code |
| latitude | float | Latitude |
| longitude | float | Longitude |
| street_address | nvarchar(150) | Customer's street address |
| street_name | nvarchar(50) | Customer's street name |
| street_number | Int | Customer's street number |
| street_suffix | nvarchar(20) | Customer's street suffix |

| Name: | Campaign | |
|---|---|---|

| Source Type: | SQL Database | |
|---|---|---|

| Column Name | Data Type | Description |
|---|---|---|
| campaign_id | Int | Unique identifier for campaign |
| campaign_type | nvarchar(5) | Type of the campaign |
| start_date | datetime | Campaign start date |
| end_date | datetime | Campaign end date |

| **Name:** Transaction | | |
|---|---|---|
| **Source Type:** SQL Database | | |
| **Column Name** | **Data Type** | **Description** |
| transaction_id | int | Unique identifier of the transaction |
| customer_id | int | Customer unique identifier |
| item_id | int | Item unique identifier |
| date | datetime | Order placement date |
| quantity | int | Order item quantity |
| selling_price | money | Item selling rice |
| other_discount | money | Discounts |
| coupon_discount | money | Coupon discount |

| **Name:** Item Category | | |
|---|---|---|
| **Source Type:** CSV file | | |
| **Column Name** | **Data Type** | **Description** |
| category_id | int | unique category identifier |
| category_name | nvarchar(50) | Name of the item category |

| Name: | Item Data | |
|---|---|---|

| Source Type: | CSV File | |
|---|---|---|

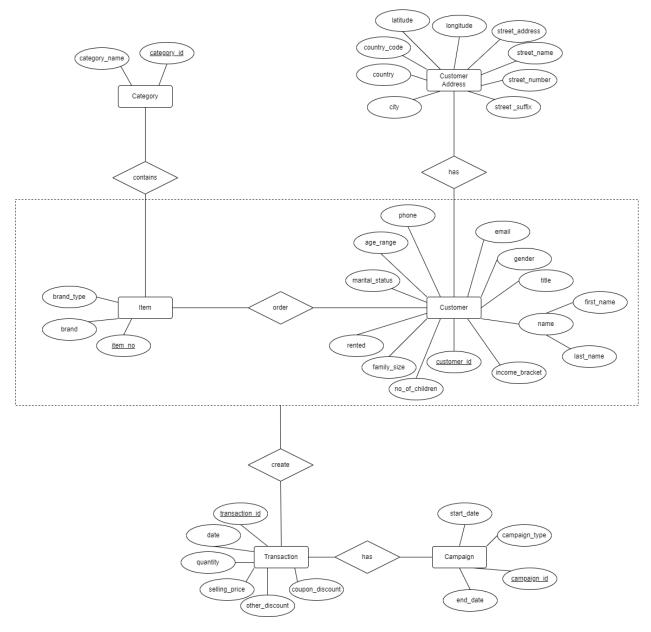| Column Name | Data Type | Description |
|---|---|---|
| item_id | int | Item unique identifier |
| brand | int | Item brand code |
| brand_type | nvarchar(50) | Item brand type |
| category | int | Item category id |

# 3. ER Diagram

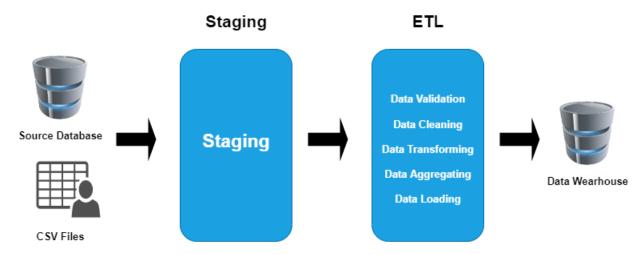The above diagram shows the connections between entities

# 4. Solution Architecture
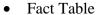


*Figure 2: High-Level BI Solution Architecture*

In the staging layer;

