# Data Warehousing & Business Intelligence

# New York City Taxi Trips

Part 1

Developed by:

**D.R.A KUMARAGE IT16058156**

Submitted to:

**MR. Sheron Dinushka**

# Table of Contents

# 1. Data set selection

## Background

This scenario is mainly based the new york city taxi trips database. In this particular scenario Customer, Vehicle and Driver details are not available, Thereby I use data from The Chicago Taxi details for connect to my main database.

Uber,Pickme like taxi services scenario are examples for the above scenario

Customer can go many number of trips , each vehicle Assigned a one driver. Trip detail data base has all trip details

## Content

The data was downloaded By exact link:

Chicago Taxi Trips :

https://data.cityofchicago.org/Transportation/Taxi-Trips/wrvz-psew

new-york-city-taxi-trips

https://www.kaggle.com/kentonnlp/2014-new-york-city-taxi-trips

Driver Details

https://data.cityofchicago.org/Community-Economic-Development/Public-Chauffeurs/97wa-y6ff

# 1. Preparation of Data Sources

In order to data extraction need to prepare the data sources. From my main data source, I extracted three types of data sources.
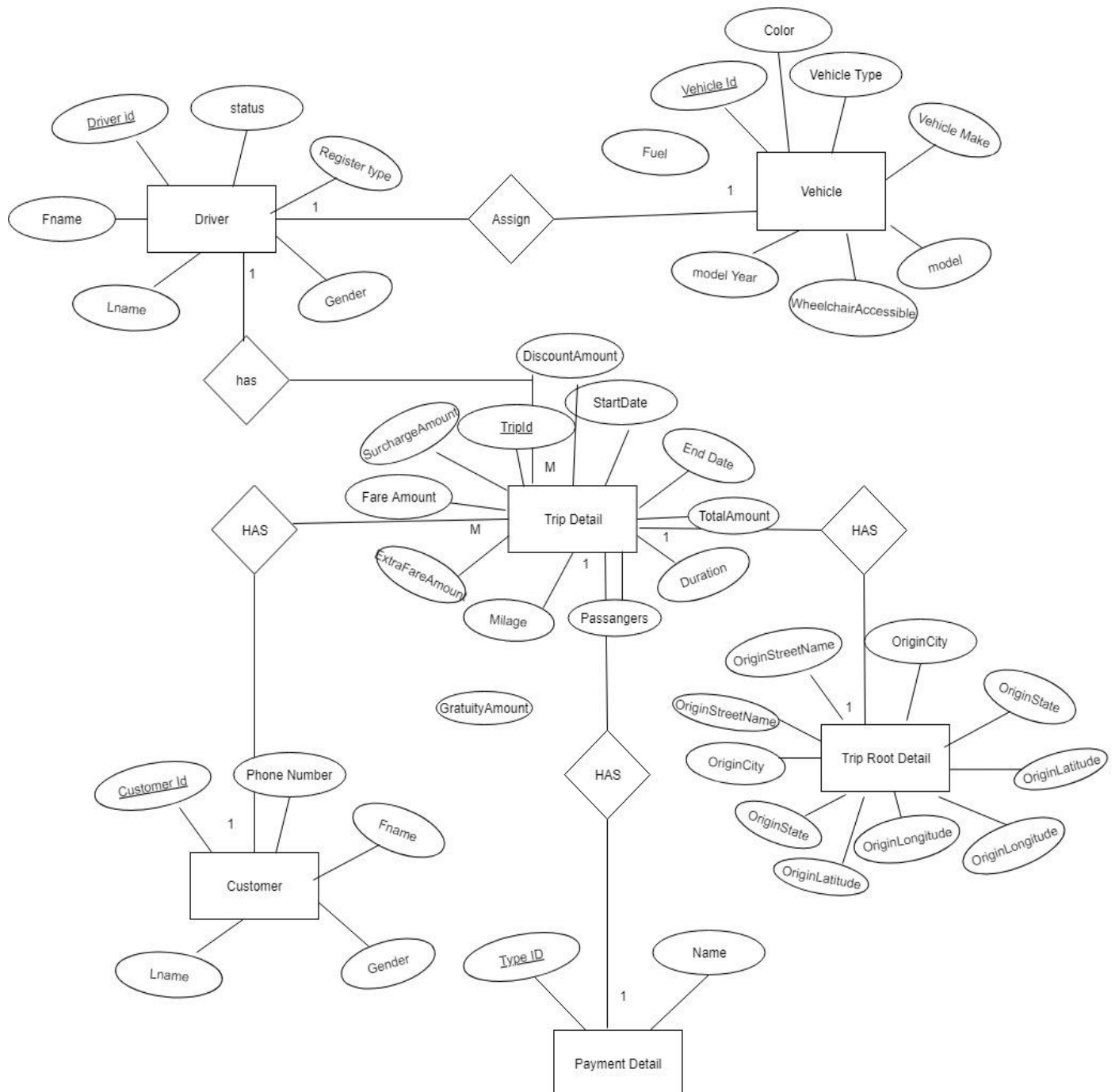
1. Database backup (.bak)

2. Text file (.txt)

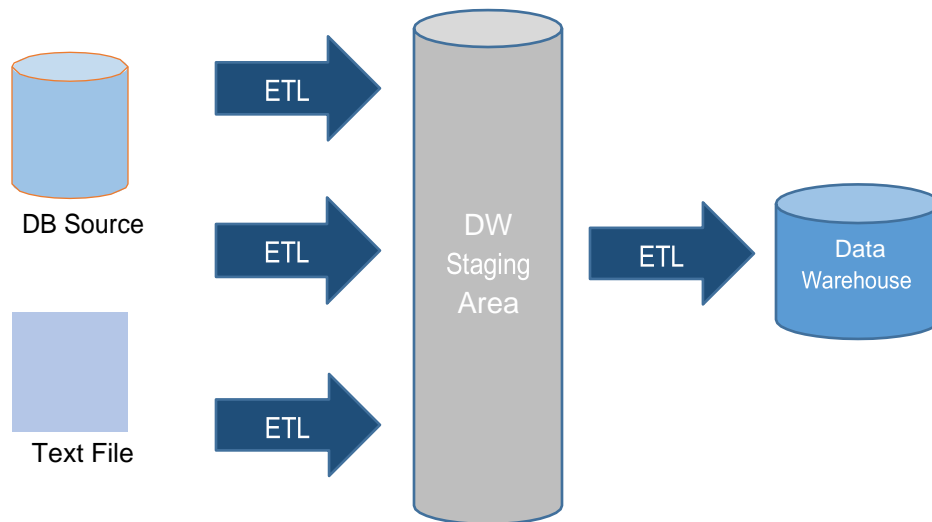3. Excel file (.xlsx)

Text file: Customer Address, Driver Address

Excel file: Customer , Driver , Trip Detail , Trip Root Detail , Vehicle .

# ER Diagram

ER Diaram

# 3. Solution Architecture



- Customer Address Staging
- Customer Staging
- Driver Staging
- Driver Address Staging
- Vehicle Staging
- Payment Type Staging
- Trip Detail Staging
- Trip Root Detail Staging

# Architecture Components

- **Data Sources**

  Operational System (Transaction)

