**Name:** Yuanlong Zhang
**Date:** March. 20, 2022
**Course:** BIDD 320: Data Migration Techniques (Etl Processing)

# Final Project
## ETL Processing Manual

### Version: 1.0

### Yuanlong Zhang

## *Introduction*

This manual describes the essential components of our ETL process. The comprehensive ETL solution includes the following components:

- *Microsoft SSIS to provide visualization, a user-friendly interface, and simplify the ETL processes;*
- *Python for data pre-processing;*
- *MongoDB for reporting file upload;*
- *SQL agent jobs for automation;*
- *Microsoft SSRS for ETL/SSIS job reporting.*

The ETL solution provides some automation for a small medical clinic. It allows the company to automatically import the Patient and visit data from CSV files to the patient database manually import doctor schedule files. All the data will be organized in the centralized data warehouse for other reporting purposes.

**Our ETL Process Manual consists of 7 sections (see below).**

- *You can click the text to jump to the topic you are interested in.*
- *You can also use "bookmark" in the PDF reader to jump to the topic.*

# Contents

*Last Updated: Mar.20, 2022*

# 1. Background and Solution Overview

The ETL solution is designed for a small medical clinic for daily operations.

## 1.1. Background

**Currently**, the business has individual clinics send data to a corporate office by uploading CSV files each day. **Those files are then added to one of the two databases**. The current ETL process is **entirely manual**.
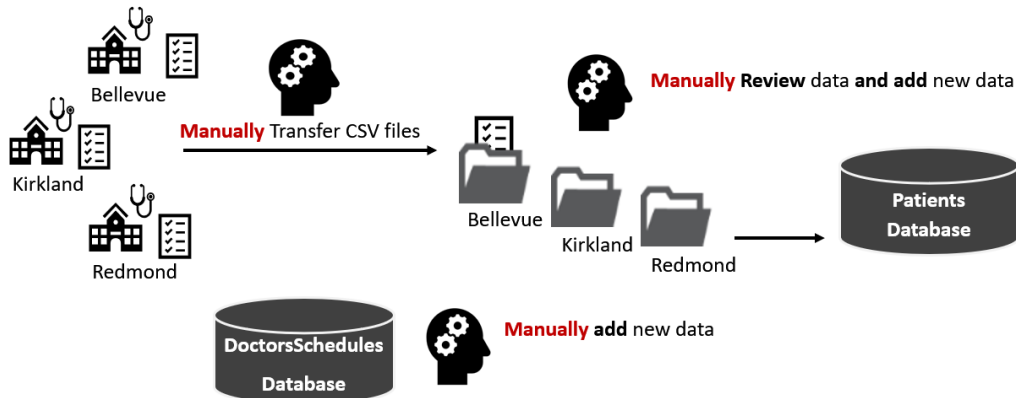


*Figure 1 General Processes of ETL (before)*

The **new** ETL process will **include some automation**. If all goes well, they may fully automate in the future.



*Figure 2  General Processes of ETL (After)*

## 1.2. ETL Solution Overview

Our ETL solutions consist of four milestones:



*Figure 3 Overall Processes of the ETL Solution*

## 1.2.1.   The detailed development processes



*Figure 4 Detailed Steps in the ETL Solution*

## 1.2.2.   The ETL files—SSIS ETL Objects

| Object Name | Milestone | Description | Location |
|---|---|---|---|
| ETLFiletoDatabases.dtsx | 1 | Set of tasks that move data from files to staging tables | _BISolutions\ETLFinal_YuanlongZhang\ETLPackages |
| DWClinicReportDataETL.dtsx | 2 | Set of tasks that move data from database to data warehouse | _BISolutions\ETLFinal_YuanlongZhang\ETLPackages |
| ETLClincReportsDocumentData.dtsx | 3 | Set of tasks that move data from the data warehouse to MongoDB | _BISolutions\ETLFinal_YuanlongZhang\ETLPackages |
| ETLJob.dtsx | 4 | Set of tasks that will call all other three packages for SSIS job | _BISolutions\ETLFinal_YuanlongZhang\ETLPackages |

## 1.2.3.   The ETL files—Non-SSIS ETL Objects

| Object Name | Milestone | Description | Location |
|---|---|---|---|
| PythonETL.py | 1 | Script that moves CSV report data to a Staging Database | _BISolutions\ETLFinal_YuanlongZhang\Scripts |
| MongoDBETL.py | 3 | Script that moves DW data to MongoDB | _BISolutions\ETLFinal_YuanlongZhang\Scripts |
| ETLDetails.rdl | 4 | Paginated report for ETL details information | _BISolutions\ETLFinal_YuanlongZhang\ETLDashboardReports |
| Maindashboard.rdl | 4 | Paginated report for the landing page (high-level summary) | _BISolutions\ETLFinal_YuanlongZhang\ETLDashboardReports |
| SSISJobDetails.rdl | 4 | Paginated report for SSIS Job information | _BISolutions\ETLFinal_YuanlongZhang\ETLDashboardReports |

2. **Checklist for ETL Solutions**  2.1 General Instructions
Yuanlong Zhang

# 2. Checklist for ETL Solutions

To make sure you can use the ETL solutions, we recommend you perform the following steps:

## 2.1. General Instructions

- All the ETL Solution files should be unzipped into C:\_BISolutions\ETLFinal_YuanlongZhang



*Figure 5 Screenshot for the ETL Solution Folder*

- The database files (in \Databases) should be copied to C:\_BISolutions\Databases; you should **restore** the Patients.bak in SSMS. The Patient database will be restored to the server.



*Figure 6 Screenshot for Restore Patient Database*

## 2.2. Checklist for Milestone 1

- Before running Milestone 1 (**ETLFilesToDatabases.dtsx**), please load the updated datafile into C:\_BISolutions\ETLFinal_YuanlongZhang\DataFiles\ClinicDailyData

*5 / 54*

*Figure 7 Screenshot for the data source (csv)*

➕ Make sure you DISABLED the "Pre-load (Databases) Sequence Container **after the first execution**:



*Figure 8 Screenshot for how to disable containers*

## 2.3. *Checklist for Milestone 2*

➕ In Milestone 2 (**DWClinicReportDataETL.dtsx**), before execution, please make sure you have a stable internet connection and have correctly set up the VPN connection to UW:



*Figure 9 Screenshot for how to connect to UW VPN*

➕ The Pre-Load (Data Warehouse) Sequence Container and the Pre-Processing DimDate Table Sequence Container should be run ONLY ONCE. After that, they should be DISABLED.

*Figure 10 Screenshot for how to disable containers*

## 2.4. Checklist for Milestone 3

➕ Before execution of Milestone 3 (**ETLClinicReportsDocumentData.dtsx**), make sure you **have stable internet connection**.

➕ We've included the MongoDB connection string into the python script. *In case you want to use MongoDB Compass*, the connection string is as follows:

mongodb+srv://BICert:BICert@clinicreportsdata.ts9ek.mongodb.net/test?retryWrites=true&w=majority

➕ If you set it up correctly, you should be able to access the MongoDB data through **MongoDB Compass** as follows after executing Milestone 3.



*Figure 11 Screenshot for MongoDB Compass*

## 2.5. *Checklist for Milestone 4*

✦ Before running the **ETLJob.sql**, please make sure you have updated the user name to match your local configuration:



*Figure 12 Screenshot for SQL Script (ETLJob.sql)*

✦ Before executing the **ETLJob.sql**, please delete all the SQL Agent Jobs, Credentials, and SSIS Proxies you've ***created from the other students' script***.

✦ You can execute the SQL Agent Job in SSMS. Before execution, please make sure you run the SQL Agent Service:



*Figure 13 Screenshot for how to start an SQL Agent Job manually*



*Figure 14 Screenshot for how to check the service status for SQL Agent Service*

✦ In case you **have any problems with the SQL Agent Job**, you can **manually execute** the **ETLJob.dtsx**



*Figure 15 Screenshot for Milestone 4: ETLJob.dtsx*

➕ Unless you've correctly installed and configured the SSRS server on your computer, you should directly open the SSRS reports in Visual Studio.

➕ Simply click on "MainDashboard.rdl" and switch to the "Preview" tab:



*Figure 16 Screenshot for MainDashboard SSRS Paginated Report*

➕ You can jump to **SSISJobDetails.rdl** and **ETLDetails.rdl** by clicking "Details."

➕ In the SSIS Job Details dashboard, we ***display all the failure records (if they exist) by default***. If you want to view the **successful one**, use the slider on the top left corner and click "view report" on the top right corner:



*Figure 17 Screenshot for SSISJobDetails SSRS Paginated Report (selected "failure")*



*Figure 18 Screenshot for "view report" button in SSRS preview mode*

# 3. ETL Process Overview

There is 4 ETL package in the solution:

## 3.1.  SSIS for Milestone 1

Here's the overall structure of the first ETL processes designed for Milestone 1: CSV file imports process.



*Figure 19 Screenshot for SSIS Package for Milestone 1*

⊣  **Note:**
- We use several Sequence containers to execute the overall processes.
- To the detailed internal logic, we will discuss in the following part.

## 3.2.  SSIS for Milestone 2

Here's the overall structure of the second ETL process, which was designed for Milestone 2: Data Warehouse ETL.

*Figure 20 Screenshot for SSIS Package for Milestone 2*

## 3.3. SSIS for Milestone 3

Here's the overall structure of the third ETL process, which was designed for Milestone 3: Non-SQL ETL.



*Figure 21 Screenshot for SSIS Package for Milestone 3*

## 3.4.  SSIS for Milestone 4

Here's the overall structure of the fourth ETL process, which was designed for Milestone 4: Automation, reports, and documentation.



*Figure 22 Screenshot for SSIS Package for Milestone 4*

# 4. ETL Process Objects

Our ETL process uses Microsoft SSIS, consisting of 4 package files.

## 4.1. SSIS Package for Milestone 1

Here's the list of files for our ETL processes for Milestone 1:

| | File Name | Description |
|---|---|---|
| 1 | ETLFilesToDatabases.dtsx | Set of tasks that move data from files to staging tables. *Note: _BISolutions\ETLFinal_YuanlongZhang\ETLPackages* |
| 2 | PythonETL.py | Script that moves CSV report data to a Staging Database. *Note: _BISolutions\ETLFinal_YuanlongZhang\Scripts* |
| 3 | ETLLoading.sql | SQL Script for staging table/store procedures definition and transformation to update patient data. *Note: _BISolutions\ETLFinal_YuanlongZhang\SQLScripts* |

## 4.2. ETL Visualization for Milestone 1

We use Microsoft's SQL Server Integration Services (SSIS) to visualize our ETL process. SSIS includes a set of commands represented by visual objects, called tasks, constraints, containers, or transformations. The most common objects used in our ETL process are Sequence Containers, Execute SQL tasks, Precedent Constraints, and Connections.



*Figure 23 Overview of SSIS package: Control Flow*

✦ **Note:**
- These visualization objects are stored and configured in an SSIS Package file.

### 4.2.1. Connections

We use a Connection object that connects to the original data source file, invalid phone number file, and data warehouse that holds our SQL ETL objects.



*Figure 24 Connections for Milestone 1*

We used OLE DB connections; here's the detailed information:

| No. | Name | Description |
|---|---|---|
| 1 | localhost.master | This is an OLE DB connection. Will connect to localhost\master |
| 2 | localhost.patients | This is an OLE DB connection. Will connect to localhost\patients |
| 3 | localhost.staging | This is an OLE DB connection. Will connect to localhost\staging |

### 4.2.2. Pre-Load (Data Warehouse) Sequence Container

There are three components inside the sequence container.



*Figure 25 Pre-Load Sequence Container in Milestone 1*

✦ **Note:**
• This container should only run ONCE.

| No. | Name | Description |
|---|---|---|
| 1 | Initialize Data Warehouse Task | This uses the script to initialize all the database files: It creates the local copy for the patient database. |

| No. | Name | Description |
|-----|------|-------------|
|  |  |  |
| 2 | Initialize Staging Database Task  | Used to create all the staging tables for ETL. The tables are designed to temporarily load the data and transform it from CSV files before loading it into the database. The script is available in ETLLoading.sql  |
| 3 | Initialize all the store procedures task  | Used to create all the views and store procedures for ETL processes. To be detailed, we created views for each table and the store procedures to load the data from views. The script is available in ETLLoading.sql |

<table>
<tr><th>No.</th><th>Name</th><th>Description</th></tr>
<tr><td></td><td></td><td>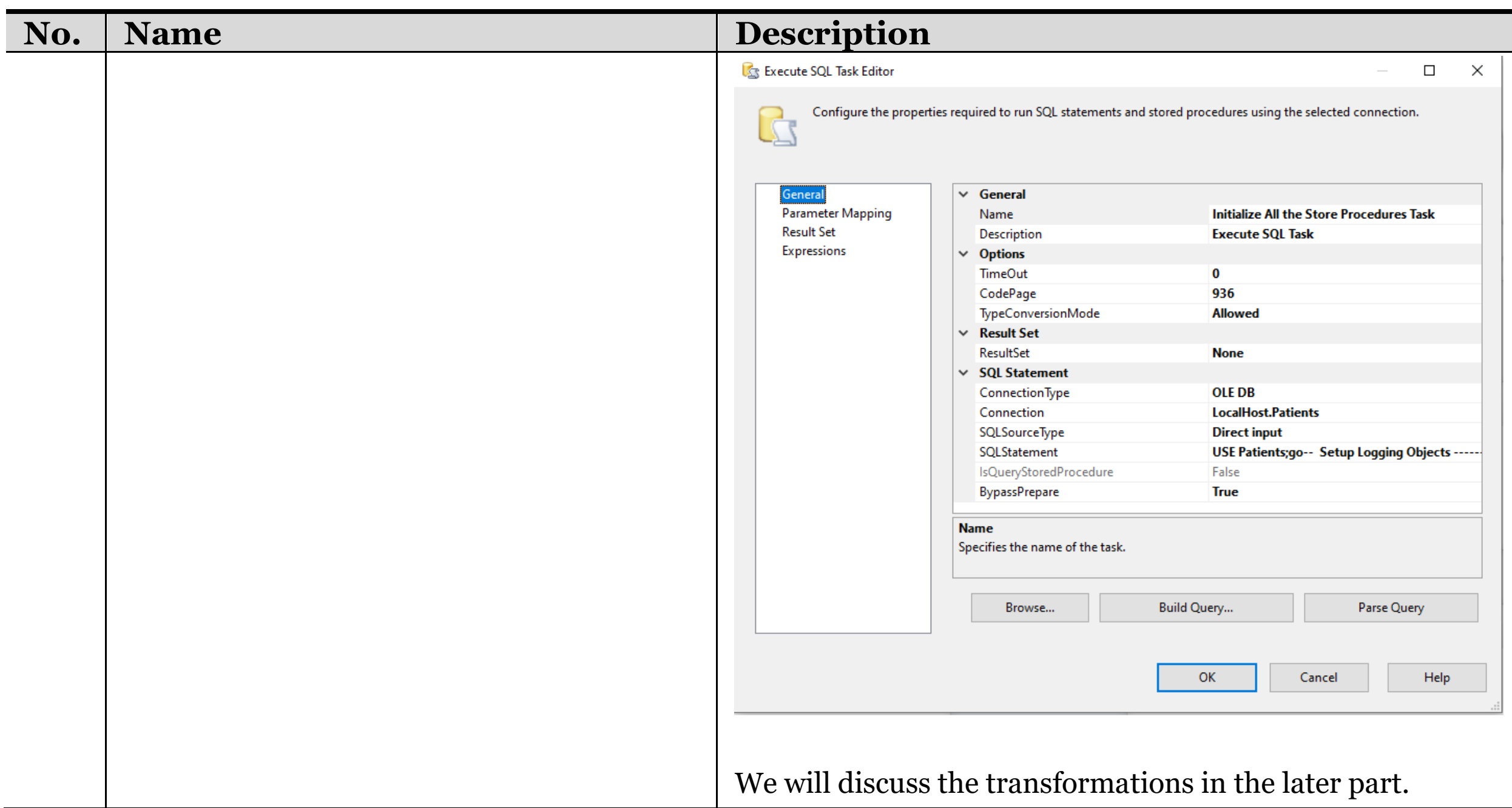

We will discuss the transformations in the later part.</td></tr>
</table>

### *4.2.3. Pre-Processing Python Sequence Container*

There is one component in the container.