- StgCustomer
- StgCustomerAddress
- StgCampaign
- StgItem
- StgItemCategory
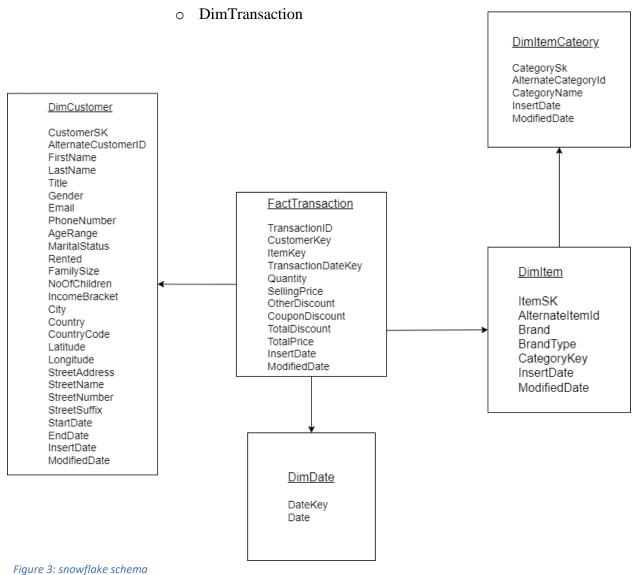- StgTransaction tables were created

# 5. Data Warehouse Design and Development

## a) Design

AmExpert_DW is designed according to the snowflake schema with one fact table and four dimension tables, including Date Dimension.

- Dimension Tables
  - DimDate
  - DimCustomer
  - DimItem
  - DimItemCategory

- Fact Table
  - DimTransaction

**DimItemCateory**

CategorySk
AlternateCategoryId
CategoryName
InsertDate
ModifiedDate

**DimCustomer**

CustomerSK
AlternateCustomerID
FirstName
LastName
Title
Gender
Email
PhoneNumber
AgeRange
MaritalStatus
Rented
FamilySize
NoOfChildren
IncomeBracket
City
Country
CountryCode
Latitude
Longitude
StreetAddress
StreetName
StreetNumber
StreetSuffix
StartDate
EndDate
InsertDate
ModifiedDate

**FactTransaction**

TransactionID
CustomerKey
ItemKey
TransactionDateKey
Quantity
SellingPrice
OtherDiscount
CouponDiscount
TotalDiscount
TotalPrice
InsertDate
ModifiedDate

**DimItem**

ItemSK
AlternateItemId
Brand
BrandType
CategoryKey
InsertDate
ModifiedDate

**DimDate**

DateKey
Date

*Figure 3: snowflake schema*

* Hierarchies

    – DimItemCategory is a hierarchical dimension of DimItem
    – DimCustomer table has hierarchical attributes about customer address

* Calculations
    – Total Discount is calculated in the FactTransaction table
        (([OtherDiscount] + [CouponDiscount]) * [Quantity])

    – Total Price is calculated in the FactTransaction table
        ( ( [SellingPrice]) * [Quantity]) – [TotalDiscount])

## b) Assumptions

- Transaction table used for creating fact table
- Transaction per customer was considered as the grain

## c) Slowly changing dimensions

- DimCustomer was considered a slowly changing dimension

| Changing Attributes | Historical Attributes |
|---|---|
| MaritalStatus | AgeRange |
| PhoneNumber | City |
| Title | Country |
| FamilySize | CountryCode |
| NoOfChildren | Latitude |
| | Longitude |
| | StreetAddress |
| | StreetName |
| | StreetNumber |
| | StreetSuffix |

# 6. ETL Development

I. Data Extraction and Loading into Staging Tables
- Data extraction is done by using Visual Studio Data Tools. The CSV files and database are used as data sources.
- OLE DB SOURCE is used to extract data from the database, and FLAT FILE SOURCE is used to extract data from CSV files.
- For the load data to the staging area, OLE DB DESTINATION was used
- Use EXECUTE SQL TASK to truncate Staging tables before load data. This will prevent data from being duplicated in tagging tables

**Control Flow**



*Attachment 1: Control Flow*

## Staging Campaign Details



Campagn data is extracted from AmExpert SourceDB and inserted into Campaign staging table.

*Attachment 2: Staging campaign data*

## Staging Customer Details



Customer data is extracted from AmExpert SourceDB and inserted into Customer staging table.

*Attachment 3: Staging customer details*

## Staging Customer  Address Details



Customer Address data is extracted from AmExpert SourceDB ans insert into Customer Address staging table

*Attachment 4: Staging customer address data*

## Staging Item Details



Item details is extracted from CSV file and insert into Item staging table

*Attachment 5: Staging item details*

## Staging Item Category Details



Item Category data is extracted from CSV file and insert into Item Category staging table

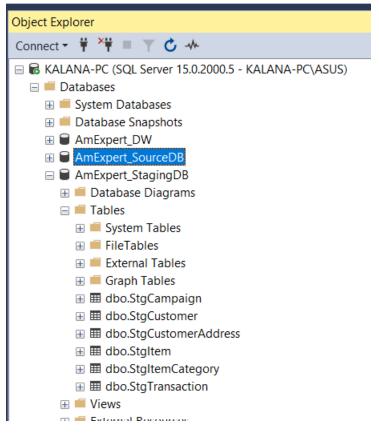*Attachment 6: Staging item category details*

## Staging Transaction Details



Transaction data is extracted from AmExpert SourceDb and insert into Transaction staging table
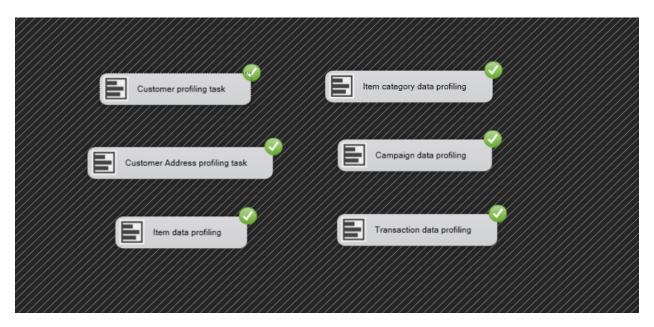
*Attachment 7: Staging transaction details*

## Created Staging Tables



*Attachment 8: Created staging tables*

## II.    Data Profiling

Data Profiling enables the analysis of enormous amounts of data using various procedures. Null values, repeated values, and data quality are all examined in this step.
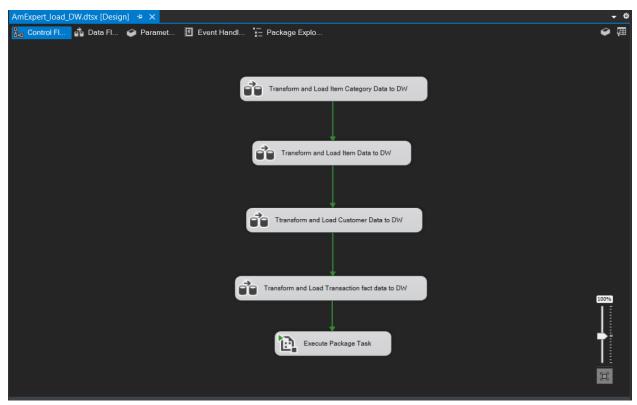


*Attachment 9: Data profiling task*

- Every staging table had run a profiling task and saved output in a selected location
- By referring these data profiles, the developer is able to identify the issues with staging data
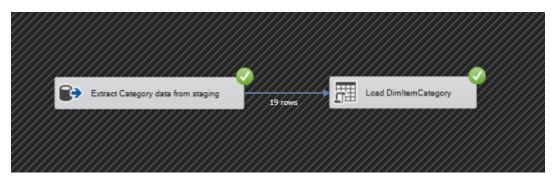
## III. Data Transforming and Loading

- Data transformation is developed according to the dimension modelling designed above (figure 3)
- Dimension tables are loaded with data from relevant staging tables in this step.



*Attachment 10: Control Flow Data Wearhouse*

## Transform and Load Item Category Data

- Item Category data is loaded into DimCategory table
- UpdateItemCategory Procedure is used to identify data should be insert or not



*Attachment 11: Load DImCategory*



```sql
UpdateCategoryPr...ANA-PC\ASUS (56))
CREATE PROCEDURE dbo.UpdateItemCategory
@CategoryID int,
@CategoryName nvarchar(50)
AS
BEGIN
if not exists(select CategorySK
from dbo.DimItemCategory
where AlternateCategoryID = @CategoryID)
BEGIN
insert into dbo.DimItemCategory
(AlternateCategoryID, CategoryName, InsertDate, ModifiedDate)
values
(@CategoryID, @CategoryName, GETDATE(), GETDATE())
END;
if exists(select CategorySK
from dbo.DimItemCategory
where AlternateCategoryID = @CategoryID)
BEGIN
update dbo.DimItemCategory
set
AlternateCategoryID = @CategoryID,
CategoryName = @CategoryName,
ModifiedDate = GETDATE()
where AlternateCategoryID = @CategoryID
END;
END;
```

*Attachment 12: UpdateItemCategory Procedure*

## Transform and Load Item Data

- Item data is loaded into DimItem table
- Lookup transformation task used for extract CategoryKey from DimCategory
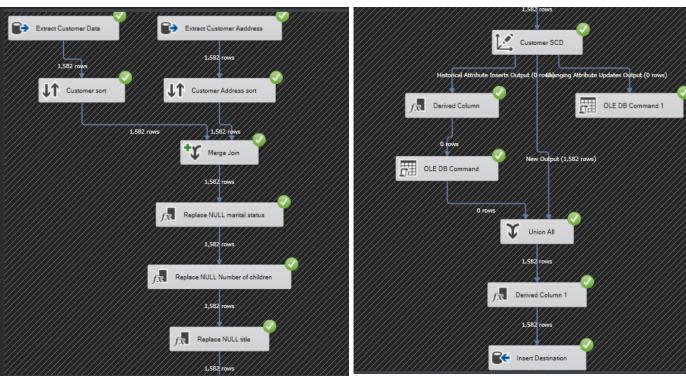- UpdateItem Procedure is used to identify data should be insert or not



*Attachment 13: Load DimItem*



```
CREATE PROCEDURE dbo.UpdateItemData
@ItemID int,
@Brand int,
@BrandType nvarchar(50),
@CategorySK int
AS
BEGIN
if not exists(select ItemSK
from dbo.DimItem
where AlternateItemID = @ItemID)
BEGIN
insert into dbo.DimItem
(AlternateItemID, Brand, BrandType, CategoryKey, InsertDate, ModifiedDate)
values
(@ItemID, @Brand, @BrandType, @CategorySK, GETDATE(), GETDATE())
END;
if exists(select ItemSK
from dbo.DimItem
where AlternateItemID = @ItemID)
BEGIN
update dbo.DimItem
set
AlternateItemID = @ItemID,
Brand = @Brand,
BrandType = @BrandType,
CategoryKey = @CategorySK,
ModifiedDate = GETDATE()
where AlternateItemID = @ItemID
END;
END;
```

*Attachment 14: UpdateItemData Procedure*

# Transform and Customer Data (Slowly Changing Dimension)

- DimCustomer is the Slowly Changing Dimension (SCD) in this modelling
- StartDate and EndDate columns ensure that the data valid at the moment
- Slowly changing dimension wizard used to implement DimCustomer model
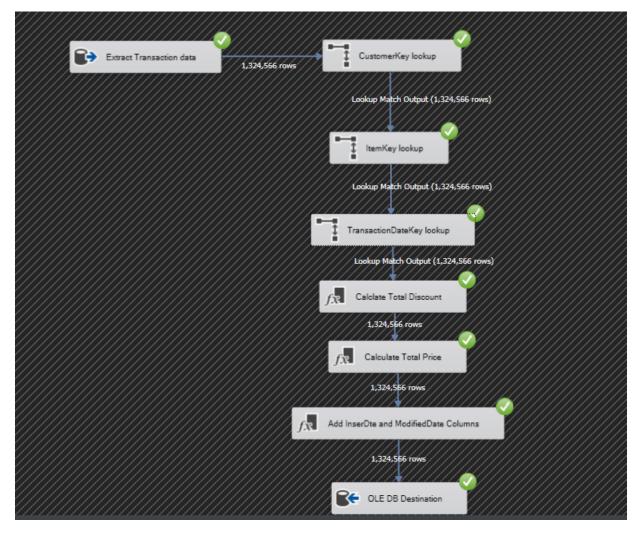- Derived Column Transformations are used to fill null values (marital status, number of children, title)



*Attachment 15 - i: Load DimCustomer*



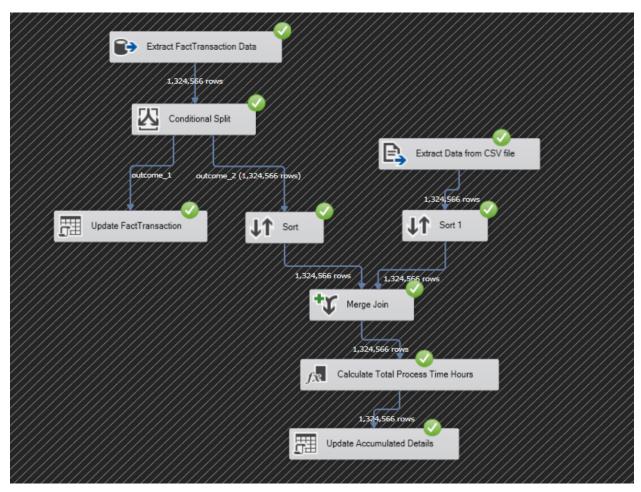*Attachment 15 - ii:Load DimCustomer*

# Load Data to Fact Table

- According to the dimension model, StgTransaction table used to insert values into FactTransaction table
- After loading all the dimension tables, lastly data was inserted in to the fact table. Below steps were followed

  1. Data extracted from the StgTransaction staging table
  2. Join operation done for the CustomerKey using lookup
  3. Join operation done for the ItemKey using lookup
  4. Join operation done for the TransactionDateKey using lookup
  5. Join operation done for the CustomerKey using lookup
  6. Calculate Total Discount and Total Price using derived column transformations
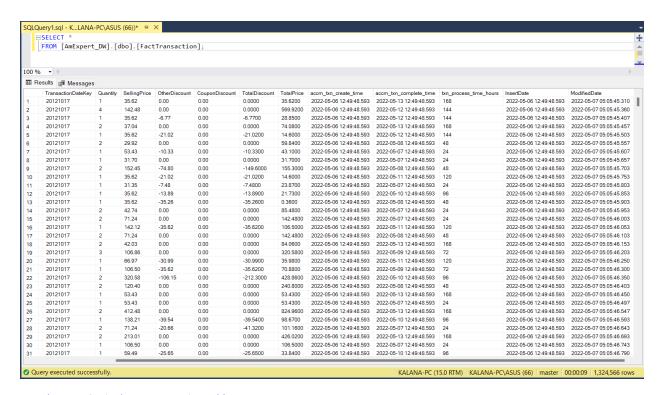  7. Insert and Modified dates were derived
  8. Load FactTransaction table



*Attachment 16: Load FactTransaction*

## Load Accumulated Data to Fact Table

- In life cycle of transaction process, Fact table should be updated with current status of transaction
- Therefore developers use accumulated type fact table to implement the solution. Below steps were followed
    1. Extract current Fact data
    2. Use conditional split to check whether accumulated data already exist or not
    3. If not exists, extract new accumulated data from source and join new data with fact data using merge join
    4. Done calculations using derived column transformation
    5. Update Fact table



*Attachment 17: Load Accumulated data*

Attachment 18: Final FactTransaction table