  External sources

- **Extract, Transform, and Load**

  Extract – reading data from source systems

  Transform – Combine data from multiple sources, De-duplicating

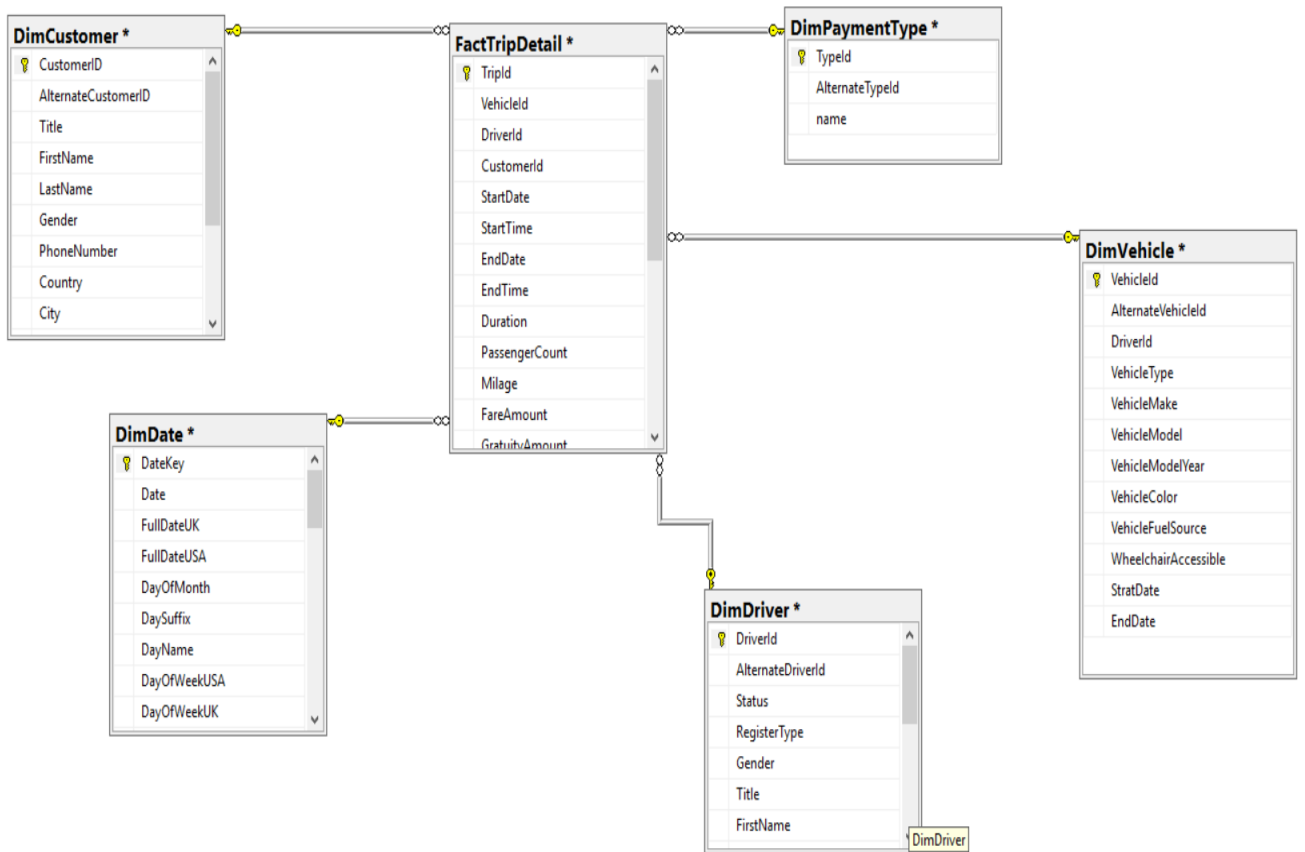  Load – loading data to destination, Surrogate key assignment, Foreign key constraint checks, Indexing

- **Data Warehouse**

  EDW vs Data Mart

  Dimensional Modeling - Facts & Dimension

  Many Schemas - In here used Star schema

## 4. Data warehouse design & development

**DimCustomer ***
| | |
|---|---|
| 🔑 CustomerID | |
| AlternateCustomerID | |
| Title | |
| FirstName | |
| LastName | |
| Gender | |
| PhoneNumber | |
| Country | |
| City | |

**FactTripDetail ***
| | |
|---|---|
| 🔑 TripId | |
| VehicleId | |
| DriverId | |
| CustomerId | |
| StartDate | |
| StartTime | |
| EndDate | |
| EndTime | |
| Duration | |
| PassengerCount | |
| Milage | |
| FareAmount | |
| GratuityAmount | |

**DimPaymentType ***
| | |
|---|---|
| 🔑 TypeId | |
| AlternateTypeId | |
| name | |

**DimVehicle ***
| | |
|---|---|
| 🔑 VehicleId | |
| AlternateVehicleId | |
| DriverId | |
| VehicleType | |
| VehicleMake | |
| VehicleModel | |
| VehicleModelYear | |
| VehicleColor | |
| VehicleFuelSource | |
| WheelchairAccessible | |
| StratDate | |
| EndDate | |

**DimDate ***
| | |
|---|---|
| 🔑 DateKey | |
| Date | |
| FullDateUK | |
| FullDateUSA | |
| DayOfMonth | |
| DaySuffix | |
| DayName | |
| DayOfWeekUSA | |
| DayOfWeekUK | |

**DimDriver ***
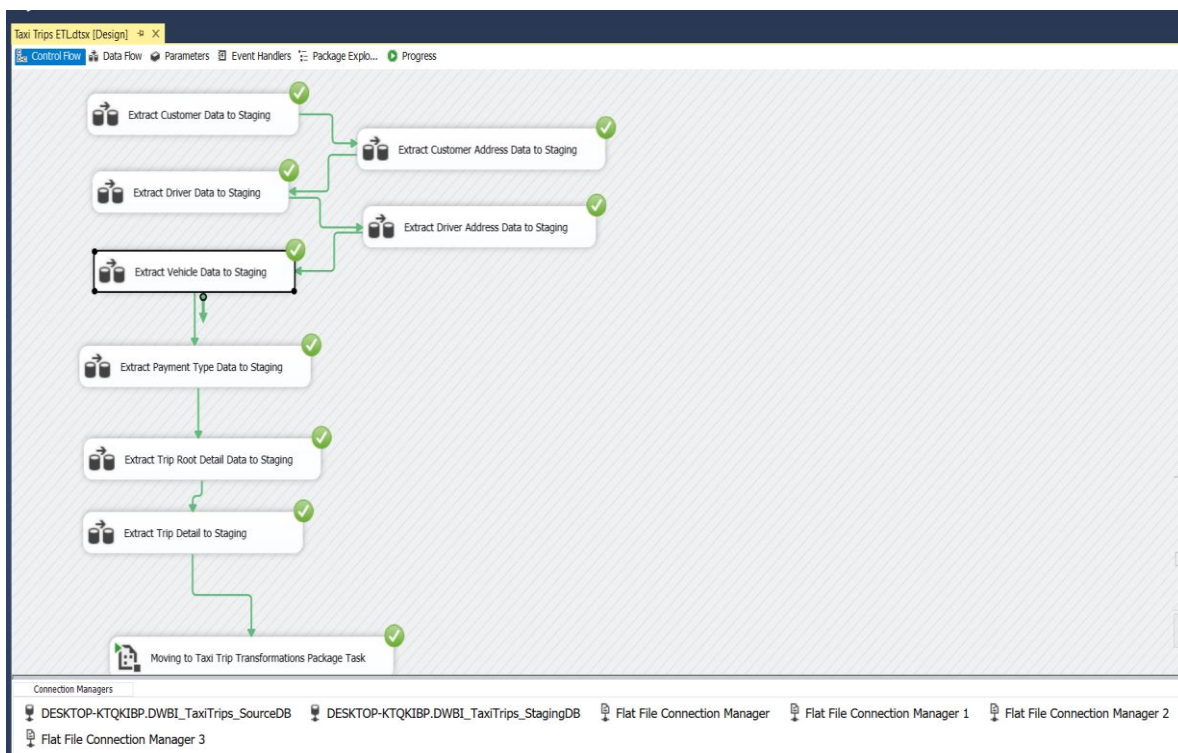| | |
|---|---|
| 🔑 DriverId | |
| AlternateDriverId | |
| Status | |
| RegisterType | |
| Gender | |
| Title | |
| FirstName | |

DimDriver

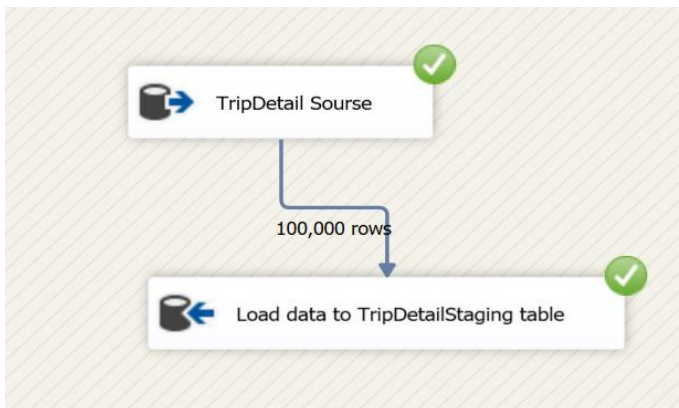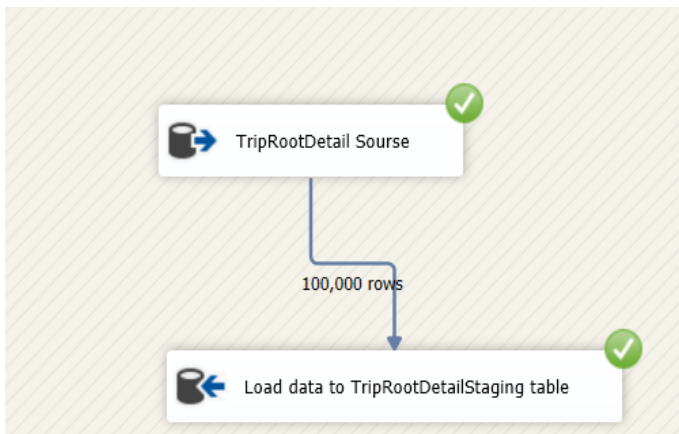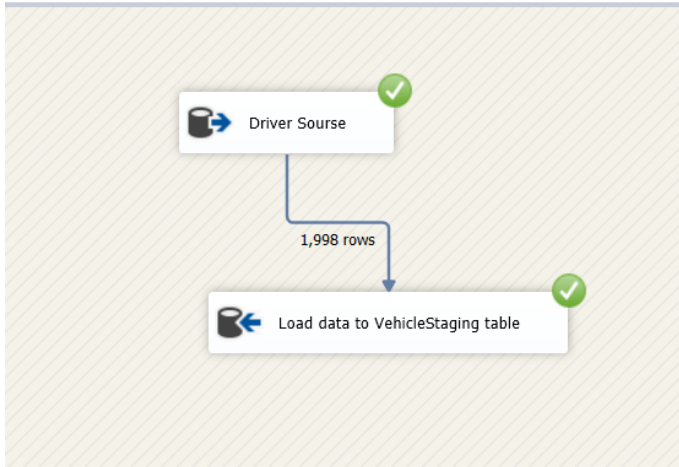## 5. ETL development

### ETL (Extract-Transform-Load)

ETL comes from Data Warehousing and stands for Extract-Transform-Load. ETL covers a process of how the data are loaded from the source system to the data warehouse. Currently, the ETL encompasses a cleaning step as a separate step. The sequence is then Extract-Clean-Transform Load. Let us briefly describe each step of the ETL process.
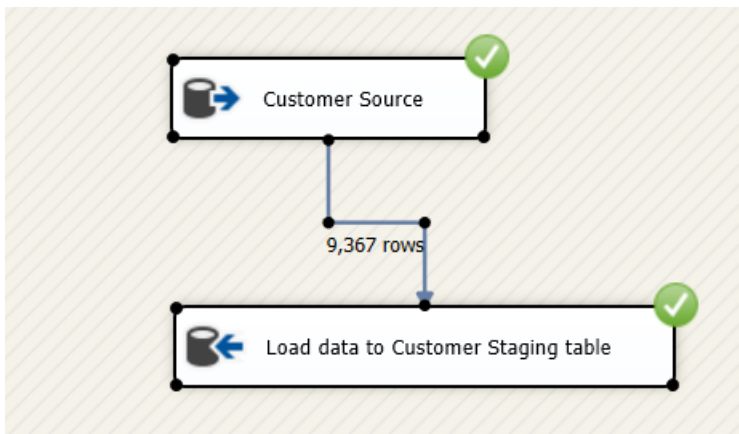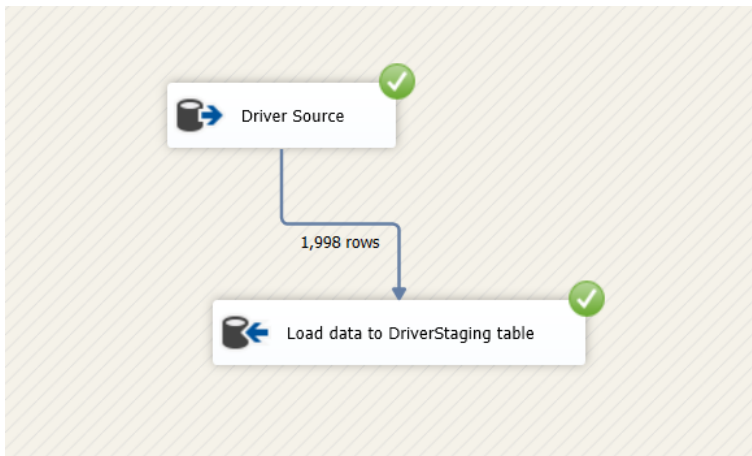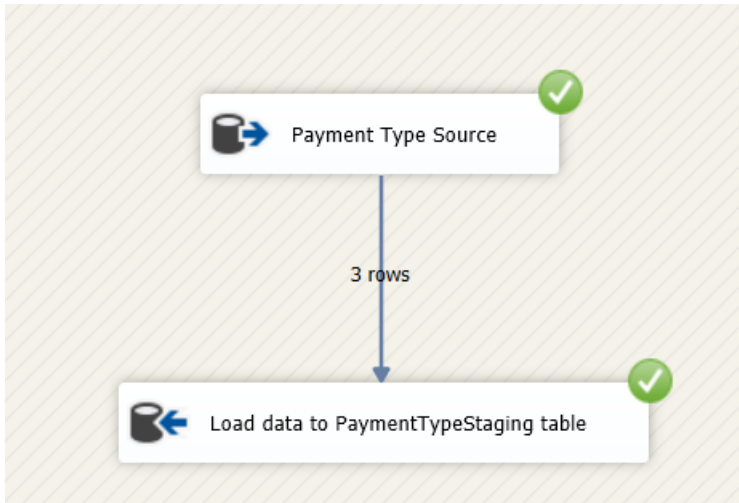
- **Extract**

The main objective of the extract step is to retrieve all the required data from the source system with as little resources as possible.

## Data Extraction from Sources to Staging Tables

Driver Sourse

1,998 rows

Load data to VehicleStaging table



TripRootDetail Sourse

100,000 rows

Load data to TripRootDetailStaging table



TripDetail Sourse

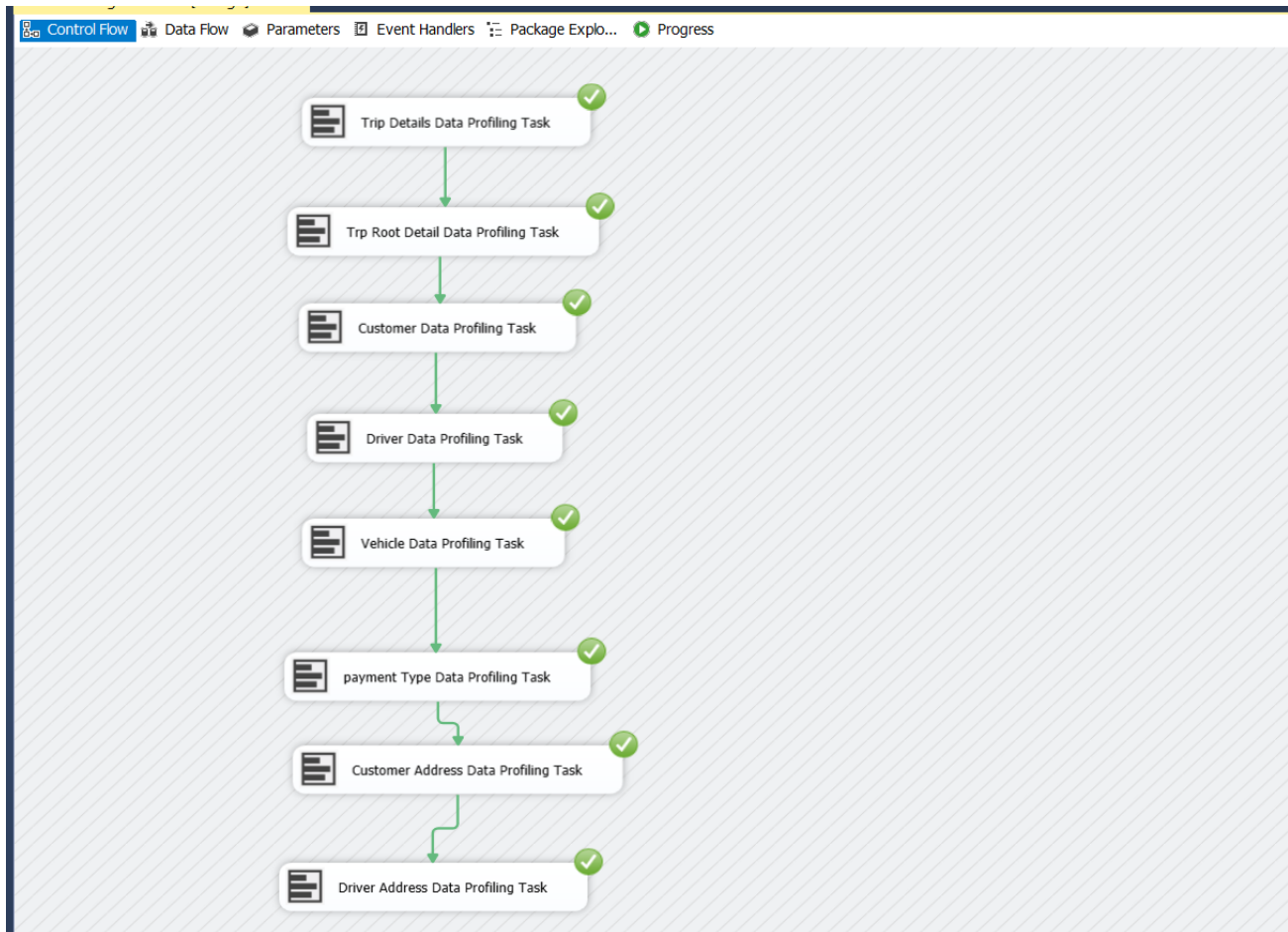100,000 rows

Load data to TripDetailStaging table

**Event Handler In Staging ETL**

EX:

# Data Profiling

## Data Transformation Control Flow



Taxi Trip Transformations.dtsx [Design]

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

Transform Customer Data

Transform Driver Data

Transform Vehical Data

Transform Payment Type Data

Load Data to FactTripDetail

# Data Flow of Loading Dimension and Loading Fact Table



**Taxi Trip Transformations.dtsx [Design]**
Control Flow | Data Flow | Parameters | Event Handlers | Package Explo... | Progress

Data Flow Task: Transform Customer Data

- Extract From CustomerStaging — 9,367 rows
- Extract from CustomerAddressStaging — 9,773 rows
- Sort — 9,367 rows
- Sort 1 — 9,773 rows
- Merge Join — 9,367 rows
- Replace Null Customer Title Column — 9,367 rows
- Slowly Changing Customer Dimension
  - Historical Attribute Inserts Output → Derived Column → OLE DB Command
  - Changing Attribute Updates Output → OLE DB Command 1
  - New Output
- Union All
- Derived Column 1
- Insert Destination



**Taxi Trip Transformations.dtsx [Design]**
Control Flow | Data Flow | Parameters | Event Handlers | Package Explo... | Progress

Data Flow Task: Transform Driver Data

- 1,998 rows
- Sort — 1,998 rows
- Sort 1 — 1,998 rows
- Merge Join — 1,998 rows
- Replace Null in Title Values — 1,998 rows
- Slowly Changing Driver Dimension
  - Historical Attribute Inserts Output → Derived Column → OLE DB Command
  - Changing Attribute Updates Output → OLE DB Command 1
  - New Output
- Union All
- Derived Column 1
- Insert Destination

Load Fact Table :