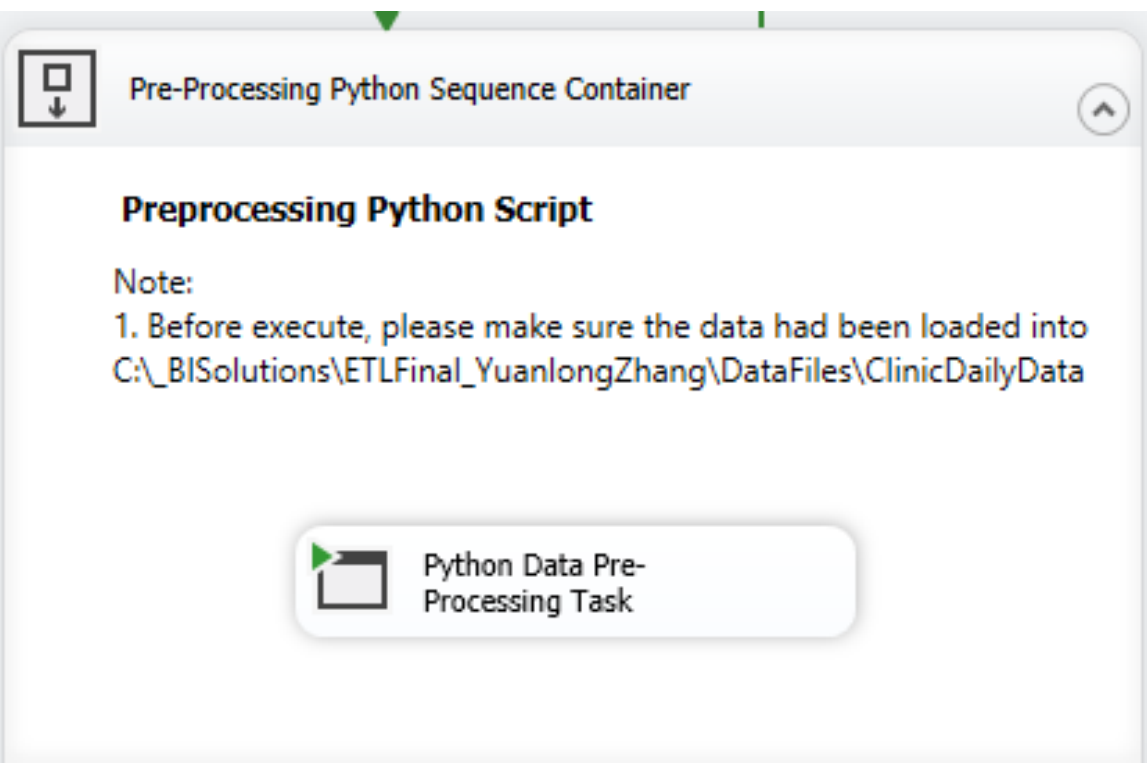


*Figure 26 Preprocessing Python Script in Milestone 1*

**Note:**
- The data file should be loaded into c:_BISolutions\ETLFinal_YuanlongZhang\DataFiles\ClinicDailyData
- We will analyze the python script in a later part.

### *4.2.4. Load Staging Tables Sequence Container*

There are two components in the container.

*Figure 27 Load Staging Tables Sequence Container in Milestone 1*

| No | Name | Description |
|----|------|-------------|
| 1 | Fill Patient Table Task  | Used to fill data into the patient table. It executes pETLSyncPatients store-procedure  We will provide more details about the script in the next section. |
| 2 | Fill Visits Table Task  | Used to fill data into visit table. It executes pETLSyncVisits store-procedure |

| No | Name | Description |
|----|------|-------------|
|    |      |  We will provide more details about the script in the next section. |

## 4.3. SSIS Package for Milestone 2

Here's the list of files for our ETL processes for Milestone 2:

| | File Name | Description |
|---|-----------|-------------|
| 1 | DWClinicReportDataETL.dtsx | Set of tasks that move data from database to data warehouse. <br> *Note: _BISolutions\ETLFinal_YuanlongZhang\ETLPackages* |
| 2 | Create DWClinicReportData.sql | SQL Script that creates the data warehouse. <br> *Note: _BISolutions\ETLFinal_YuanlongZhang\SQLScripts* |
| 3 | DWETLObjects.sql | SQL Script for views/store procedures definition and transformation to ETL the data from database to data warehouse. <br> *Note: _BISolutions\ETLFinal_YuanlongZhang\SQLScripts* |

## 4.4. *Data Warehouse Design for Milestone 2*



*Figure 28 Diagram for Data Warehouse in Milestone 2*

✦ **Note:**
- All the tables (except DimPatients) are a Type-1 SCD design.
- DimPatients is a Type-2 SCD design.

## 4.5. *ETL Visualization for Milestone 2*

We use Microsoft's SQL Server Integration Services (SSIS) to visualize our ETL process. SSIS includes a set of commands represented by visual objects, called tasks, constraints, containers, or transformations. The most common objects used in our ETL process are Sequence Containers, Execute SQL tasks, Precedent Constraints, and Connections.
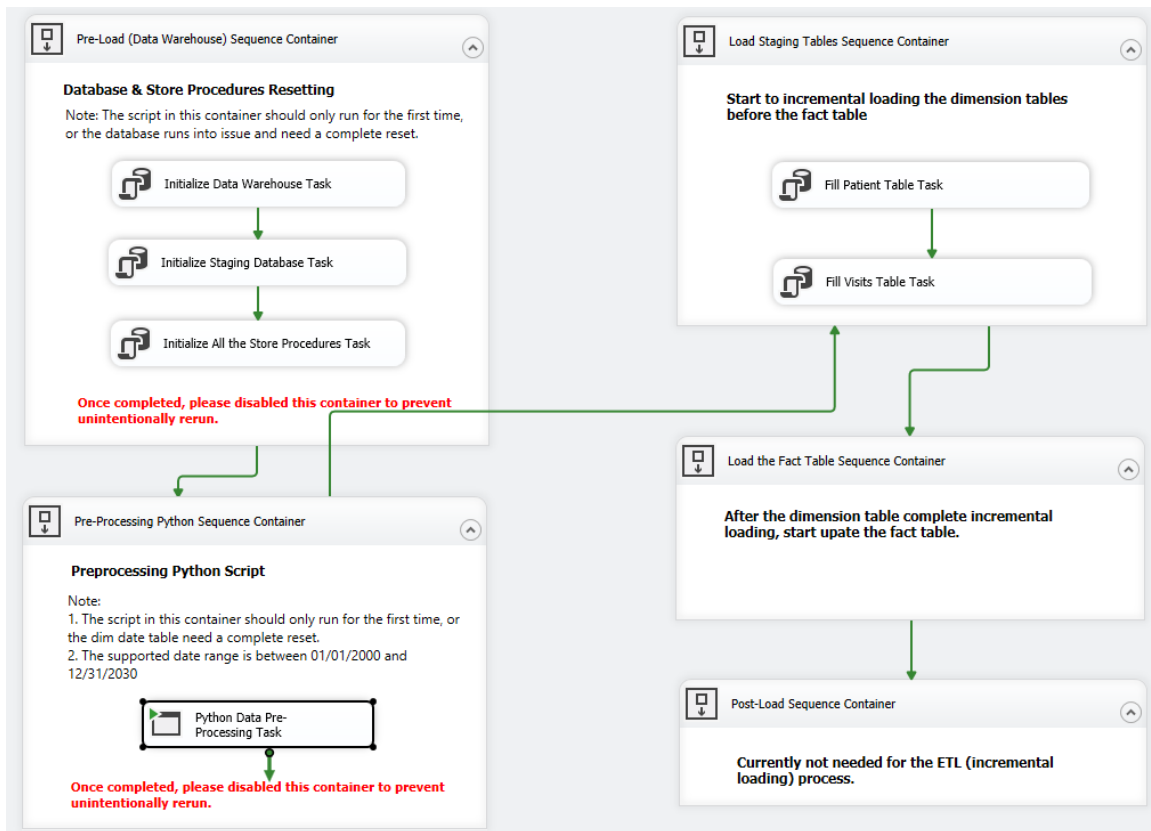
*Figure 29 Overview of SSIS package in Milestone 2*

✦ **Note:**

- These visualization objects are stored and configured in an SSIS Package file.

### 4.5.1. SSIS Connections

We use a Connection object that connects to the original data source file, invalid phone number file, and data warehouse that holds our SQL ETL objects.



*Figure 30 Connections in Milestone 2*

We used OLE DB connections; here's the detailed information:

| No. | Name | Description |
|-----|------|-------------|
| 1 | localhost.master | This is an OLE DB connection. Will connect to localhost\master |
| 2 | localhost.DWClinicReportDataYuanlongZhang | This is an OLE DB connection. |

| No. | Name | Description |
|-----|------|-------------|
|     |      | Will connect to localhost\DWClinicReportDataYuanlongZhang |

### 4.5.2. Database Connections

There are two external connections in milestone 2:



*Figure 31 High-level Diagram for Milestone 2*

- *The **DoctorsSchedule** database is on our class's **Azure cloud Server***
  - o Server: **continuumsql.westus2.cloudapp.azure.com**
  - o Login: BICert
  - o Password: BICert
- *The **Patients** database is on **our class's** **UW Server***
  - o Server: **is-root01.ischool.uw.edu\BI**
  - o Login: BICert
  - o Password: BICert

➕ **Note:**
- To access my UW SQL Server, **you must connect to the UW network** using "OnNet" Virtual Private **Network software:** Husky-OnNet

The database diagrams are as follows:



*Figure 32 Database Diagram for DoctorsSchedule*

*Figure 33 Database Diagram for Patients*

### 4.5.3. Pre-Load (Data Warehouse) Sequence Container
There are three components inside the sequence container.



*Figure 34 Pre-Load Sequence Container for Milestone 2*

✦ **Note:**
• This container should only run ONCE.

| No. | Name | Description |
|---|---|---|
| 1 | Initialize Data Warehouse Task  | This uses the script to initialize all the data warehouse tables and constraints: It creates the local copy for the data warehouse. The script is available in Create DWClinicReportData.sql |
| 2 | Initialize All the Store Procedures Task | Used to create all the views and store procedures definition and transformation to ETL the data from database to data warehouse. |

| No. | Name | Description |
|-----|------|-------------|
| |  Initialize All the Store Procedures Task | The script is available in DWETLObjects.sql |

### 4.5.4.  Pre-Load DimDate Table Sequence Container
There is one component inside the sequence container.



*Figure 35 Preprocessing DimDate Table Sequence Container in Milestone 2*

**Note:**
- This container should only run ONCE.
- The date range is from 01/01/2000 ~ 12/31/2030. Please adjust the time range if needed.

### 4.5.5.  Loading Staging Table  Sequence Container
There are five components inside the sequence container.



*Figure 36 Load Staging Tables Sequence Container in Milestone 2*

✚ **Note:**
- These items are executed in the store-procedures created in the Pre-Load stage. It had a similar logic as we explained in milestone 1.
- For any transformation we performed, please refer to the later section.

| No. | Name | Executed Store Procedure | SCD Design |
|---|---|---|---|
| 1 | Fill DimClinic Table Task | pETLSyncDimClinics | Type 1 |
| 2 | Fill DimDoctors Table Task | pETLSyncDimDoctors | Type 1 |
| 3 | Fill DimPatients Table Task | pETLSyncDimPatients | Type 2 |
| 4 | Fill DimProcedures Table Task | pETLSyncDimProcedures | Type 1 |
| 5 | Fill DimShifts Table Task | pETLSyncDimShifts | Type 1 |

### 4.5.6. *Loading Fact Table Sequence Container*
There are two components inside the sequence container.



*Figure 37 Load the Fact Table Sequence Container in Milestone 2*

✚ **Note:**
- These items are executed in the store-procedures created in the Pre-Load stage. It had a similar logic as we explained in milestone 1.
- We will explain all the complex transformations in the later section.

| No. | Name | Executed Store Procedure | SCD Design |
|---|---|---|---|
| 1 | Fill FactDoctorShifts Table Task | pETLSyncFactDoctorShifts | Type 1 |
| 2 | Fill FactVisits Table Task | pETLSyncFactVisits | Type 1 |

### 4.6. *SSIS Package for Milestone 3*
Here's the list of files for our ETL processes for Milestone 1:

| | File Name | Description |
|---|---|---|
| 1 | ETLClinicReportsDocumentData.dtsx | Set of tasks that define views and move data from the data warehouse to MongoDB. <br> *Note: _BISolutions\ETLFinal_YuanlongZhang\ETLPackages* |
| 2 | MongoDBETL.py | Script that defines views and moves data from the data warehouse to MongoDB. <br> *Note: _BISolutions\ETLFinal_YuanlongZhang\Scripts* |
| 3 | DWMongoDBView.sql | SQL Script for creating views used for MongoDB upload. <br> *Note: _BISolutions\ETLFinal_YuanlongZhang\SQLScripts* |

| File Name | Description |
|-----------|-------------|
|           | *Note2: The script had already been included in Python Script. There's no need to run this SQL Script separately.* |

## 4.7. ETL Visualization for Milestone 2

We use Microsoft's SQL Server Integration Services (SSIS) to visualize our ETL process. SSIS includes a set of commands represented by visual objects, called tasks, constraints, containers, or transformations. The most common objects used in our ETL process are Sequence Containers, Execute SQL tasks, Precedent Constraints, and Connections.
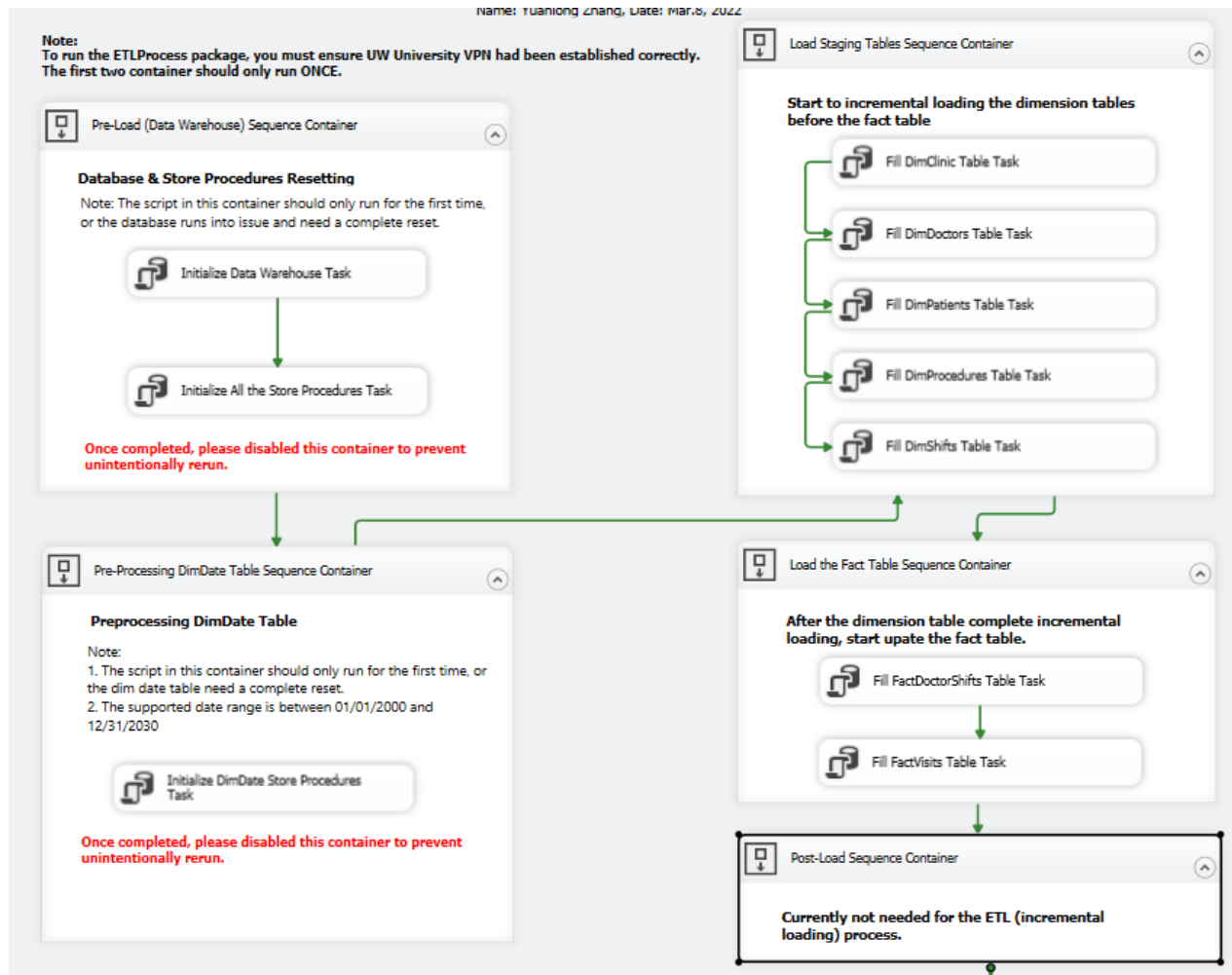


*Figure 38 Overview of SSIS package: Control Flow*

**⊥ Note:**
- These visualization objects are stored and configured in an SSIS Package file.

### 4.7.1. External Connections

We use external connections to MongoDB.



*Figure 39 High-level Diagram for Milestone 3*

The connection string is as follows:
`'mongodb+srv://BICert:BICert@clinicreportsdata.ts9ek.mongodb.net/test?retryWrites=true&w=majority'`

To ensure access, we configure the 0.0.0.0/0 in the IP Address list:

*Figure 40 Screenshot for MongoDB IP Address Configurations*

If you configure everything correctly, after executing SSIS for milestone 3, you can access the database by **connection string** in MongoDB Compass:



*Figure 41 Screenshot for MongoDB Compass*

### 4.7.2. *MongoDB Python Sequence Container*



*Figure 42 MongoDB Python Sequence Container in Milestone 3*

✦ **Note:**
- The SSIS package only has one python script task. We will explain the python script in a later section.

## 4.8. SSIS Package for Milestone 4

Here's the list of files for our ETL processes for Milestone 4:

| | File Name | Description |
|---|---|---|
| 1 | ETLJob.dtsx | Set of tasks triggered by SQL Agent Job. |
| | | *Note: _BISolutions\ETLFinal_YuanlongZhang\ETLPackages* |
| 2 | ETLJob.sql | SQL Script, which defined the SQL Agent Jobs and SSIS Proxy. |
| | | *Note: _BISolutions\ETLFinal_YuanlongZhang\SQLScripts* |
| 3 | ETLViews.sql | SQL Script, which defined the views for SSRS paginated report. |
| | | *Note: _BISolutions\ETLFinal_YuanlongZhang\SQLScripts* |
| 4 | MainDashboard.rdl | Paginated report for the landing page (high-level summary) |
| | | *Note: _BISolutions\ETLFinal_YuanlongZhang\ETLDashboardReports* |
| 5 | ETLDetails.rdl | ReportPaginated report for ETL details information |
| | | *Note: _BISolutions\ETLFinal_YuanlongZhang\ETLDashboardReports* |
| 6 | SSISJobDetails.rdl | Paginated report for SSIS Job information |
| | | *Note: _BISolutions\ETLFinal_YuanlongZhang\ETLDashboardReports* |

## 4.9. ETL Visualization for Milestone 4

We use Microsoft's SQL Server Integration Services (SSIS) to visualize our ETL process. SSIS includes a set of commands represented by visual objects, called tasks, constraints, containers, or transformations. The most common objects used in our ETL process are Sequence Containers, Execute SQL tasks, Precedent Constraints, and Connections.



*Figure 43 Overview of SSIS package: Control Flow*

✦ **Note:**
- These visualization objects are stored and configured in an SSIS Package file.
- This package contains 3 SSIS packages we created for Milestone 1 ~ Milestone 3.
- There are technical issues when using SQL Agent Jobs to execute Python Scripts. Thus we disabled the SSIS package for milestone 1 and milestone 3.

### 4.9.1.  Connections

We use a Connection object that connects to the original data source file, invalid phone number file, and data warehouse that holds our SQL ETL objects.

Connection Managers

DWClinicReportDataETL.dtsx    ETLClincReportsDocumentData.dtsx    ETLFilesToDatabases.dtsx    LocalHost.master

*Figure 44 Connections in Milestone 4*

We used 1 OLE DB connection and 3 file connections; here's the detailed information:

| No. | Name | Description |
| --- | --- | --- |
| 1 | localhost.master | This is an OLE DB connection. Will connect to localhost\master |
| 2 | DWClinicReportDataETL.dtsx | This is a file connection. Will connect to DWClinicReportDataETL.dtsx (SSIS packge for Milestone 2) |
| 3 | ETLClinicReportsDocumentData.dtsx | This is a file connection. Will connect to ETLClinicReportsDocumentData.dtsx (SSIS package for Milestone 3) *Note: This component had been disabled in the SSIS package because it included the Python Script task.* |
| 4 | ETLFilesToDatabases.dtsx | This is a file connection. Will connect to ETLFilesToDatabases.dtsx (SSIS package for Milestone 1) *Note: This component had been disabled in the SSIS package because it included the Python Script task.* |

### 4.9.2.  ETL Job Sequence Container

There are four components inside the sequence container.



ETL Job Sequence Container

**ETL Job (Package Container)**

Execute ETLFilesToDatabases Package Task

Execute DWClinicReportDataETL Package Task

Execute ETLClincReportsDocumentData Package Task

Execute Create ETL View Task

*Figure 45 Screenshot for SSIS Package for Milestone 4*

➕ **Note:**
- This container should only run ONCE.

| No. | Name | Description |
|---|---|---|
| 1 | The first three package tasks:<br><br>Execute ETLFilesToDatabases Package Task<br><br>Execute DWClinicReportDataETL Package Task<br><br>Execute ETLClincReportsDocumentData Package Task | Which refer to the SSIS packages for Milestone 1~ 3:<br>Please refer to the previous section if you need to know more details.<br><br>We configured the reference type as an External reference, ensuring SQL Agent Job can find the file and execute it correctly. |
| 2 | Execute Create ETL View Task | Used to execute the SQL Script, which defined the views for SSRS paginated report.<br>The script is available in ETLViews.sql |

### 4.9.3.   SQL Agent Job for Milestone 4

To make sure the SQL Job works, we create an SSIS proxy as follows:

*Figure 46 Screenshot for SSIS Proxy configurations*

Then we created the SQL Agent Job by ETLJob.sql script as follows:

Set the Job Owner as SA:



*Figure 47 Screenshot for ETLDWClinicReportData job configurations (General)*



*Figure 48 Screenshot for ETLDWClinicReportData job configurations (Steps)*

Configured to run each night at 1 AM:

*Figure 49 Screenshot for ETLDWClinicReportData job configurations (Schedules)*

## 4.10. SSRS ETL Dashboard for Milestone 4
The package has three dashboards: MainDashboard, SSISJobDetails, and ETL Details.

You should always start from MainDashboard.

### 4.10.1.  MainDashboard



*Figure 50 Screenshot for MainDashboard.rdl (design mode)*

The dashboard is used as a landing page for all high-level summary information; it also allows you to navigate to the other two dashboards **by clicking on "(details)":**

*Figure 51 Screenshot for MainDashboard.rdl (preview mode)*

### 4.10.2. SSISJobDetails Dashboard



*Figure 52 Screenshot for SSISJobDetails.rdl (design mode)*

This dashboard provides a detailed view of the SSIS Job Log (which is the view from MSDB):

✦ **Note:**
• You need to select the status (success/failure) and click on "view report."

To view the result "success," please use the filter at the left top corner and click "view report" at the right top corner.



*Figure 53 Screenshot for SSISJobDetails.rdl (preview mode: Failure)*



*Figure 54 Screenshot for SSISJobDetails.rdl (preview mode: Success)*

### 4.10.3. *ETLDetails Dashboard*



*Figure 55 Screenshot for ETLDetails.rdl (design mode)*

This dashboard provides a detailed view of the ETL Log table and the row count of the table in the data warehouse:



*Figure 56 Screenshot for ETLDetails.rdl (preview mode)*

# 5. Data Transformation and Considerations

We performed transformation at milestone 1 (CSV file to the database) and milestone 2 (database to the data warehouse):

## 5.1. Files to Source Databases

| Source Data | Source Type | Destination | Destination Type | Transformations |
|---|---|---|---|---|
| Folder Name | Directory Name | Staging.dbo.StagingVisits.Clinic | New Column | Extract folder name |
| Filename | Filename | Staging.dbo.StagingVisits.Date | New Column | Extract the first 8 digits of the filename (date part) |
| Email | csv | Staging.dbo.StagingNewPatient.Email | Column | Python Script to validate the valid email address and remove the invalid one |
| Staging.dbo.StagingVisits.Clinic | New Column | Patients.dbo.Patients.Date | Datetime | Combine the Date and Time column and convert it to datetime datatype |
| Staging.dbo.StagingVisits.Date | New Column | Patients.dbo.Patients.Charge | Decimal | Transform the int to decimal (18,2) |
| Staging.dbo.StagingNewPatient.Email | Column | Patients.dbo.Patients.Clinic | Int | Select the Clinic ID by joining the Patient.dbo.Clinics table |

## 5.2. Source DB to Data Warehouse

| Source Data | Source Type | Destination | Destination Type | Transformations |
|---|---|---|---|---|
| Patients.dbo.Patients.Date | Datetime | | | |
| Patients.dbo.Patients.Charge | Decimal | | | |
| Patients.dbo.Patients.Clinic | Int | | | |
| Patients.dbo.Patients.ID | int | DWClinicReportDataYuanlongZhang.dbo.DimPatients.PatientID | int | |
| Patients.dbo.Patients.FName | varchar(28) | DWClinicReportDataYuanlongZhang.dbo.DimPatients.PatientFullName | varchar(100) | Concate First Name and Last Name |
| Patients.dbo.Patients.LName | varchar(29) | DWClinicReportDataYuanlongZhang.dbo.DimPatients.PatientFullName | varchar(100) | Concate First Name and Last Name |
| Patients.dbo.Patients.Email | varchar(100) | | | |
| Patients.dbo.Patients.Address | varchar(97) | | | |
| Patients.dbo.Patients.City | varchar(72) | DWClinicReportDataYuanlongZhang.dbo.DimPatients.PatientCity | varchar(100) | |
| Patients.dbo.Patients.State | varchar(50) | DWClinicReportDataYuanlongZhang.dbo.DimPatients.PatientState | varchar(100) | |
| Patients.dbo.Patients.ZipCode | int | DWClinicReportDataYuanlongZhang.dbo.DimPatients.PatientZipCode | int | |
| Patients.dbo.Clinics.ID | int | | | |
| Patients.dbo.Clinics.Name | varchar(50) | | | |
| Patients.dbo.Clinics.Address | varchar(100) | | | |
| Patients.dbo.Clinics.City | varchar(50) | | | |
| Patients.dbo.Doctors.ID | int | | | |

| Source Data | Source Type | Destination | Destination Type | Transformations |
|---|---|---|---|---|
| Patients.dbo.Doctors.FName | varchar(28) | | | |
| Patients.dbo.Doctors.LName | varchar(19) | | | |
| Patients.dbo.Procedures.ID | int | DWClinicReportDataYuanlongZhang.dbo.DimProcedures.ProcedureID | int | |
| Patients.dbo.Procedures.Name | varchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimProcedures.ProcedureName | varchar(100) | |
| Patients.dbo.Procedures.Desc | varchar(1000) | DWClinicReportDataYuanlongZhang.dbo.DimProcedures.ProcedureDesc | varchar(1000) | |
| Patients.dbo.Procedures.Charge | money | DWClinicReportDataYuanlongZhang.dbo.DimProcedures.ProcedureCharge | money | |
| Patients.dbo.Visits.ID | int | DWClinicReportDataYuanlongZhang.dbo.FactVisits.VisitKey | int | |
| Patients.dbo.Visits.Date | datetime | DWClinicReportDataYuanlongZhang.dbo.FactVisits.DateKey | int | Extact Date from datetime |
| Patients.dbo.Visits.Clinic | int | DWClinicReportDataYuanlongZhang.dbo.FactVisits.ClinicKey | int | |
| Patients.dbo.Visits.Patient | int | DWClinicReportDataYuanlongZhang.dbo.FactVisits.PatientKey | int | |
| Patients.dbo.Visits.Doctor | int | DWClinicReportDataYuanlongZhang.dbo.FactVisits.DoctorKey | int | Replace null value as 'No Doctor' |
| Patients.dbo.Visits.Procedure | int | DWClinicReportDataYuanlongZhang.dbo.FactVisits.ProcedureKey | int | |
| Patients.dbo.Visits.Charge | money | DWClinicReportDataYuanlongZhang.dbo.FactVisits.ProcedureVistCharge | money | |
| DoctorsSchedules.dbo.Clinics.ClinicID | int | DWClinicReportDataYuanlongZhang.dbo.DimClinics.ClinicID | int | Convert 1,2,3 into 100, 200, 300 |
| DoctorsSchedules.dbo.Clinics.ClinicName | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimClinics.ClinicName | nvarchar(100) | |
| DoctorsSchedules.dbo.Clinics.Address | nvarchar(100) | | | |
| DoctorsSchedules.dbo.Clinics.City | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimClinics.ClinicCity | nvarchar(100) | |
| DoctorsSchedules.dbo.Clinics.State | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimClinics.ClinicState | nvarchar(100) | |
| DoctorsSchedules.dbo.Clinics.Zip | nvarchar(5) | DWClinicReportDataYuanlongZhang.dbo.DimClinics.ClinicZip | nvarchar(5) | |
| DoctorsSchedules.dbo.Doctors.DoctorID | int | DWClinicReportDataYuanlongZhang.dbo.DimDoctors.DoctorID | int | |
| DoctorsSchedules.dbo.Doctors.FirstName | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimDoctors.DoctorFullName | int | Concate First Name and Last Name |
| DoctorsSchedules.dbo.Doctors.LastName | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimDoctors.DoctorFullName | nvarchar(200) | Concate First Name and Last Name |
| DoctorsSchedules.dbo.Doctors.EmailAddress | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimDoctors.DoctorEmailAddress | nvarchar(200) | |
| DoctorsSchedules.dbo.Doctors.Address | nvarchar(100) | | | |
| DoctorsSchedules.dbo.Doctors.City | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimDoctors.DoctorCity | nvarchar(100) | Trim the space |
| DoctorsSchedules.dbo.Doctors.State | nvarchar(100) | DWClinicReportDataYuanlongZhang.dbo.DimDoctors.DoctorState | nvarchar(100) | Marked bad data as "(Error)" |
| DoctorsSchedules.dbo.Doctors.Zip | nvarchar(5) | DWClinicReportDataYuanlongZhang.dbo.DimDoctors.DoctorZip | nvarchar(5) | |
| DoctorsSchedules.dbo.Shifts.ShiftID | int | DWClinicReportDataYuanlongZhang.dbo.DimShifts.ShiftID | int | |
| DoctorsSchedules.dbo.Shifts.ShiftStart | time | DWClinicReportDataYuanlongZhang.dbo.DimShifts.ShiftStart | time | Transformations of time format: 12:00 to 24:00 |
| DoctorsSchedules.dbo.Shifts.ShiftEnd | time | DWClinicReportDataYuanlongZhang.dbo.DimShifts.ShiftEnd | time | Transformations of time format: 12:00 to 24:00 |

| Source Data | Source Type | Destination | Destination Type | Transformations |
|---|---|---|---|---|
| DoctorsSchedules.dbo.DoctorShifts.DoctorsShiftID | int | DWClinicReportDataYuanlongZhang.dbo.FactDoctorShifts.DoctorsShiftID | int | |
| DoctorsSchedules.dbo.DoctorShifts.ShiftDate | datetime | DWClinicReportDataYuanlongZhang.dbo.FactDoctorShifts.ShiftDateKey | int | Extract Shiftkey from DimShift |
| DoctorsSchedules.dbo.DoctorShifts.ClinicID | int | DWClinicReportDataYuanlongZhang.dbo.FactDoctorShifts.ClinicKey | int | |
| DoctorsSchedules.dbo.DoctorShifts.ShiftID | int | DWClinicReportDataYuanlongZhang.dbo.FactDoctorShifts.ShiftKey | int | Convert 1,2,3 into 100, 200, 300 |
| DoctorsSchedules.dbo.DoctorShifts.DoctorID | int | DWClinicReportDataYuanlongZhang.dbo.FactDoctorShifts.DoctorKey | int | |
| | | DWClinicReportDataYuanlongZhang.dbo.FactDoctorShifts.HoursWorked | int | Calculated based on Shiftstart and ShiftEnd |
| | | DWClinicReportDataYuanlongZhang.dbo.DimPatients.StartDate | date | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimPatients.EndDate | date | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimPatients.IsCurrent | int | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimDates.DateKey | int | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimDates.FullDate | datetime | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimDates.FullDateName | nvarchar(50) | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimDates.MonthID | int | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimDates.MonthName | nvarchar(50) | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimDates.YearID | int | |
| | | DWClinicReportDataYuanlongZhang.dbo.DimDates.YearName | nvarchar(50) | |
| | | DWClinicReportDataYuanlongZhang.dbo.ETLLog.ETLLogID | int | |
| | | DWClinicReportDataYuanlongZhang.dbo.ETLLog.ETLDateAndTime | datetime | |
| | | DWClinicReportDataYuanlongZhang.dbo.ETLLog.ETLAction | varchar(100) | |
| | | DWClinicReportDataYuanlongZhang.dbo.ETLLog.ETLLogMessage | varchar(2000) | |
| | | DWClinicReportDataYuanlongZhang.dbo.vETLLog.ETLLogID | int | |
| | | DWClinicReportDataYuanlongZhang.dbo.vETLLog.ETLDate | nvarchar(4000) | |
| | | DWClinicReportDataYuanlongZhang.dbo.vETLLog.ETLTime | nvarchar(4000) | |
| | | DWClinicReportDataYuanlongZhang.dbo.vETLLog.ETLAction | varchar(100) | |
| | | DWClinicReportDataYuanlongZhang.dbo.vETLLog.ETLLogMessage | varchar(2000) | |

# 6. ETL Scripts Analysis & Explanation

In this section, we will demonstrate and provide more details about the script file:

## 6.1. Python Script

### 6.1.1. PythonETL.py (Milestone 1)

```python
1
2   # import required module
3   import os
4   import pandas as pd
5   import re
6   import pyodbc
7
8   # assign directory
9   directory = 'C:\_BISolutions\ETLFinal_YuanlongZhang\DataFiles\ClinicDailyData'
10
11  StagingData = pd.DataFrame()
12  StagingNewPatient = pd.DataFrame()
13  StagingVisits = pd.DataFrame()
14
15
16  # iterate over files in
17  # that directory
18  for root, dirs, files in os.walk(directory):
19      for filename in files:
20          StagingData = pd.read_csv(root + "\\" + filename)
21          StagingData["Date"] = filename[:8]
22          StagingData["Clinic"] = root.split("\\")[-1]
23          if filename[8:11] == "New":
24              frames = [StagingNewPatient, StagingData]
25              StagingNewPatient = pd.concat(frames)
26
27          if filename[8:11] == "Vis":
28              frames = [StagingVisits, StagingData]
29              StagingVisits = pd.concat(frames)
30
31  dfnew = pd.DataFrame()
32  dferror = pd.DataFrame()
33  i, j = StagingNewPatient.shape
34  for row in range(0,i):
35      temp = StagingNewPatient.iloc[row,:]
36      if(re.search('^\w+([\.-]?\w+)*@\w+([\.-]?\w+)*(\.\w{2,3})+$',temp["Email"])):
37          frames = [dfnew, temp]
38          dfnew = pd.concat(frames,axis=1)
39      else:
40          def calculate_area(row):
41              return "(error)" + row['Email']
42          temp["Email"] = temp["Email"] + " (error)"
43          frames = [dfnew, temp]
44          dfnew = pd.concat(frames, axis=1)
45
46  # print(StagingNewPatient)
47  df_corrected = dfnew.transpose()
48  df_error = dferror.transpose()
49
50  # Load the data to SQL Server
51  # Connect to SQL
52  con_str = ("Driver={SQL Server Native Client 11.0};"
53             "Server=localhost;"
54             "Database=master;"
55             "Trusted_Connection=yes;")
56  con_obj = pyodbc.connect(con_str)
57
58  # Create (update or delete) uses a cursor
59  cusor_obj = con_obj.cursor()
60  # cusor_obj.execute('''
61  #         USE [master];
62  #         If Exists (Select * from Sysdatabases Where Name = 'Staging')
```

*Read the CSV file from ClinicDailyData folder*

*Combine all the data from CSV file into the dataframes*

*Split data based on new patients/visits data; Add additional columns based on file name and folder name.*

*Use regular expression to determine the valid email-address.*
*Marked invalid email address with suffix (error).*

*Established SQL Server connections.*

```
63   #              Begin
64   #                  ALTER DATABASE [Staging] SET SINGLE_USER WITH ROLLBACK IMMEDIATE
65   #                  DROP DATABASE [Staging]
66   #              End
67   #          Create Database [Staging];
68   #          ALTER DATABASE Staging SET SINGLE_USER WITH ROLLBACK IMMEDIATE;
69   #          ALTER DATABASE Staging collate SQL_Latin1_General_CP1_CI_AS;
70   #          ALTER DATABASE Staging SET MULTI_USER;
71   #          ''')
72   cusor_obj.execute("Use Staging;If (object_id('StagingNewPatient') is not null) Drop Table
73   StagingNewPatient;")
74   SQL_str = '''Create Table StagingNewPatient
75              (FName nvarchar(50)
76              ,LName nvarchar(50)
77              ,Email nvarchar(100)
78              ,Address nvarchar(100)
79              ,City nvarchar(50)
80              ,State nvarchar(50)
81              ,ZipCode nvarchar(20)
82              ,Date nvarchar(20)
83              ,Clinic nvarchar(20) )
84              '''
85   cusor_obj.execute(SQL_str)
86
87   cusor_obj.execute("Use Staging;If (object_id('StagingVisits') is not null) Drop Table
88   StagingVisits;")
89   SQL_str = '''Create Table StagingVisits
90              (Time time
91              ,Patient int
92              ,Doctor int
93              ,[Procedure] int
94              ,Charge int
95              ,Date int
96              ,Clinic nvarchar(20) )
97              '''
98   cusor_obj.execute(SQL_str)
99
100  # Insert Dataframe into SQL Server:
101  SQL_str = ''' INSERT INTO Staging.dbo.StagingNewPatient
102              (FName, LName, Email, Address, City, State, ZipCode, Date, Clinic)
103              values(?,?,?,?,?,?,?,?,?) '''
104  for index, row in df_corrected.iterrows():
105      cusor_obj.execute(SQL_str, row.FName, row.LName, row.Email, row.Address, row.City,
106  row.State,
107      str(row.ZipCode), str(row.Date), row.Clinic)
108
109  con_obj.commit()
110
111  # Insert Dataframe into SQL Server:
112  SQL_str = ''' INSERT INTO Staging.dbo.StagingVisits
113              (Time, Patient, Doctor, [Procedure], Charge, Date, Clinic)
114              values(?,?,?,?,?,?,?) '''
115  for index, row in StagingVisits.iterrows():
116      cusor_obj.execute(SQL_str, row.Time, row.Patient, row.Doctor, row.Procedure,
117  row.Charge, row.Date, row.Clinic)
118
119  con_obj.commit()
120  cusor_obj.close()
121
122  # Save ErrorReport data to a new file.
123  # df_error.to_csv('C:\_BISolutions\ETLFinal_YuanlongZhang\ErrorReport\ErrorReport.csv')
124
125  # print(df_corrected)
126  # print(df_error)
127  # print(StagingVisits)
```

*To avoid dirty read, we reset the table every time when ETL process execute.*

*Insert the records from dataframe into staging database*

### 6.1.2. MongoDBETL.py (Milestone 3)

```
1   '''*************************************************************************--
2   -- Desc: This script connects to MongoDB Atlas cloud and Perform ETL procedures
3   -- Change Log: When,Who,What
4   -- 2020-05-15,RRoot,Created File
5   -- 2022-02-21,Yuanlong Zhang, Updated the File to add new functions for ETL.
6   -- 2022-03-09,Yuanlong Zhang, Updated the file for the final project.
7   --*************************************************************************'''
8
9   # Please install the package:
10  # Python -m pip install pyodbc
11
12  import pymongo
13  import pandas as pd
14  import pyodbc
15
16  try:
17          # Connect to SQL
18          con_str = ("Driver={SQL Server Native Client 11.0};"
19                     "Server=localhost;"
20                     "Database=DWClinicReportDataYuanlongZhang;"
21                     "Trusted_Connection=yes;")
22          con_obj = pyodbc.connect(con_str)
23
24          # Create (update or delete) uses a cursor
25          cusor_obj = con_obj.cursor()
26          cusor_obj.execute("IF (Object_ID('vRptDoctorShifts') is not null) Drop View
27  vRptDoctorShifts;")
28          SQL_str = '''Create View vRptDoctorShifts
29                              AS
30                              Select
31                              [ShiftDate] = Cast(Cast([FullDate] as date) as varchar(100))
32                              ,[ClinicName] = dc.ClinicName
33                              ,[ClinicCity] = dc.ClinicCity
34                              ,[ClinicState] = dc.ClinicState
35                              ,[ShiftID] = ds.ShiftID
36                              ,[ShiftStart] = ds.ShiftStart
37                              ,[ShiftEnd] = ds.ShiftEnd
38                              ,[DoctorFullName] = ddo.DoctorFullName
39                              ,[HoursWorked] = fds.HoursWorked
40                              FROM
41                              dbo.FactDoctorShifts fds
42                              JOIN dbo.DimDates dd
43                              ON fds.ShiftDateKey = dd.DateKey
44                              JOIN dbo.DimClinics dc
45                              ON fds.ClinicKey = dc.ClinicID
46                              JOIN dbo.DimShifts ds
47                              ON fds.ShiftKey = ds.ShiftID
48                              JOIN dbo.DimDoctors ddo
49                              ON fds.DoctorKey = ddo.DoctorKey;'''
50          cusor_obj.execute(SQL_str)
51
52          # Select from SQL Server using the Pandas module!
53          vRptDoctorShifts = pd.read_sql("select * from vRptDoctorShifts;", con_obj)
54          #print(vRptDoctorShifts)
55
56          cusor_obj.execute("IF (Object_ID('vRptPatientVisits') is not null) Drop View
57  vRptPatientVisits;")
58          SQL_str = '''Create View vRptPatientVisits
59                              AS
60                              Select
61                              [VisitDate] = Cast(Cast([FullDate] as date) as varchar(100))
62                              ,[ClinicName] = dc.ClinicName
63                              ,[ClinicCity] = dc.ClinicCity
64                              ,[ClinicState] = dc.ClinicState
65                              ,[ProcedureName] = dp.ProcedureName
66                              ,[PatientFullName] = dpt.PatientFullName
67                              ,[PatientCity] = dpt.PatientCity
68                              ,[PatientState] =dpt.PatientState
69                              ,[DoctorFullName] = ddo.DoctorFullName
```

*Establish connections to Data Warehouse*

*To avoid dirty read, reset and recreate the view every time*

*Load the data from SQL view into the dataframe.*

```
70                            ,[ProcedureVisitCharge] = fv.ProcedureVistCharge
71                            FROM
72                            dbo.FactVisits fv
73                            JOIN dbo.DimDates dd
74                            ON fv.DateKey = dd.DateKey
75                            JOIN dbo.DimClinics dc
76                            ON fv.ClinicKey = dc.ClinicID
77                            JOIN dbo.DimProcedures dp
78                            ON fv.ProcedureKey = dp.ProcedureID
79                            JOIN dbo.DimDoctors ddo
80                            ON fv.DoctorKey = ddo.DoctorID
81                            JOIN dbo.DimPatients dpt
82                            ON fv.PatientKey = dpt.PatientID;'''
83          cusor_obj.execute(SQL_str)
84
85          # Select from SQL Server using the Pandas module!
86          vRptPatientVisits = pd.read_sql("select * from vRptPatientVisits;", con_obj)
87          #print(vRptPatientVisits)
88
89          # Always clean up your objects when done!
90          cusor_obj.close()
91          con_obj.close()
92
93
94   vRptPatientVisits.to_csv('C:\\_BISolutions\\ETLFinal_YuanlongZhang\\DataFiles\\Staging\\vRptPatie
95   ntVisitStaging.csv')
96
97   vRptDoctorShifts.to_csv('C:\\_BISolutions\\ETLFinal_YuanlongZhang\\DataFiles\\Staging\\vRptDoctor
98   ShiftsStaging.csv')
99
100          MongoRptPatient =
101  pd.read_csv('C:\\_BISolutions\\ETLFinal_YuanlongZhang\\Data
102  g.csv')
103          MongoRptDoctor =
104  pd.read_csv('C:\\_BISolutions\\ETLFinal_YuanlongZhang\\DataFiles\\Staging\\vRptDoctorShiftsStagin
105  g.csv')
106
107          # 1. Whitelist your IP address and create a user on
108          # ( https://account.mongodb.com/account/login )
109          # 2. Create a connection string
110          # Note: the connection string can change without notice!
111          strCon =
112  'mongodb+srv://BICert:BICert@clinicreportsdata.ts9ek.mongodb.net/test?retryWrites=true&w=majority
113  '
114
115          objCon = pymongo.MongoClient(strCon)
116          db = objCon["ClinicReportsData"]
117          db.drop_collection('DoctorsShifts')
118          objCol = db.create_collection('DoctorsShifts')
119          MongoRptDoctor['_id'] = MongoRptDoctor['Unnamed: 0']
120          MongoRptDoctor = MongoRptDoctor.drop(columns=['Unnamed: 0'])
121          MongoRptDoctor = MongoRptDoctor.to_dict('records')
122          objCol.insert_many(MongoRptDoctor)
123
124          db.drop_collection('PatientsVisits')
125          objCol = db.create_collection('PatientsVisits')
126          MongoRptPatient['_id'] = MongoRptPatient['Unnamed: 0']
127          MongoRptPatient = MongoRptPatient.drop(columns=['Unnamed: 0'])
128          MongoRptPatient = MongoRptPatient.to_dict('records')
129          objCol.insert_many(MongoRptPatient)
130
131  except Exception as e:
132          print(e)
```

*Export the data into the external CSV file and preparing for uploading.*

*Read the CSV file into dataframe.
(This step can solve the data type issue)*

*MongoDB connection string*

*Drop the collection every time before upload*

*Reset the id column and mapping to dataframe id*

*Transform the data type and insert into MongoDB.*

## 6.2. SQL Script

➕ **Note:**
* The SQL Script is very long. Thus, we hide the repeatable parts, help files, and the code provided by Randal.

### 6.2.1. ETLLoading.sql (Milestone 1)

We use vETLNewPatient as an example:

```
1   --****************************************************************--
2   -- Title: DWFinal-ETL Loading
3   -- Author: Yuanlong Zhang
4   -- Desc: This file used to load the new data into the database
5   -- Change Log: When,Who,What
6   -- 2021-01-17,RRoot,Created File
7   -- 2022-01-24,Yuanlong Zhang, Completed File
8   -- 2022-03-01,Yuanlong Zhang, Updated the file for Final Projects
9
10  --****************************************************************--
11
12
13  USE Patients;
14  go
15
16  --  Setup Logging Objects ----------------------------------------------
17
18  If NOT Exists(Select * From Sys.tables where Name = 'ETLLog')
19    Create -- Drop
20    Table ETLLog
21    (ETLLogID int identity Primary Key
22    ,ETLDateAndTime datetime Default GetDate()
23    ,ETLAction varchar(100)
24    ,ETLLogMessage varchar(2000)
25    );
26  go
27
28  Create or Alter View vETLLog
29  As
30    Select
31     ETLLogID
32    ,ETLDate = Format(ETLDateAndTime, 'D', 'en-us')
33    ,ETLTime = Format(Cast(ETLDateAndTime as datetime2), 'HH:mm', 'en-us')
34    ,ETLAction
35    ,ETLLogMessage
36    From ETLLog;
37  go
38
39
40  Create or Alter Proc pInsETLLog
41   (@ETLAction varchar(100), @ETLLogMessage varchar(2000))
42  --****************************************************************--
43  -- Desc: This Sproc creates an admin table for logging ETL metadata.
44  -- Change Log: When,Who,What
45  -- 2020-01-01,RRoot,Created Sproc
46  --****************************************************************--
47  As
48  Begin
49    Declare @RC int = 0;
50    Begin Try
51      Begin Tran;
52        Insert Into ETLLog
53          (ETLAction,ETLLogMessage)
54        Values
55          (@ETLAction,@ETLLogMessage)
56      Commit Tran;
57      Set @RC = 1;
58    End Try
59    Begin Catch
60      If @@TRANCOUNT > 0 Rollback Tran;
61      Set @RC = -1;
62    End Catch
```

*Define ETL Log Table*

```
63      Return @RC;
64    End
65    Go
66
67
68    --***********************************************************************--
69    -- A) Drop the FOREIGN KEY CONSTRAINTS and Clear the tables
70     -- NOT NEEDED FOR INCREMENTAL LOADING:
71    --***********************************************************************--
72
73
74    --***********************************************************************--
75    -- B) Synchronize the Tables
76    --***********************************************************************--
77
78    /****** [dbo].[DimCustomers] ******/
79    go
80    Create or Alter View vETLNewPatient
81    /* Author: Yuanlong Zhang
82    ** Desc: Extracts and transforms data for DimCustomers
83    ** Change Log: When,Who,What
84    ** 2022-01-25,Yuanlong Zhang,Created Sproc (MERGE).
85    */
86    As
87      Select [FName] = [FName]
88             ,[LName] = [LName]
89             ,[Email] = [Email]
90             ,[Address] = [Address]
91             ,[City] = [City]
92             ,[State] = [State]
93                   ,[ZipCode] = [ZipCode]
94        FROM [Staging].[dbo].[StagingNewPatient]
95    go
96    /* Testing Code:
97     SELECT * FROM vETLNewPatient;
98    */
99
100
101
102    go
103    Create or Alter Procedure pETLSyncPatients
104    /* Author: Yuanlong Zhang
105    ** Desc: Inserts data into DimCustomers
106    ** Change Log: When,Who,What
107    ** 2022-01-24,Yuanlong Zhang, