

Sri Lanka Institute of Information Technology

Data warehousing and Business Intelligence

Assignment 1



Student Registration No: IT20186142

Student Name: Wijesooriya H.M.A.H.

Step 1: Data Set Selection

This data set contains shopping online analytics of a famous super store which centers are distributed globally. They reach out to customers who do orders online. Each online order has a relevant customer, a product, a market and a shipment mode. Through this data set, it is able to profile the customers based on their of purchases and to profile countries based on the sales.

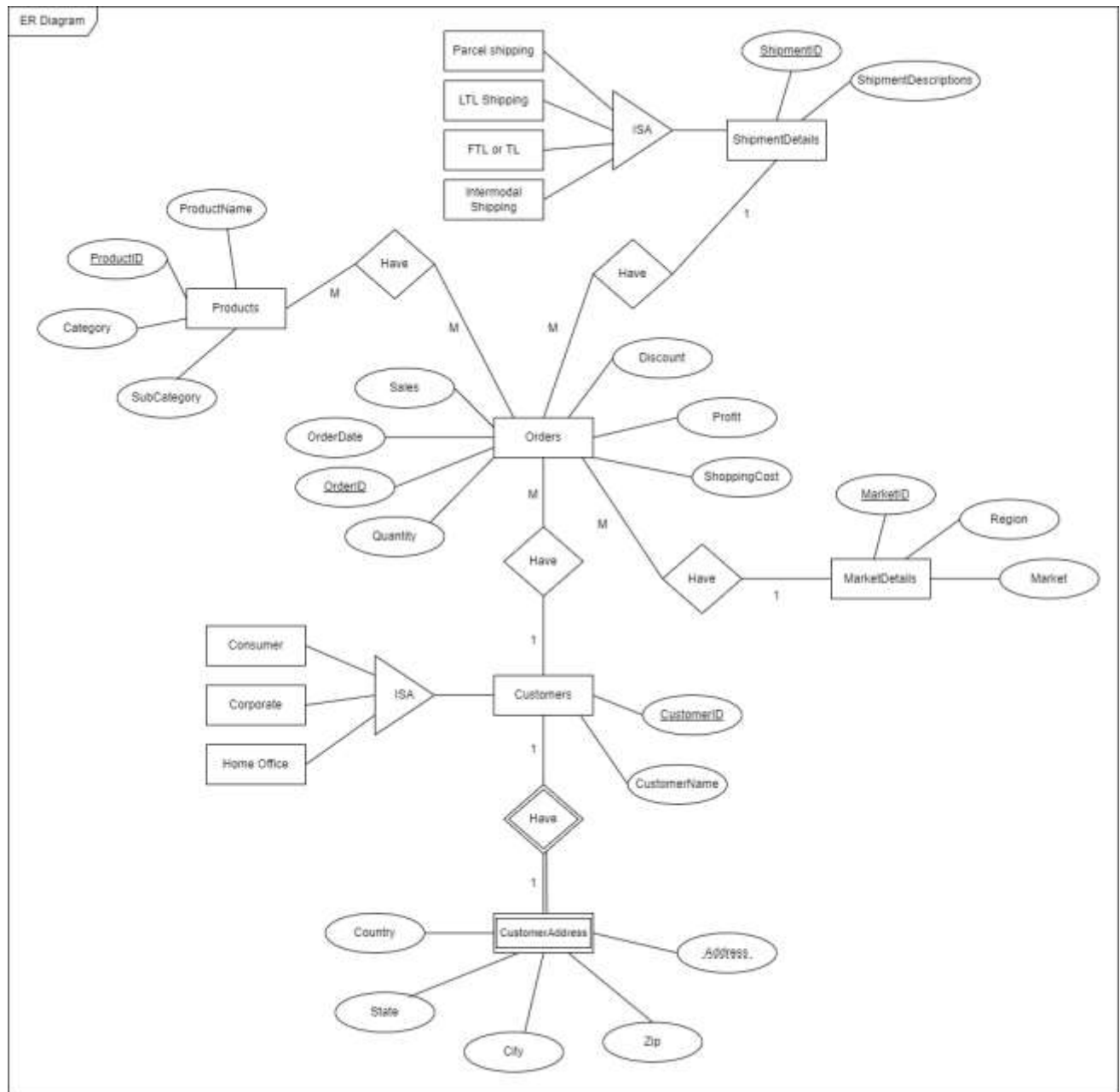
This data set includes details about 20,000 orders which were happened throughout 4 years and the number of customers involved in these orders are over 1500. In those orders there are about 2600 significant products which varies with category and subcategory.

This dataset contains Super store details,

- Customer details
- Customers segment details
- Customer addresses
- Order details
- Market details
- Product details
- Shipment details
- Product category details
- Product subcategory details

Also, there are some added details to this database.

Following ER- diagram will describe the scenario of the selected dataset



Step 2: Preparation of Data Sources

The whole of data was in 'csv' file type and they were separated into the following data sources, Database, Text, Excel and csv. And they were used to create the following,

1.Database(.bak)

Category.xls, MarketDetails.xls, OrderPriorityDetails.xls, Products.xls, SubCategory.xls and ShipmentDetails.xls files were imported to the Globa_Super_Store Database.

2.Text(.txt)

CustomerAddress.txt was used directly.

3.Excel(.xls)

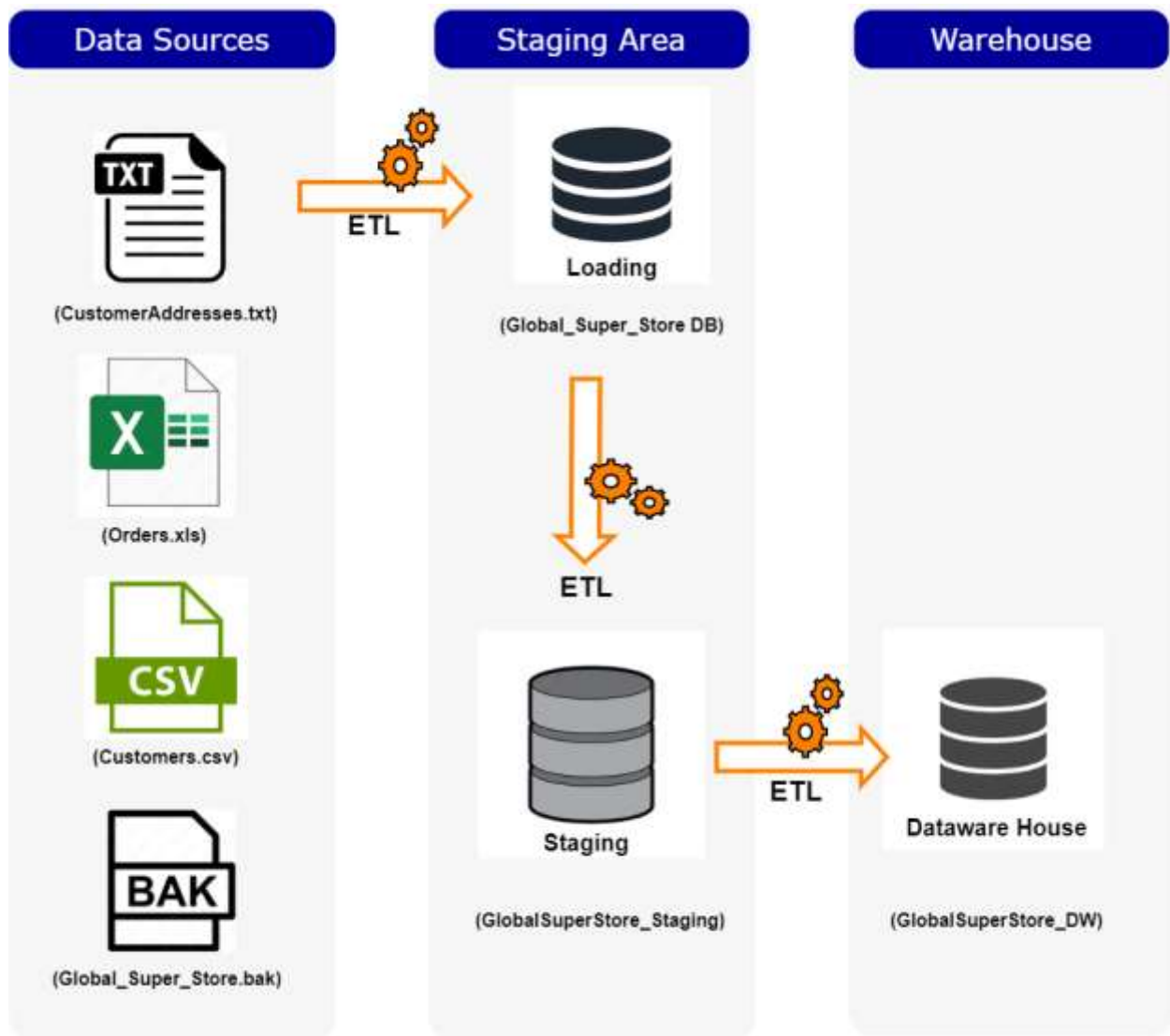
Orders.xls was used.

4.Comma Separated Values (.csv)

Customers.csv was used.

Data Source Type	Source Name	Column Name	Data Type	Description
Database File (.bak)	dbo.Products	Product_ID	nvarchar(50)	Unique ID
		Product_Name	nvarchar(150)	Name of Product
		CategoryID	int	Product Category ID
		Sub_CategoryID	int	Product Subcategory ID
	dbo.Category	CategoryID	int	Unique ID
		Category	nvarchar(50)	Product Category name
	Dbo.SubCategory	Sub_CategoryID	int	Unique ID
		Sub_Category	nvarchar(50)	Product SubCategory name
	dbo.MarketDetails	MarketID	int	Unique ID
		Market	nvarchar(50)	Market Name
		Region	nvarchar(50)	Region of market
	dbo.OrderPriorityDetails	OrderPriorityID	int	Unique ID
		OrderPriority	nvarchar(50)	Order Priority types
	Dbo.ShipmentDetails	ShipmentID	int	Unique ID
		Ship_Mode	nvarchar(50)	Shipment types
		ShipmentDescriptions	nvarchar(100)	Shipment type descriptions
Excel File	Orders.xls	Row ID	nvarchar(50)	Unique ID
		OrderID	nvarchar(50)	ID of the order
		OrderDate	date	Order placed date
		OrderPriority ID	int	ID of Order Priority
		Customer ID	nvarchar(50)	ID of Customer
		Sales	money	Sales
		Quantity	int	Quantity
		Discount	decimal(4, 3)	Discount of Product
		Profit	money	Profit
		Shipping Cost	money	Shipping Cost
		MarketID	int	ID of Market
		Product ID	nvarchar(50)	ID of Product
		ShipmentID	int	ID of Shipment mode
		ShipmentDate	date	Shipment placed date
CSV File	Customers.csv	Customer ID	varchar(50)	Unique ID
		Customer Name	varchar(50)	
		Segment	varchar(50)	
Text File	CustomerAddresses.txt	Customer ID	varchar(50)	Unique ID
		Country	varchar(50)	Customer's Country
		State	varchar(50)	Customer's State
		City	varchar(50)	Customer's City
		ZIP	varchar(50)	ZIP code of the Customer
		Address	varchar(60)	Customer's Address

Step 3: Solution Architecture



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

Data Sources

‘.txt’ component represents Text files, ‘.xls’ component is used to represent Excel files, ‘.csv’ component is used to display Comma Separated files and ‘.bak’ component represents database files.

Staging Area

Loading DB component represents the process of the creating database tables. Test, Patient, TestPrices, AdmissionFees and Attendance text files was imported to the database and was used to create the tables. And these tables were used as the DB source data.

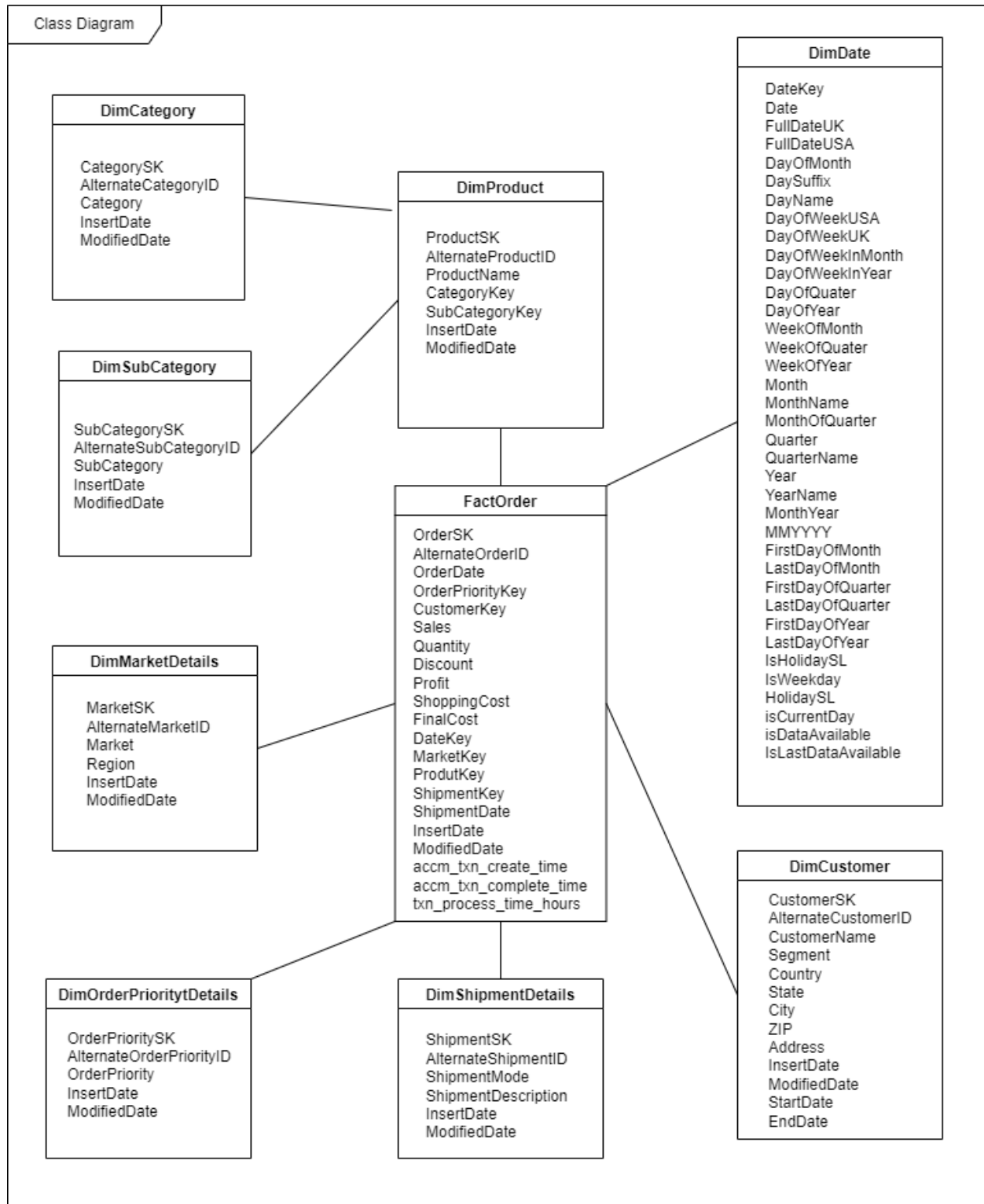
Staging DB component represents creating staging level tables through the 'Extract'.

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.



Schema Type

For this scenario, snowflake schema type was used.

Dimension Types

- Hierarchical Dimension
 - Date – all the hierarchies in date
 - Product – product → category → subcategory
 - CustomerAddress – country → state → city → ZIP → address
- Slowly Changing Dimension
 - Customer – used type 2.
 - Following columns were set as changing attributes.
 - Segment
 - Country
 - State
 - City
 - ZIP code
 - Address
- Fact Table
 - Numbers – Sales, Quantity, Discount, Profit, ShippingCost, FinalCost
 - FK - CustomerKey, OrderDateKey, OrderPriorityKey, MarketKey, ProductKey, ShipmentKey

Assumptions

- Customer dimension was considered as a slowly changing dimension.

Step 5: ETL development

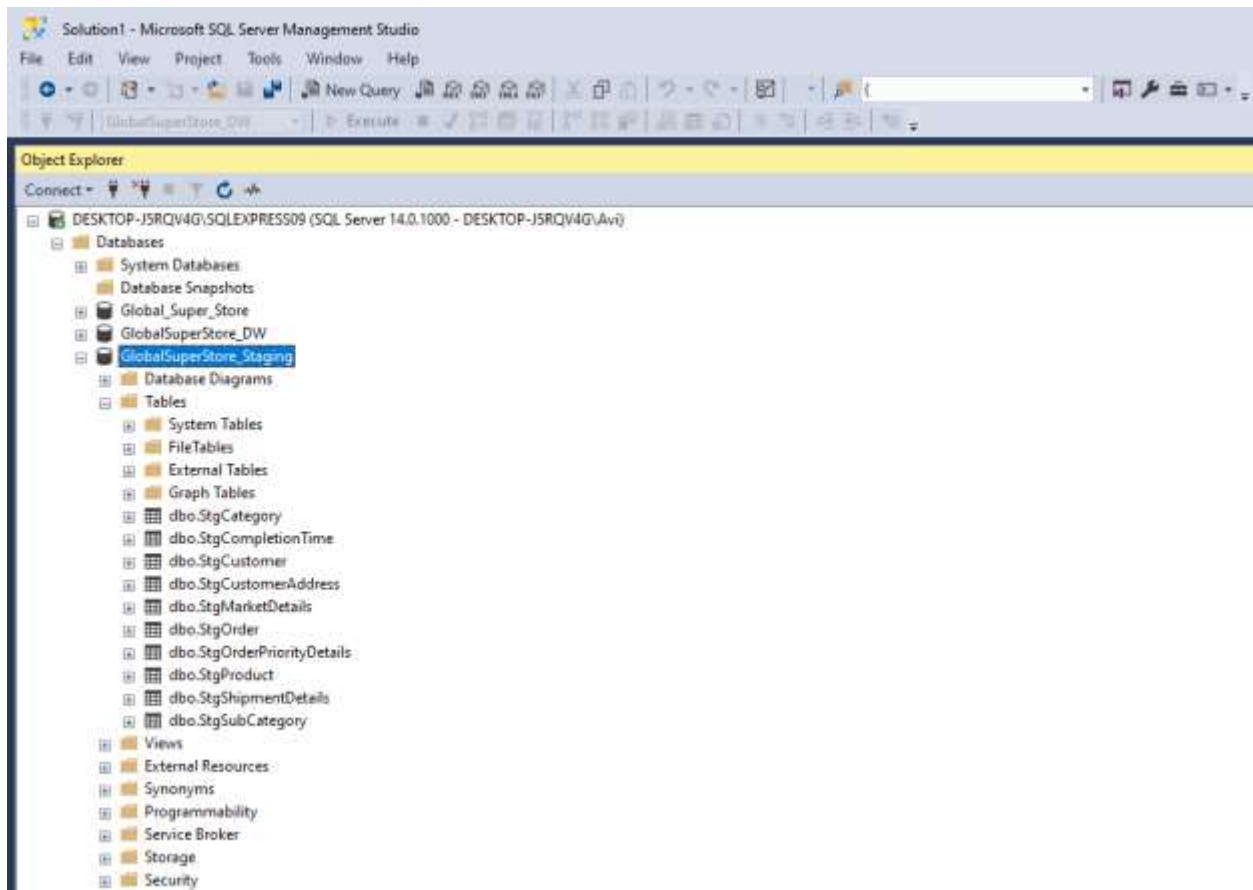
1.Extract

In this step, All the data sources were imported to the staging tables by using the relevant Data connection.

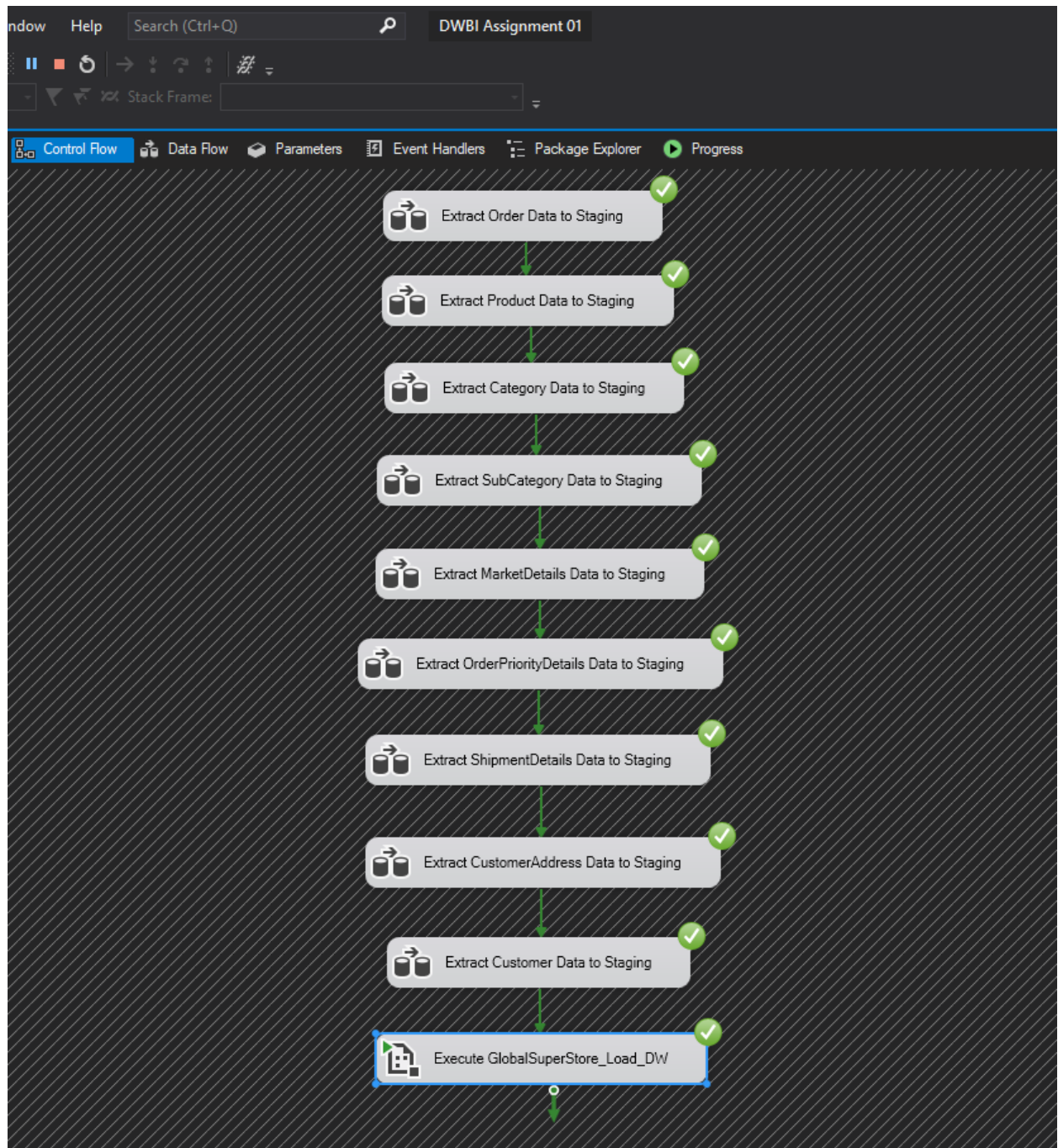
Flat file connection was used for text files and csv files, Excel file connections for excel file, DB source connection for DB file. All those tables were imported to the GlobalSuperStore_Staging DB, which contains the below tables,

1. StgCategory
2. StgCustomer
3. StgCustomerAddress
4. StgMarketDetails
5. StgOrder
6. StgOrderPriorityDetails
7. StgProduct
8. StgShipmentDetails
9. SgtSubCategory

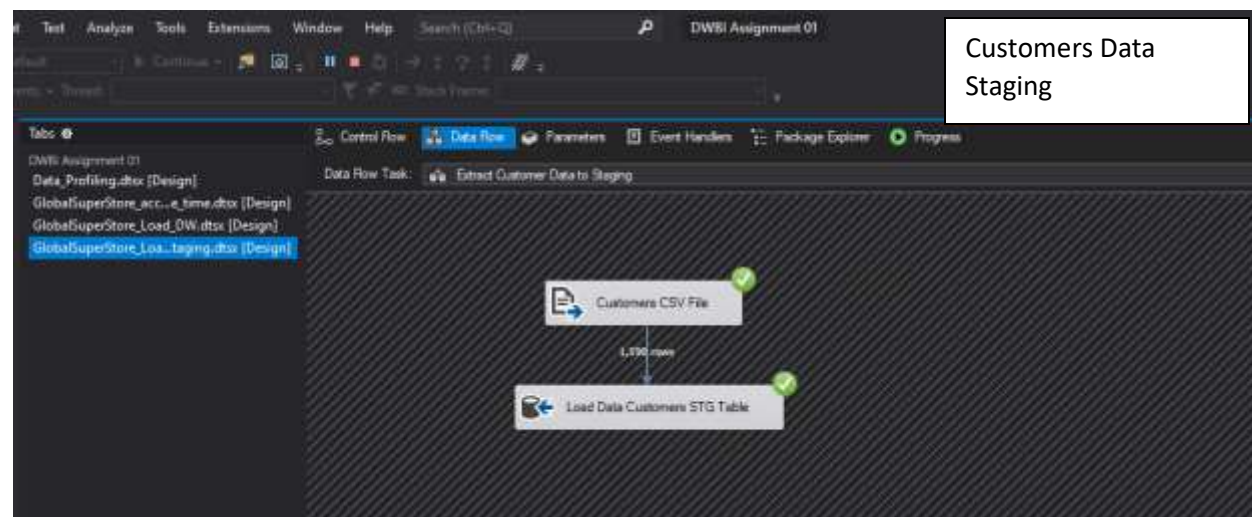
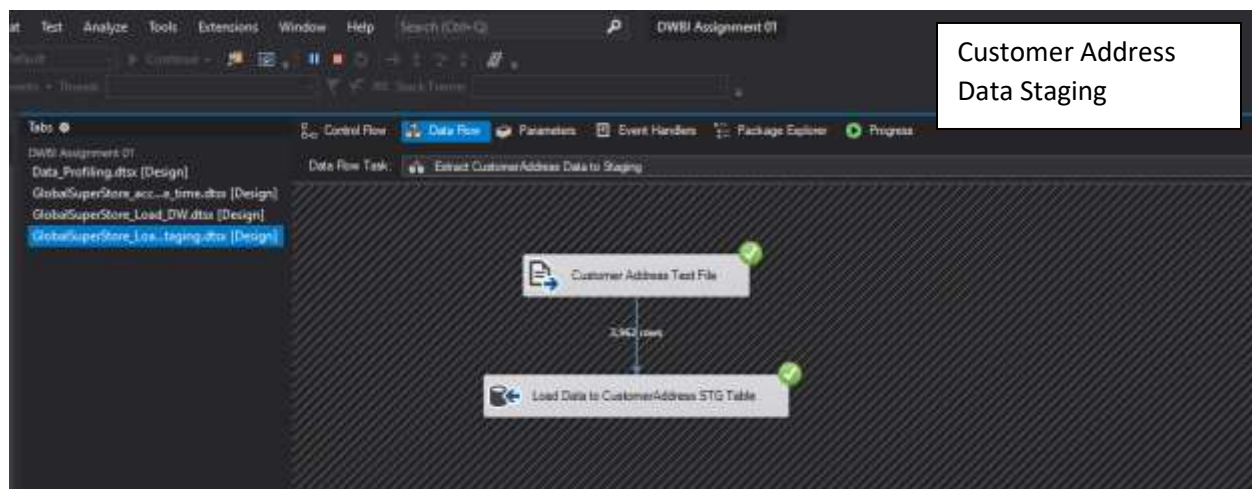
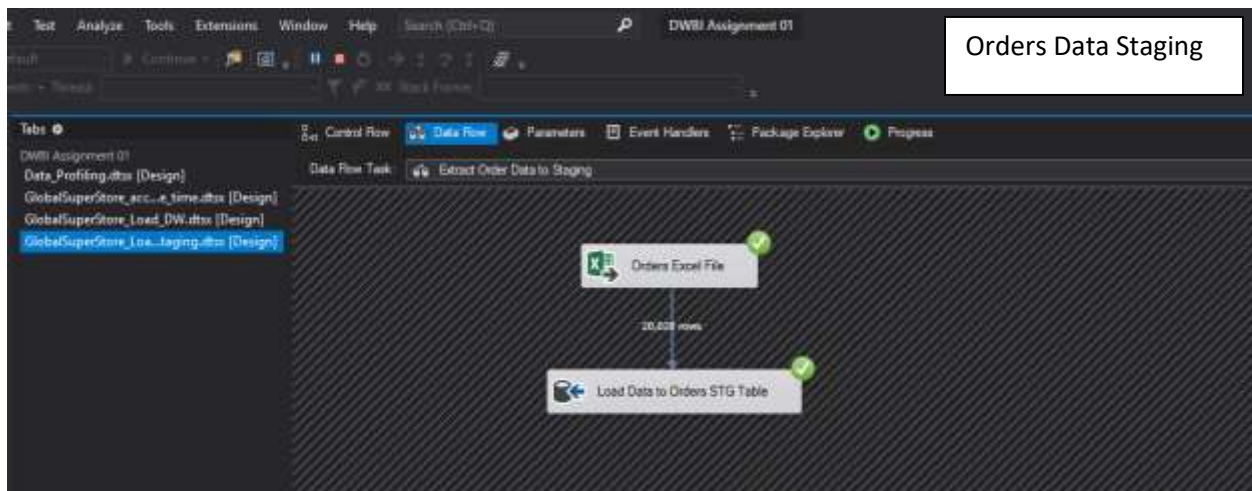
- **Snapshot of SSMS Staging Database**



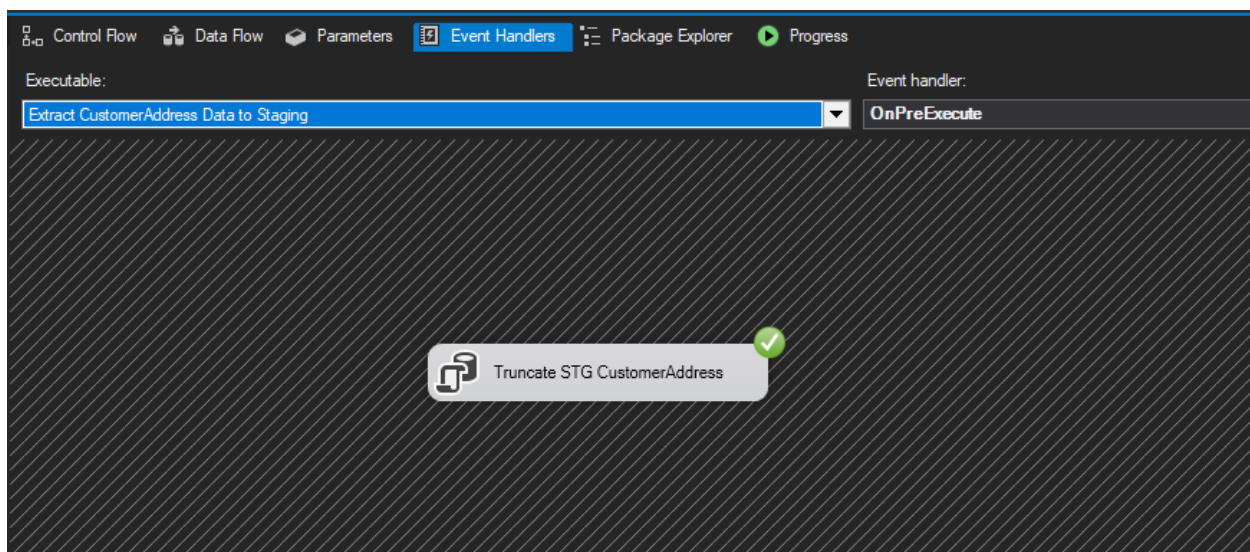
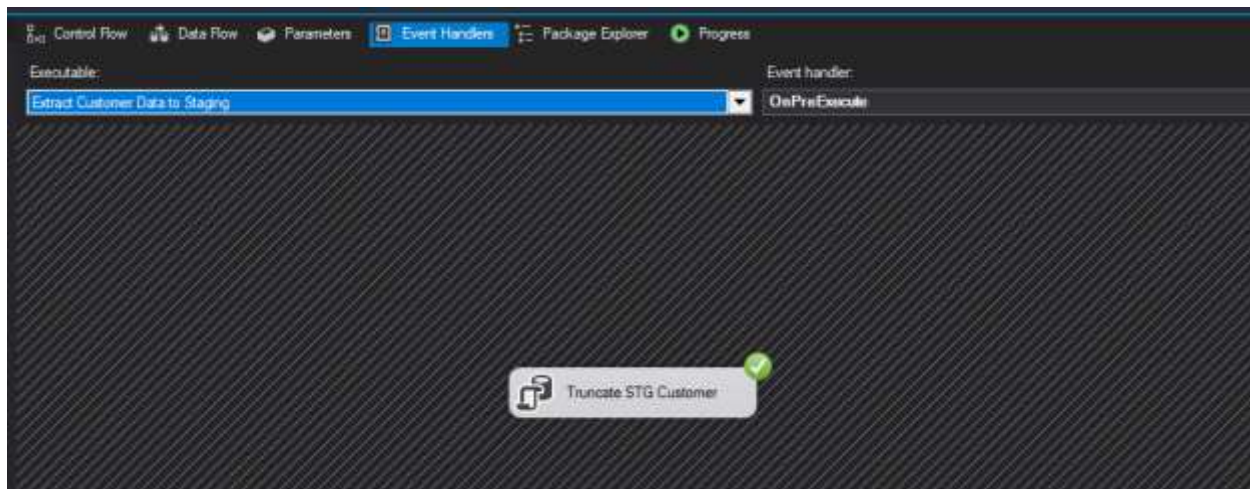
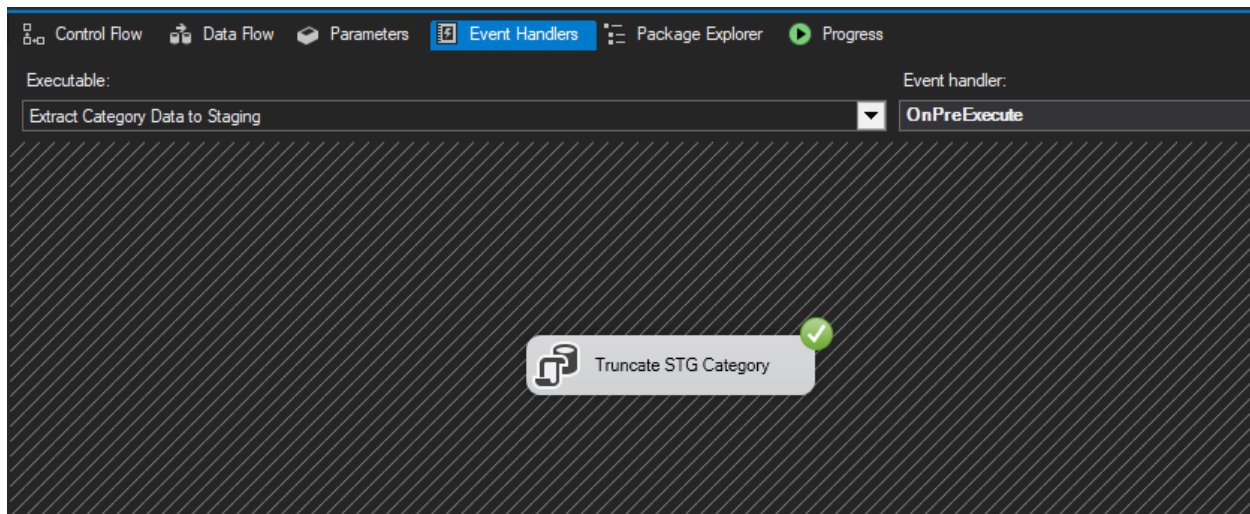
- **Snapshot of Visual Studio Control Flow of Extract**

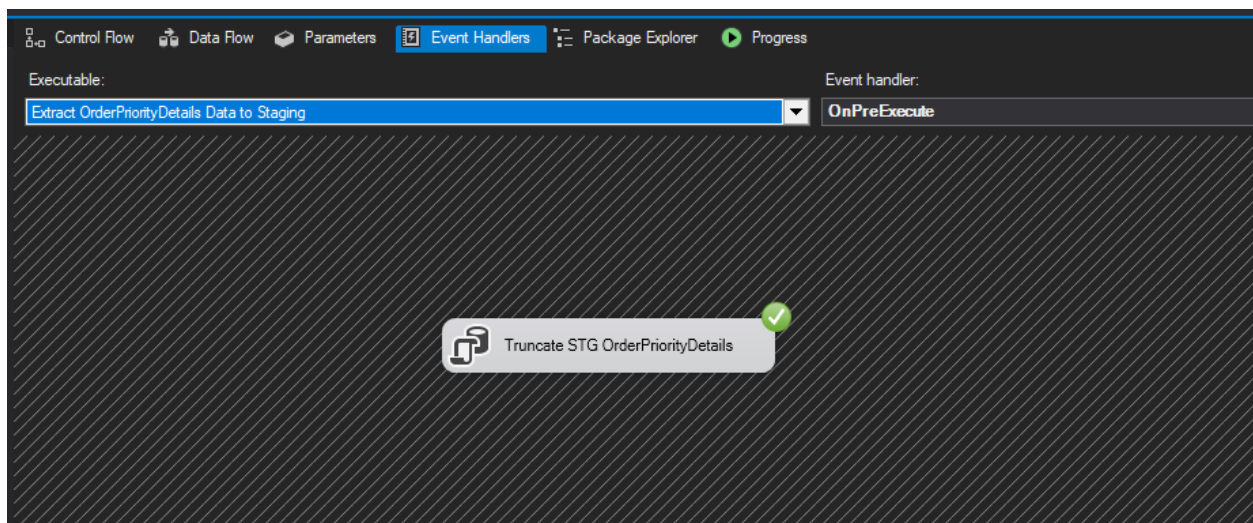
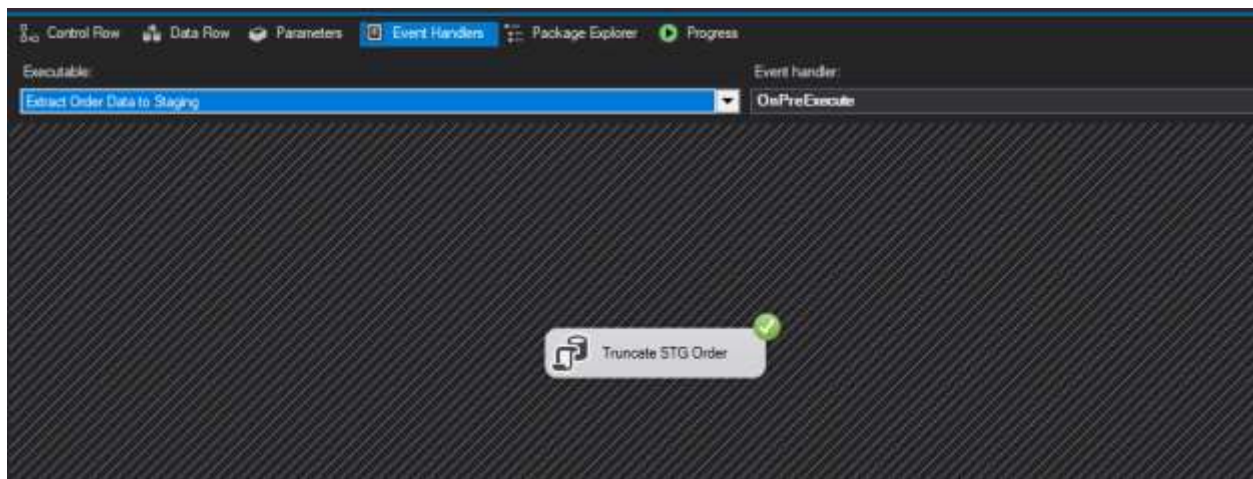
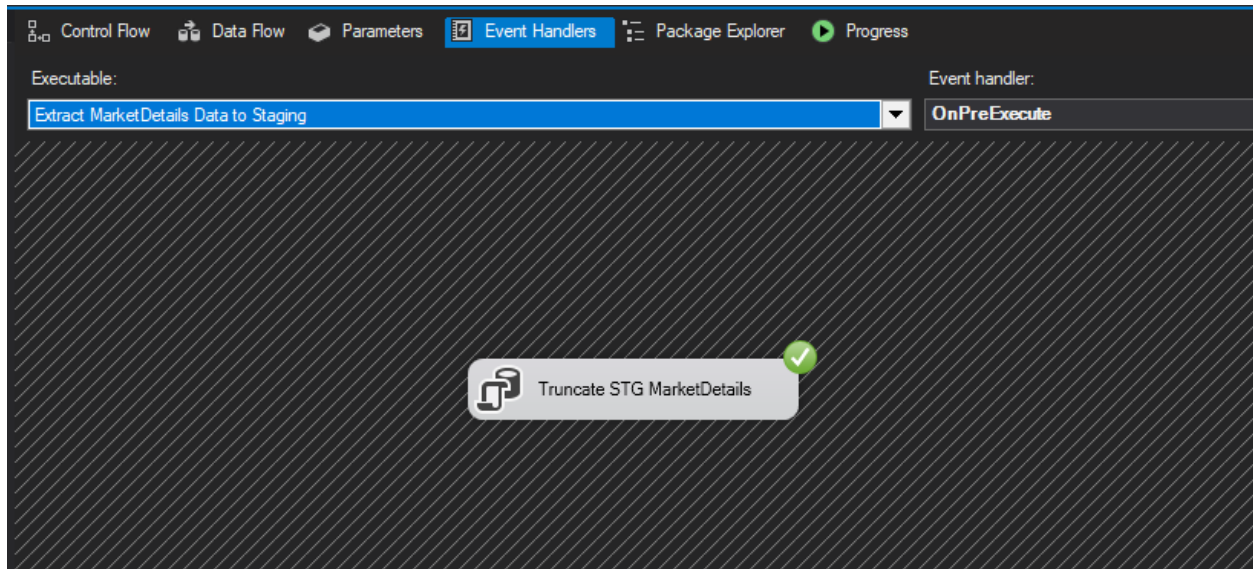


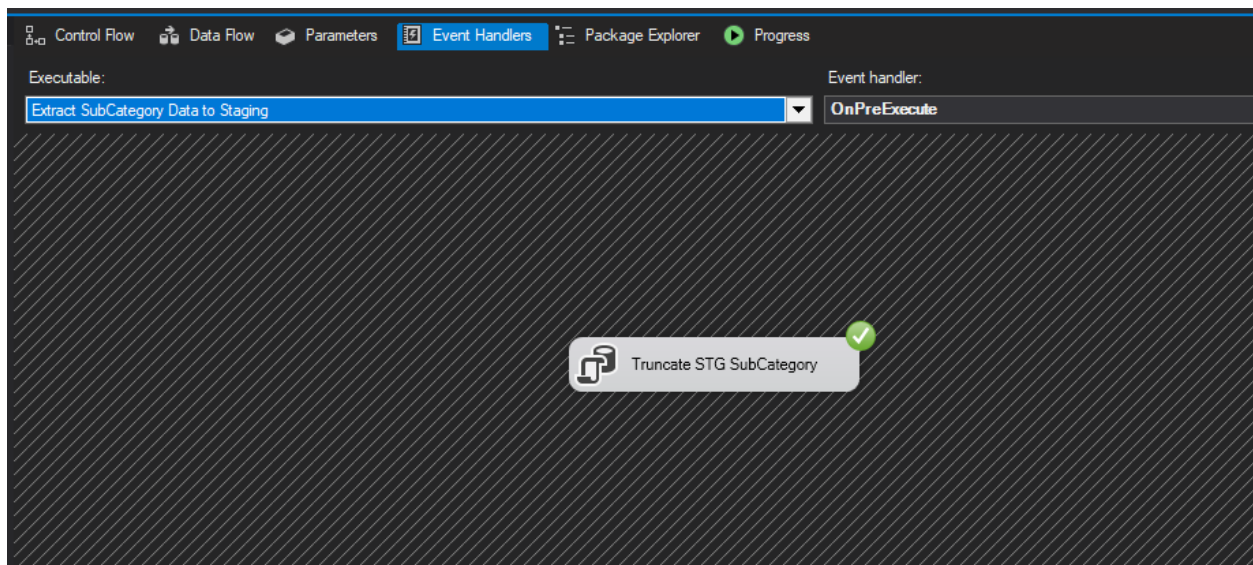
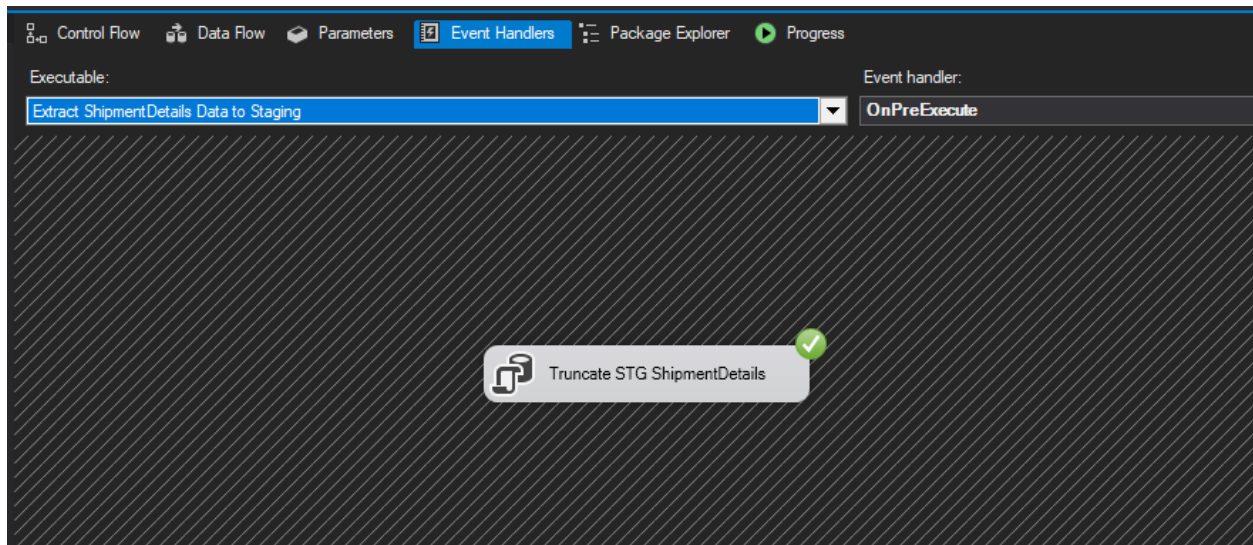
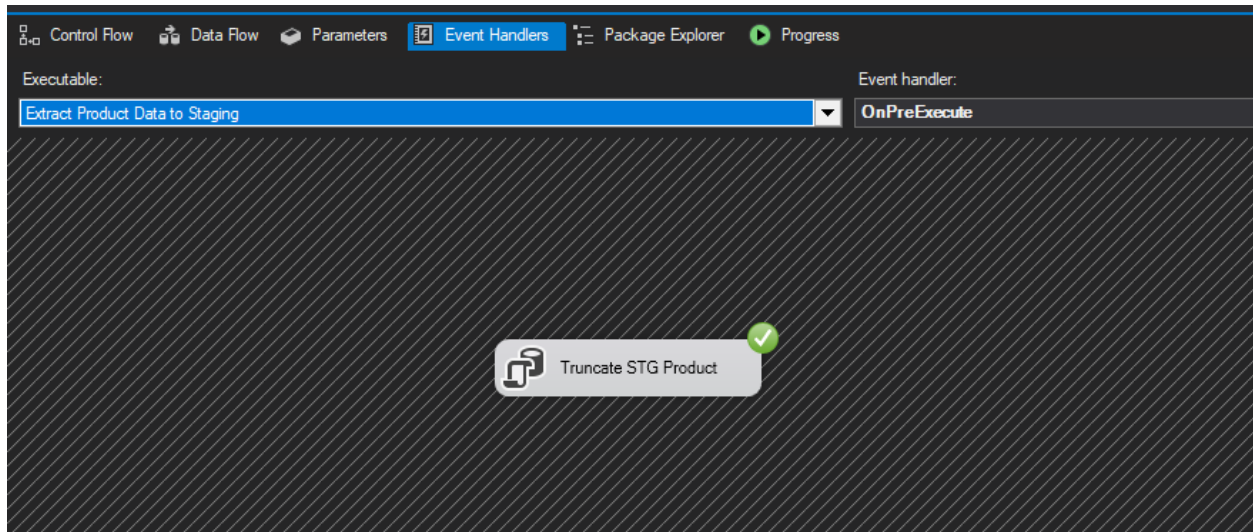
- **Snapshots of several data types of Data Flows**



- **Event Handling (Truncate Staging Data)**

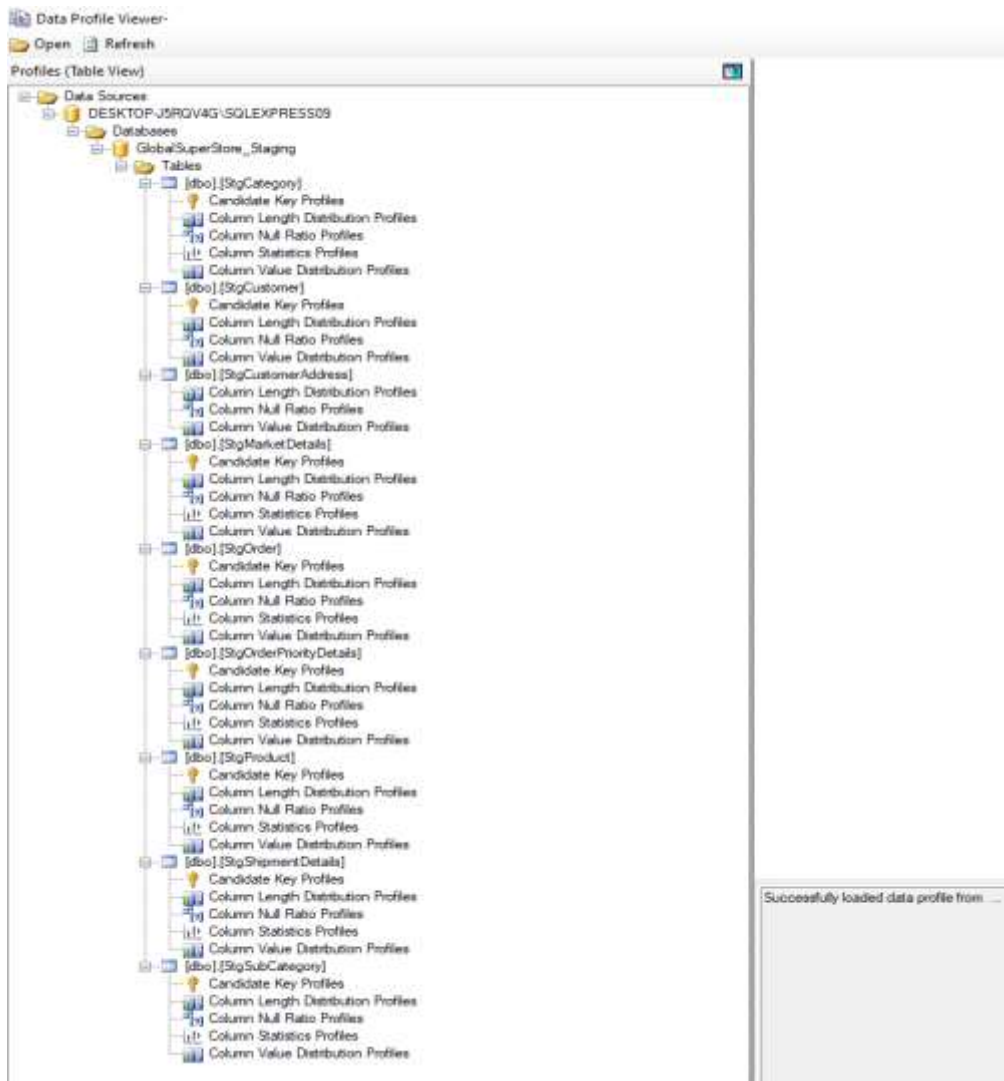
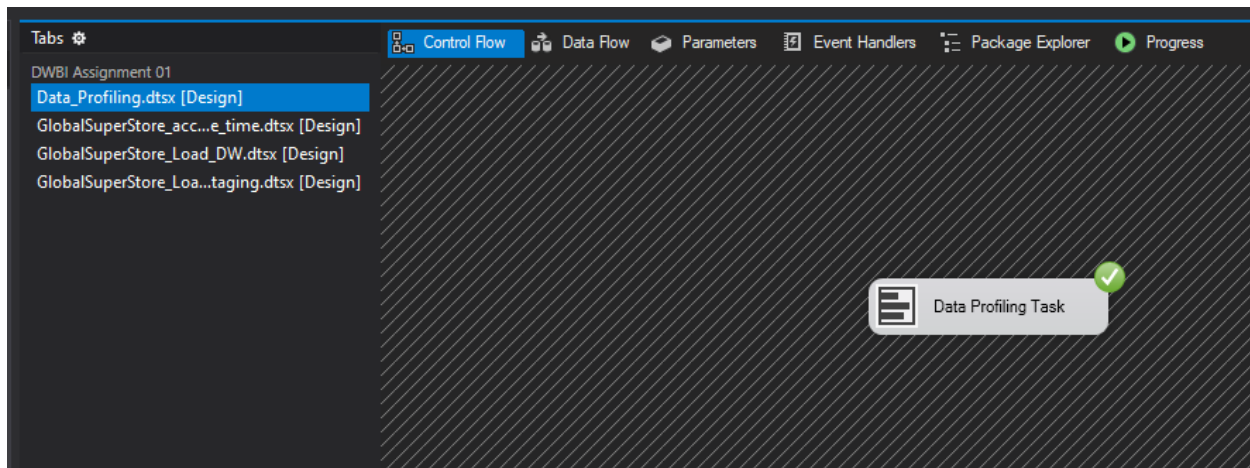






- **Data profiling**

Used Data_Profiling package to profiling the staging tables



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, GlobalSuperStore_DW, which contains the below tables.

1. DimCategory
2. DimCustomer
3. DimDate
4. DimMarketDetails
5. DimOrderPriorityDetails
6. DimProduct
7. DimShipmentDetails
8. DimSubCategory
9. FactOrders

Used Transformation Tasks

1. Lookups

DimCustomer's CustomerID is looked when loading using DimOrder
DimProduct's ProductSK is looked when loading using DimMarket

2. Derived Columns

Derived column is used in FactOrders to derive both StartDate and EndDate by using GETDATE() expression and to derive the FinalCost too.

3. Union

Union is used in the Extract step to combine and get all the data from data files.

4. Sort and Merge

'Sort' is used in DimProduct to sort out the Product and Category data and they are merged using CategoryID.

'Sort' is used sort out the Product and Sub-Category data and they are merged using SubCategoryID.

Update Functions

- DimMarketDetails

```
SQLQuery7.sql - D:\_FRQ\40\Avt (35) - X
1 USE [GlobalSuperStore_DW]
2 GO
3 /***** Object: StoredProcedure [dbo].[UpdateDimMarketDetails]    Script Date: 5/12/2022 11:38:08 PM *****/
4 SET ANSI_NULLS ON
5 GO
6 SET QUOTED_IDENTIFIER ON
7 GO
8 ALTER PROCEDURE [dbo].[UpdateDimMarketDetails]
9     @MarketID int,
10    @Market nvarchar(50),
11    @Region nvarchar(50)
12 AS BEGIN
13     IF NOT EXISTS (select MarketSK from dbo.DimMarketDetails where AlternateMarketID = @MarketID)
14     BEGIN
15         insert into dbo.DimMarketDetails (AlternateMarketID, Market, Region, InsertDate, ModifiedDate) values (@MarketID, @Market, @Region, GETDATE(), GETDATE())
16     END
17     IF EXISTS (select MarketSK from dbo.DimMarketDetails where AlternateMarketID = @MarketID)
18     BEGIN
19         update dbo.DimMarketDetails set Market = @Market, Region = @Region where AlternateMarketID = @MarketID AND (Market != @Market OR Region != @Region)
20     END
21 END
22
```

- DimProductCategory

```
SQLQuery11.sql - D:\_FRQ\40\Avt (32) - X
1 USE [GlobalSuperStore_DW]
2 GO
3 /***** Object: StoredProcedure [dbo].[UpdateDimProductCategory]    Script Date: 5/12/2022 11:46:57 PM *****/
4 SET ANSI_NULLS ON
5 GO
6 SET QUOTED_IDENTIFIER ON
7 GO
8 ALTER PROCEDURE [dbo].[UpdateDimProductCategory]
9     @ProductCategoryID int,
10    @ProductCategoryName nvarchar(50)
11 AS BEGIN
12     IF NOT EXISTS (select CategorySK from dbo.DimCategory where AlternateCategoryID = @ProductCategoryID)
13     BEGIN
14         insert into dbo.DimCategory (AlternateCategoryID, Category, InsertDate, ModifiedDate) values (@ProductCategoryID, @ProductCategoryName, GETDATE(), GETDATE())
15     END
16     IF EXISTS (select CategorySK from dbo.DimCategory where AlternateCategoryID = @ProductCategoryID)
17     BEGIN
18         update dbo.DimCategory set Category = @ProductCategoryName, ModifiedDate = GETDATE() where AlternateCategoryID = @ProductCategoryID AND (Category != @ProductCategoryName)
19     END
20 END
21
```

- DimProductSubCategory

```
SQLQuery12.sql - D:\_FRQ\40\Avt (33) - X
1 USE [GlobalSuperStore_DW]
2 GO
3 /***** Object: StoredProcedure [dbo].[UpdateDimProductSubCategory]    Script Date: 5/12/2022 11:56:54 PM *****/
4 SET ANSI_NULLS ON
5 GO
6 SET QUOTED_IDENTIFIER ON
7 GO
8 ALTER PROCEDURE [dbo].[UpdateDimProductSubCategory]
9     @ProductSubCategoryID int,
10    @ProductSubCategoryName nvarchar(50)
11 AS BEGIN
12     IF NOT EXISTS (select SubCategorySK from dbo.DimSubCategory where AlternateSubCategoryID = @ProductSubCategoryID)
13     BEGIN
14         insert into dbo.DimSubCategory (AlternateSubCategoryID, SubCategory, InsertDate, ModifiedDate) values (@ProductSubCategoryID, @ProductSubCategoryName, GETDATE(), GETDATE())
15     END
16     IF EXISTS (select SubCategorySK from dbo.DimSubCategory where AlternateSubCategoryID = @ProductSubCategoryID)
17     BEGIN
18         update dbo.DimSubCategory set SubCategory = @ProductSubCategoryName, ModifiedDate = GETDATE() where AlternateSubCategoryID = @ProductSubCategoryID AND (SubCategory != @ProductSubCategoryName)
19     END
20 END
21
```

- DimShipmentDetails

```

SQLQuery12.sql - D:\_PROD\40\An (31) - * X
1 USE [GlobalSuperStore_DW]
2 GO
3 /***** Object: StoredProcedure [dbo].[UpdateDimShipmentDetails]    Script Date: 5/13/2022 12:04:14 AM *****/
4 SET ANSI_NULLS ON
5 GO
6 SET QUOTED_IDENTIFIER ON
7 GO
8 ALTER PROCEDURE [dbo].[UpdateDimShipmentDetails]
9     @ShipmentID int,
10    @ShipmentNode nvarchar(50),
11    @ShipmentDescription nvarchar(100)
12 AS
13 BEGIN
14     IF NOT EXISTS (select ShipmentSK from dbo.DimShipmentDetails where AlternateShipmentID = @ShipmentID)
15     BEGIN
16         insert into dbo.DimShipmentDetails (AlternateShipmentID, ShipmentNode, ShipmentDescription, InsertDate, ModifiedDate)
17         values (@ShipmentID, @ShipmentNode, @ShipmentDescription, GETDATE(), GETDATE())
18     END
19
20 IF EXISTS (select ShipmentSK from dbo.DimShipmentDetails where AlternateShipmentID = @ShipmentID)
21 BEGIN
22     update dbo.DimShipmentDetails set ShipmentNode = @ShipmentNode, ShipmentDescription = @ShipmentDescription, ModifiedDate = GETDATE()
23     where AlternateShipmentID = @ShipmentID AND (ShipmentNode != @ShipmentNode OR ShipmentDescription != @ShipmentDescription)
24 END
25
26 END

```

- DimOrderPriorityDetails

```

SQLQuery13.sql - D:\_PROD\40\An (33) - * X
1 USE [GlobalSuperStore_DW]
2 GO
3 /***** Object: StoredProcedure [dbo].[UpdateDimOrderPriorityDetails]    Script Date: 5/13/2022 12:11:26 AM *****/
4 SET ANSI_NULLS ON
5 GO
6 SET QUOTED_IDENTIFIER ON
7 GO
8 ALTER PROCEDURE [dbo].[UpdateDimOrderPriorityDetails]
9     @OrderPriorityID int,
10    @OrderPriority nvarchar(50)
11 AS
12 BEGIN
13     IF NOT EXISTS (select OrderPrioritySK from dbo.DimOrderPriorityDetails where AlternateOrderPriorityID = @OrderPriorityID)
14     BEGIN
15         insert into dbo.DimOrderPriorityDetails (AlternateOrderPriorityID, OrderPriority, InsertDate, ModifiedDate)
16         values (@OrderPriorityID, @OrderPriority, GETDATE(), GETDATE())
17     END
18
19 IF EXISTS (select OrderPrioritySK from dbo.DimOrderPriorityDetails where AlternateOrderPriorityID = @OrderPriorityID)
20 BEGIN
21     update dbo.DimOrderPriorityDetails set OrderPriority = @OrderPriorityID, ModifiedDate = GETDATE()
22     where AlternateOrderPriorityID = @OrderPriorityID AND (OrderPriority != @OrderPriority)
23 END
24
25 END

```

- FactOrders

```

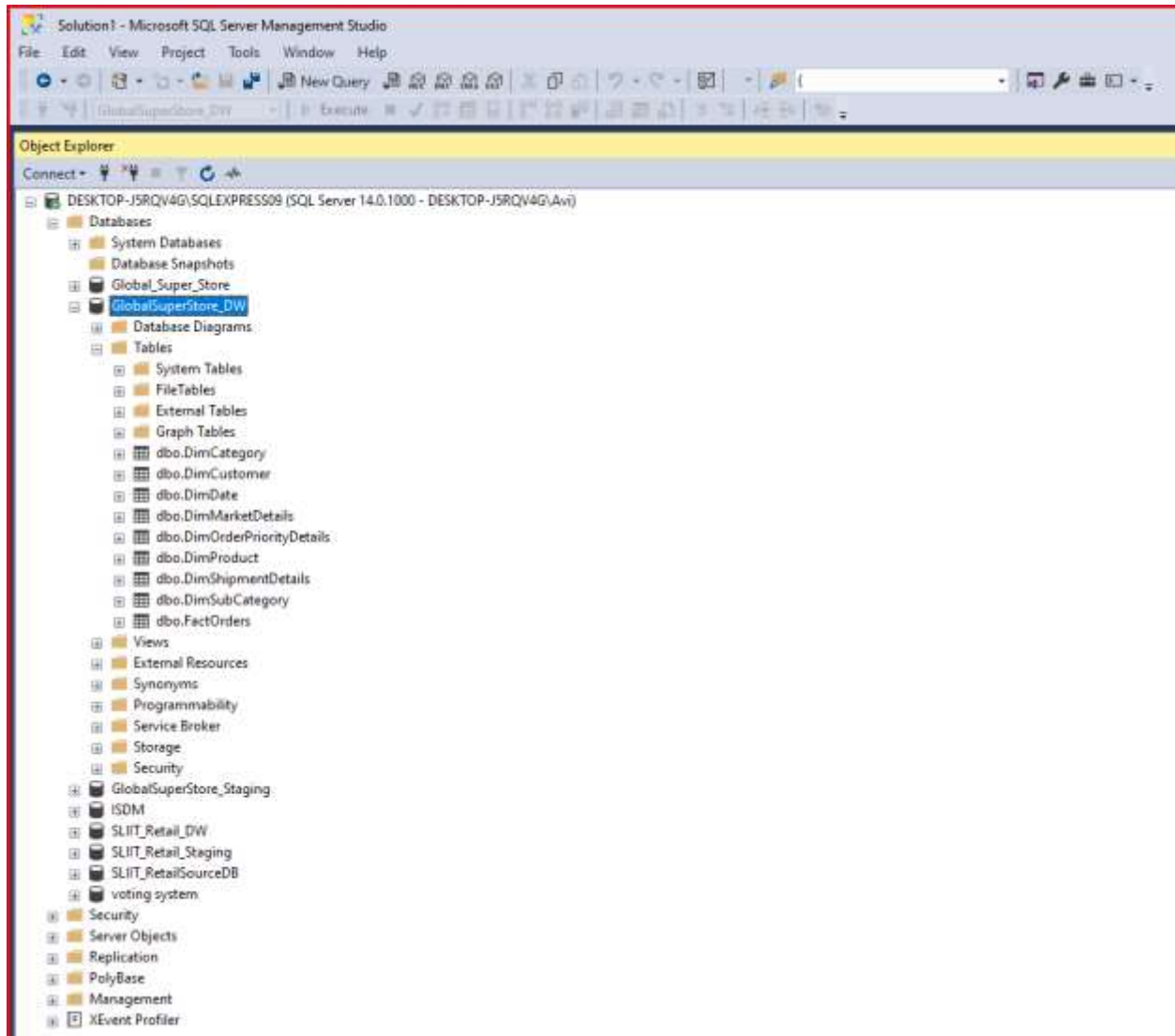
SQLQuery14.sql - D:\_PROD\40\An (35) - * X
1 USE [GlobalSuperStore_DW]
2 GO
3 /***** Object: StoredProcedure [dbo].[UpdateFactOrders]    Script Date: 5/13/2022 12:18:42 AM *****/
4 SET ANSI_NULLS ON
5 GO
6 SET QUOTED_IDENTIFIER ON
7 GO
8 ALTER PROCEDURE [dbo].[UpdateFactOrders]
9     @RowID int,
10    @acct_txn_complete_time datetime,
11    @txn_process_time_hours int
12 AS
13 BEGIN
14     IF EXISTS (select RowID
15     from dbo.FactOrders
16     where RowID = @RowID)
17     BEGIN
18         update dbo.FactOrders
19         set
20         acct_txn_complete_time=@acct_txn_complete_time,
21         txn_process_time_hours=@txn_process_time_hours
22         where RowID = @RowID
23     END
24 END

```

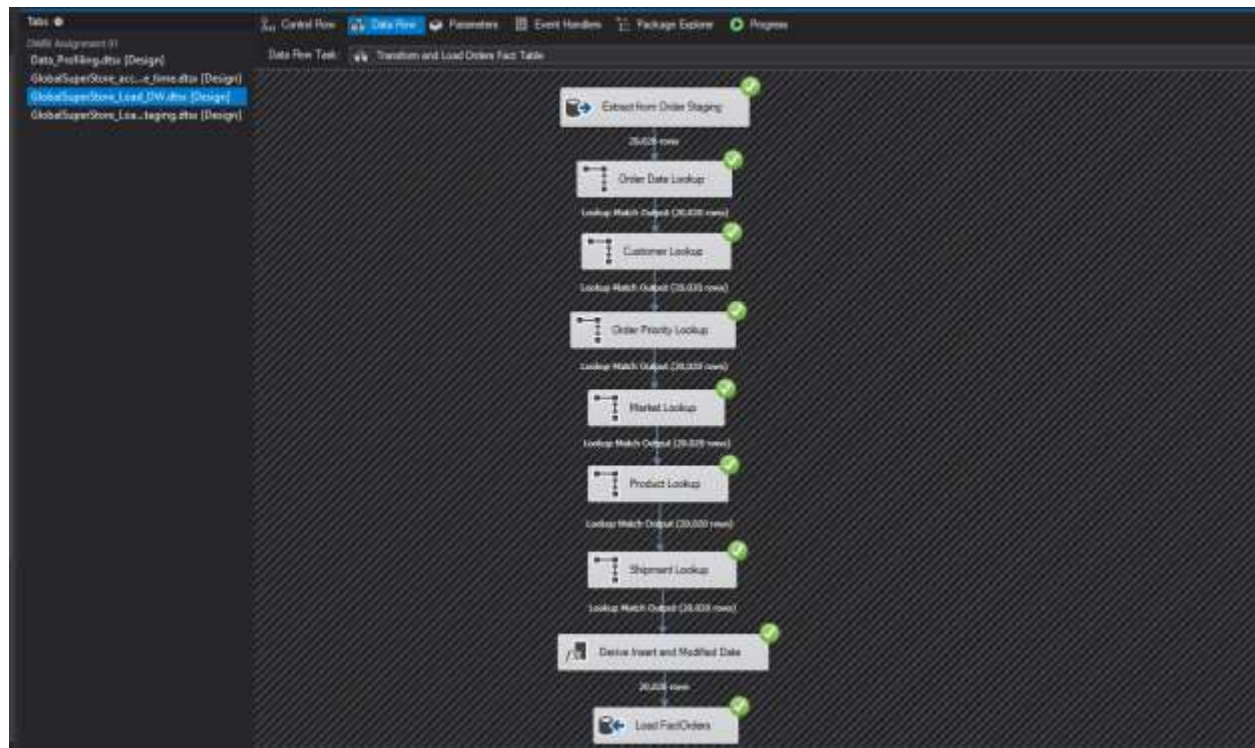
- DimProduct

```
SQLQuery19.sql - D:\-PROJ\40\A\ (M) - >
1 USE [GlobalSuperStore_DW]
2 GO
3 /***** Object: StoredProcedure [dbo].[UpdateDimProduct]    Script Date: 5/15/2022 12:55:05 AM *****/
4 SET ANSI_NULLS ON
5 GO
6 SET QUOTED_IDENTIFIER ON
7 GO
8 ALTER PROCEDURE [dbo].[UpdateDimProduct]
9     @ProductID nvarchar(255),
10     @ProductName nvarchar(99),
11     @ProductCategoryKey int,
12     @ProductSubCategoryKey int
13
14 AS
15 BEGIN
16     IF NOT EXISTS (select ProductSK from dbo.DimProduct where AlternateProductID = @ProductID)
17     BEGIN
18         insert into dbo.DimProduct (AlternateProductID, ProductName, CategoryKey, SubCategoryKey, InsertDate, ModifiedDate)
19         values (@ProductID, @ProductName, @ProductCategoryKey, @ProductSubCategoryKey, GETDATE(), GETDATE())
20     END
21
22     IF EXISTS (select ProductSK from dbo.DimProduct where AlternateProductID = @ProductID)
23     BEGIN
24         update dbo.DimProduct set ProductName = @ProductName, CategoryKey = @ProductCategoryKey, SubCategoryKey = @ProductSubCategoryKey, ModifiedDate = GETDATE()
25         where AlternateProductID = @ProductID AND (ProductName != @ProductName OR CategoryKey != @ProductCategoryKey)
26     END
27
28 END
```

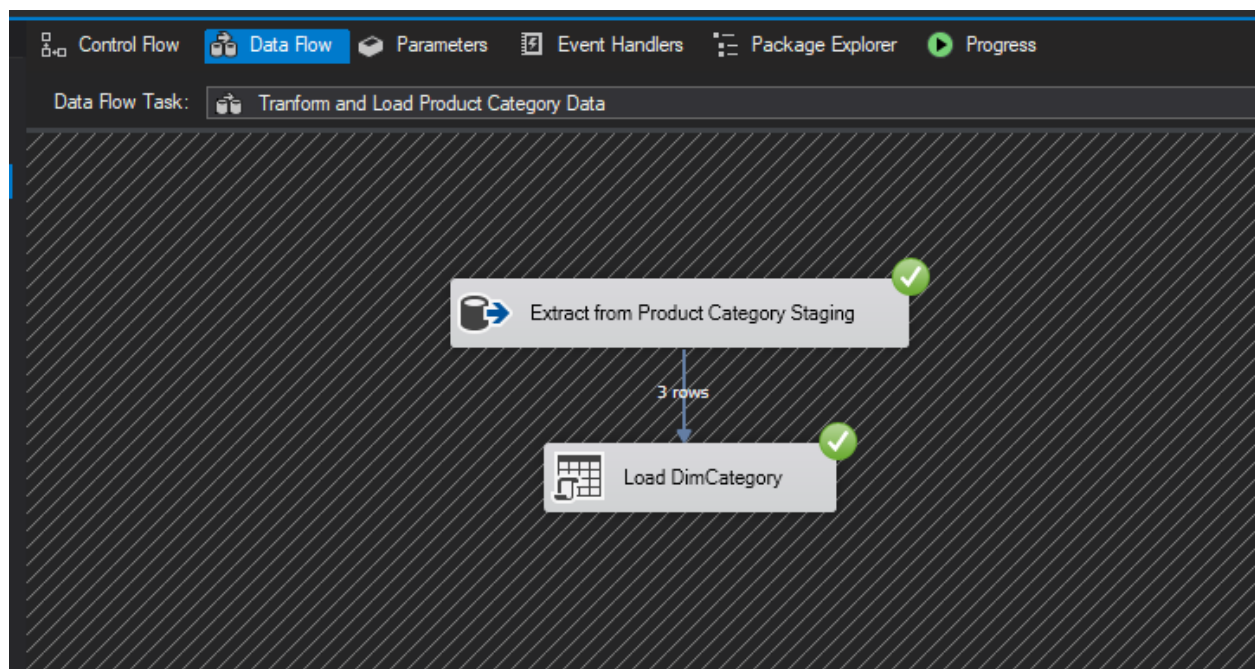
- Snapshot of SQL server Data warehouse Database



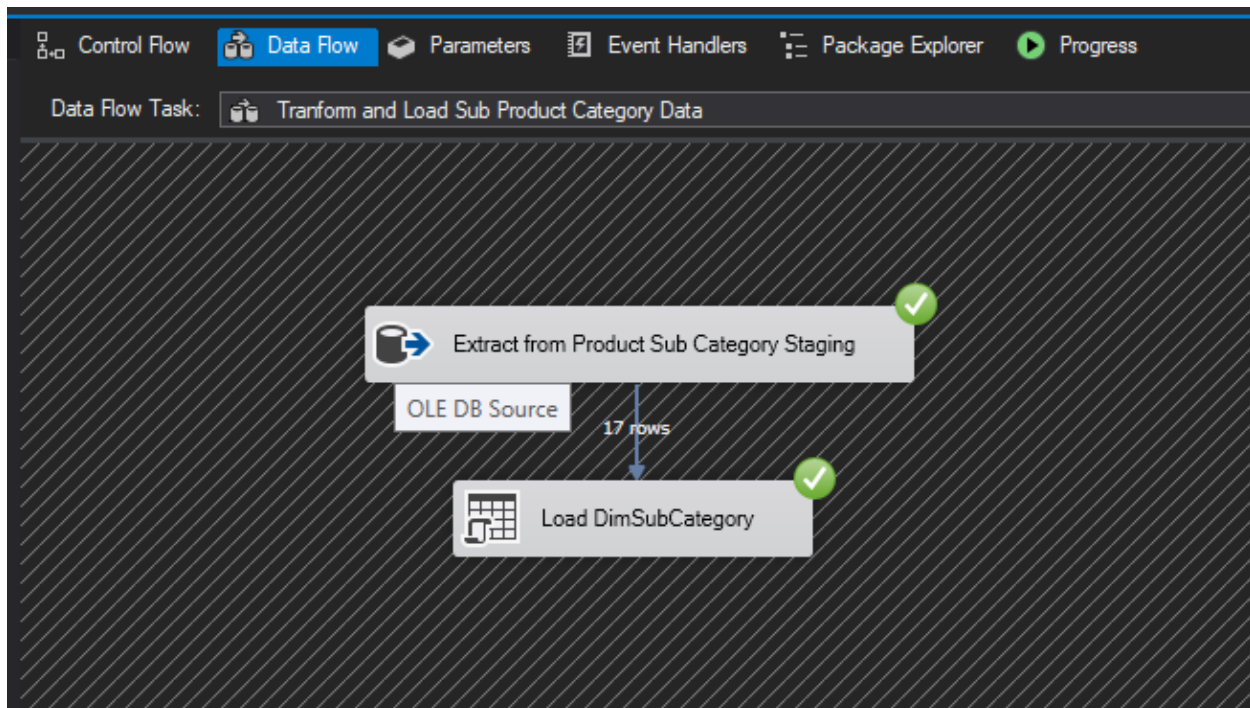
- **Snapshot of Visual Studio Control Flow of Extraction**



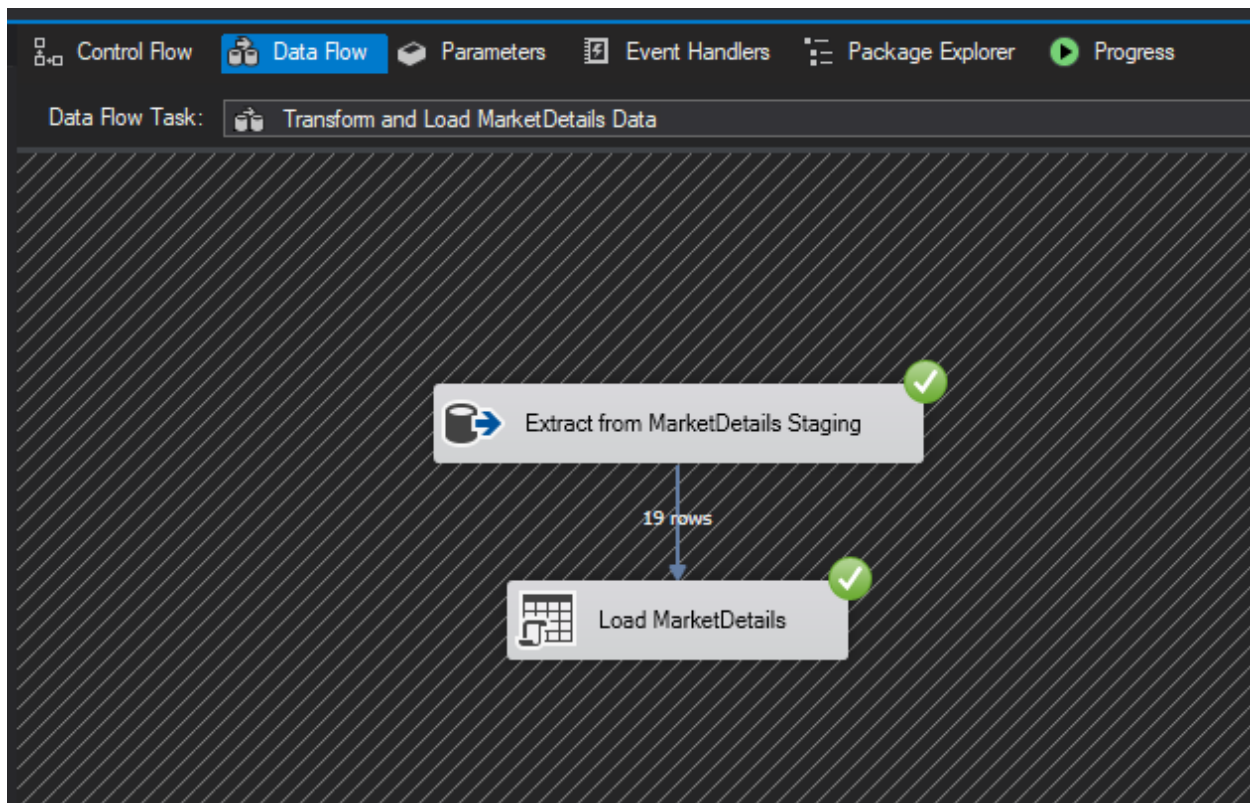
- **Product Category Data Transform and Load**



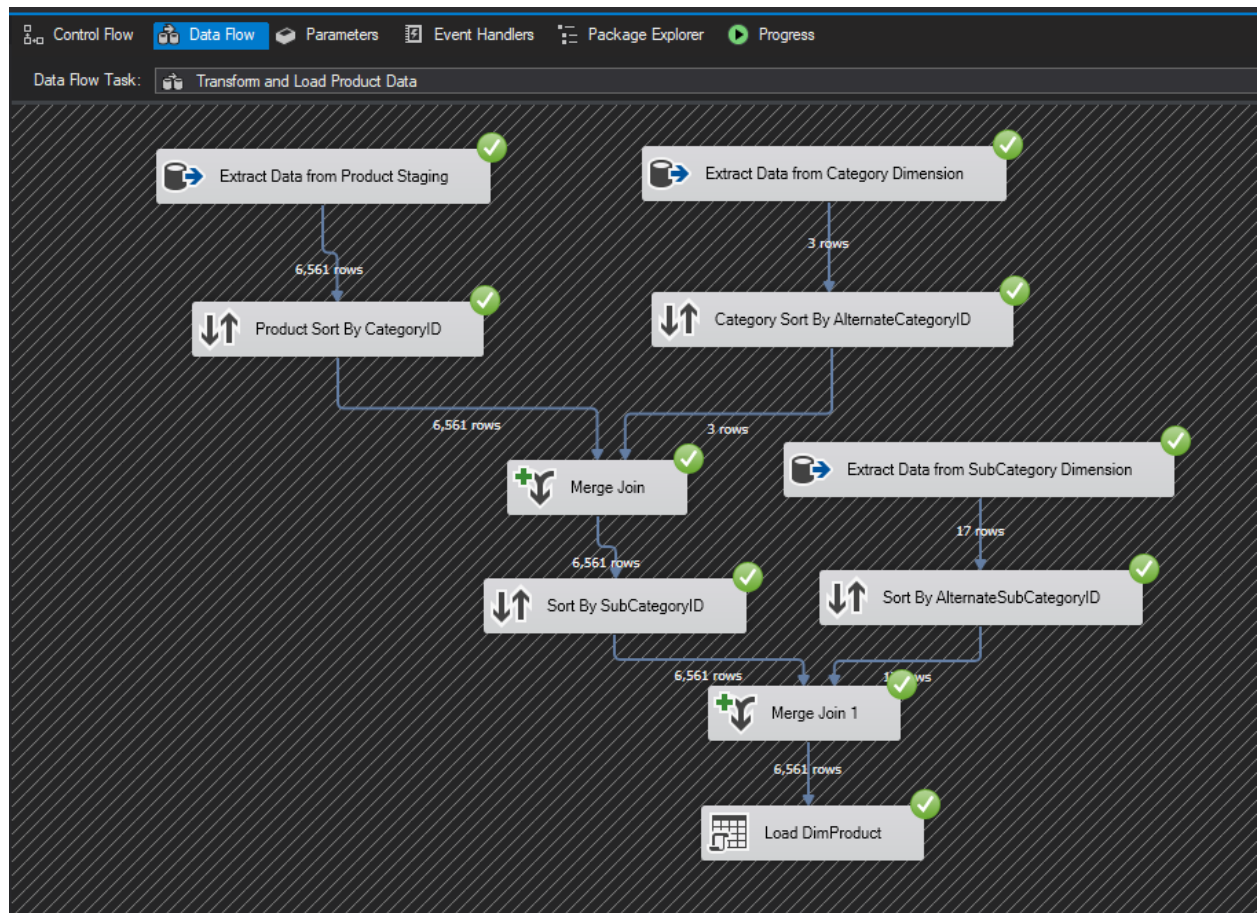
- **Product Subcategory Data Transform and Load**



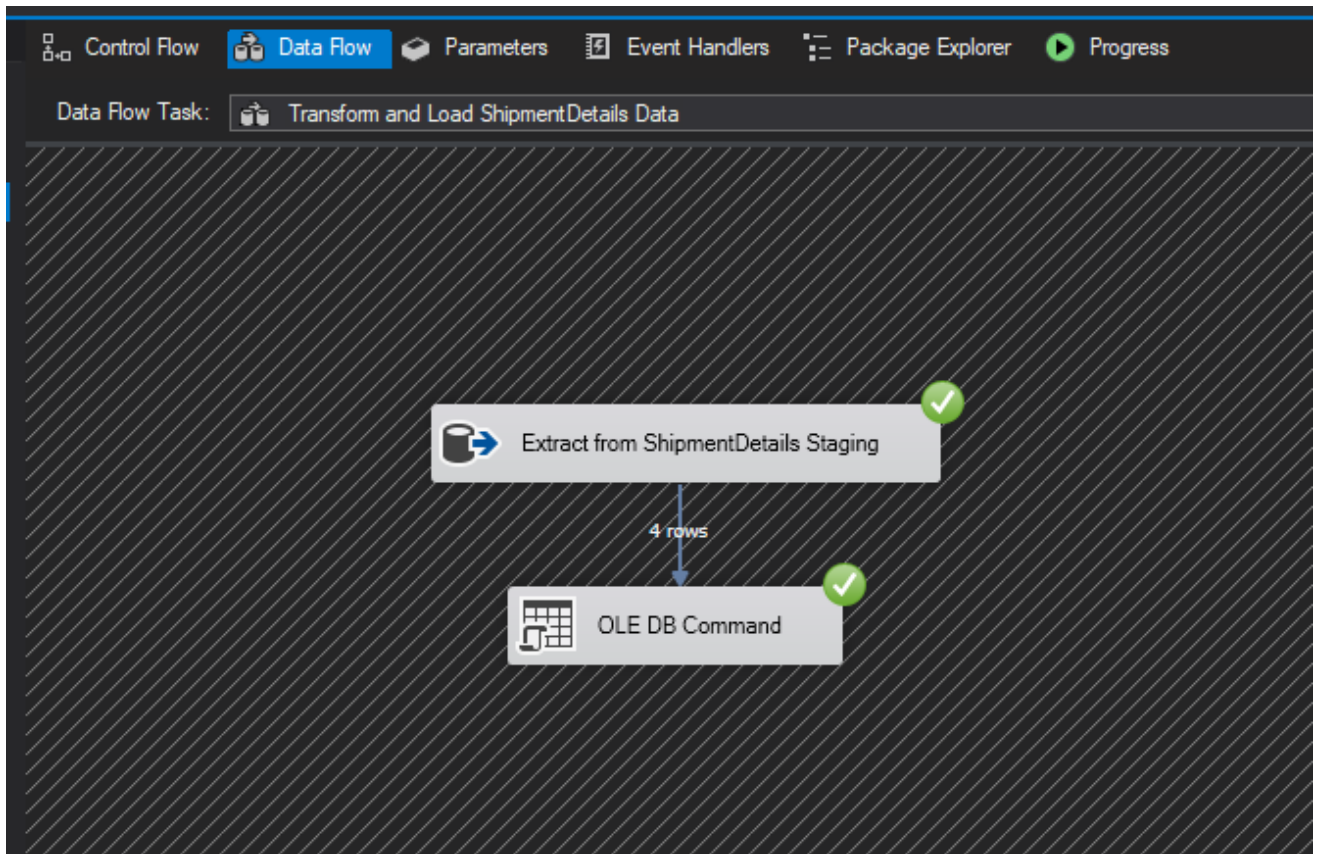
- **Market Details Data Transform and Load**



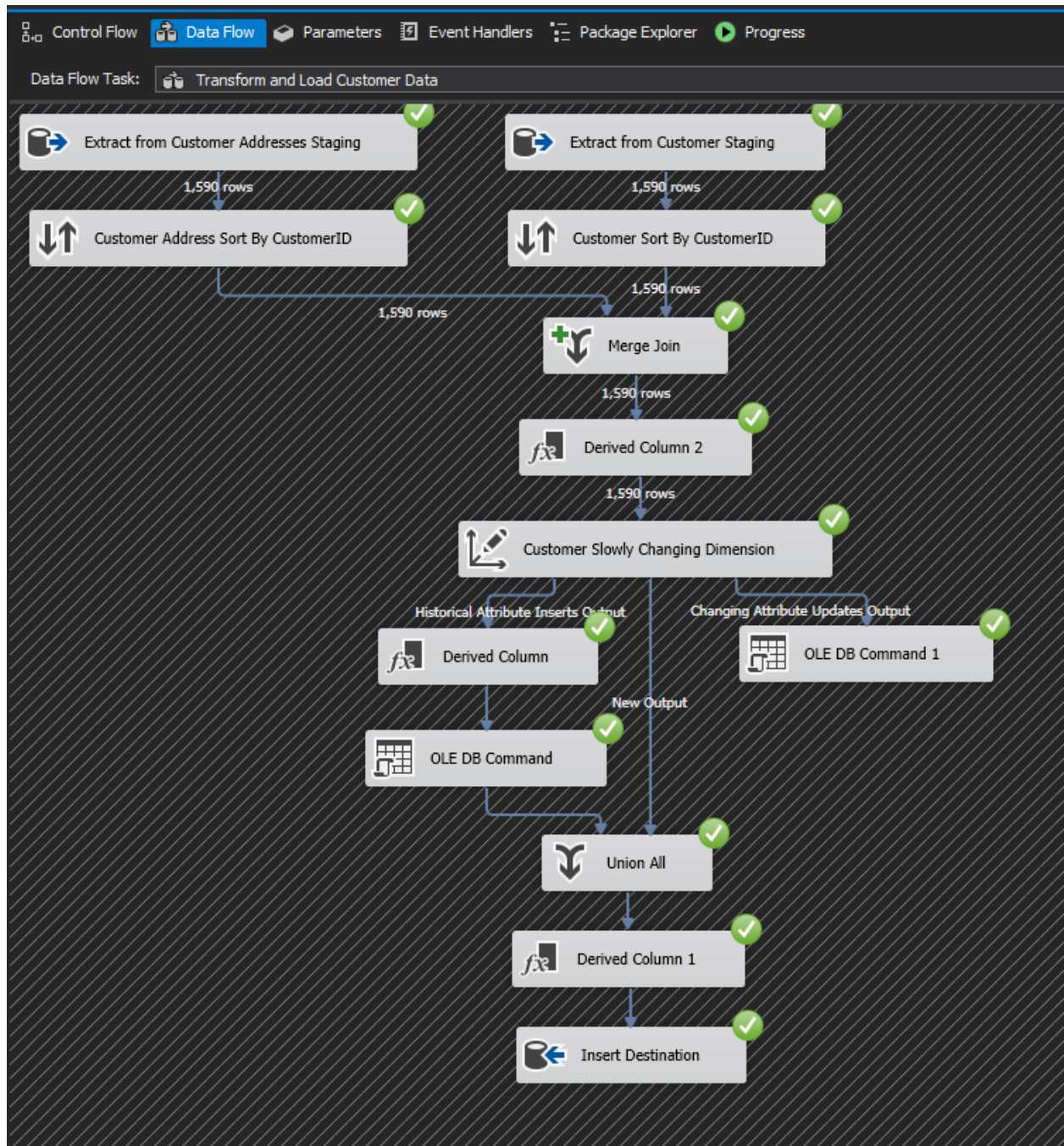
- **Product Data Transform and Load**



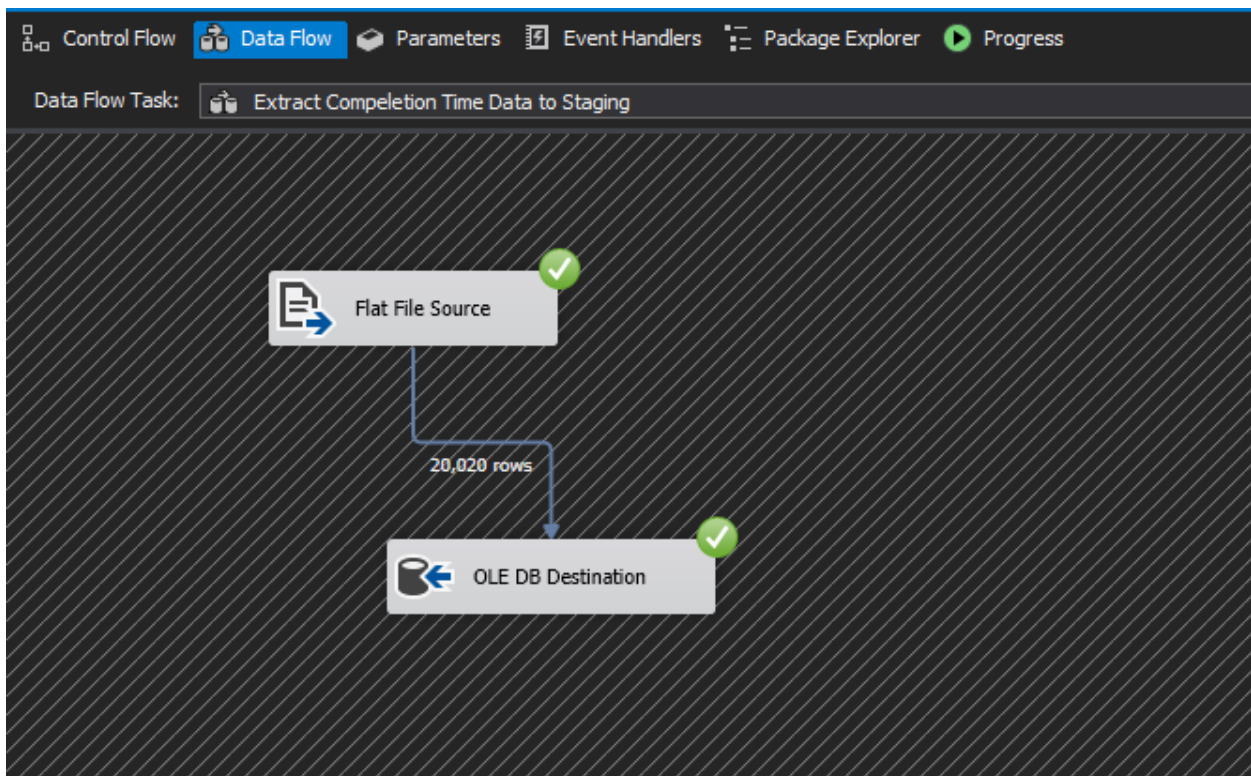
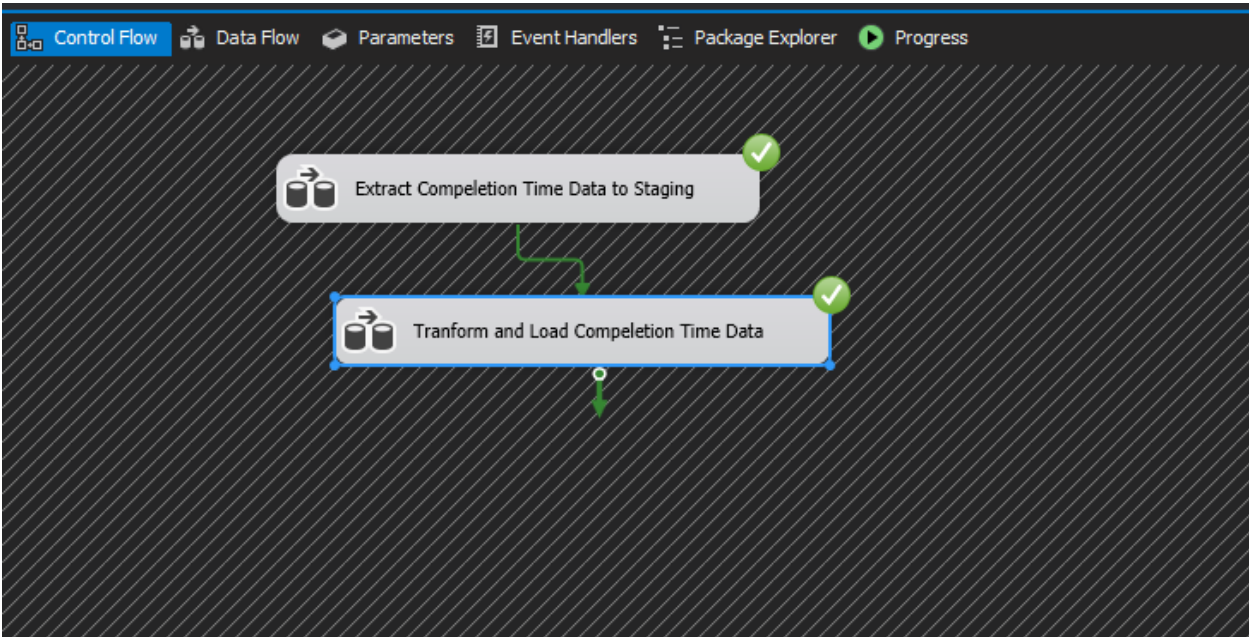
- **Shipment Details Data Transform and Load**



- **Customer Data Transform and Load**



- **Accumulating Fact Table**



Control Flow Data Flow Parameters Event Handlers Package Explorer Progress

Data Flow Task: Transform and Load Completion Time Data

Derived Column Transformation Editor

Specify the expressions used to create new column values, and indicate whether the values update existing columns or populate new columns.

+ Variables and Parameters
+ Columns

+ Mathematical Functions
+ String Functions
+ Date/Time Functions
+ NULL Functions
+ Type Casts
+ Operators

Description:

Derived Column Name	Derived Column	Expression	Data Type
txn_process_time_hours	<add as new column>	DATEDIFF("hh",accm_txn_create_time,accm_txn_complete_time)	four-byte signed in

Configure Error Output... OK Cancel Help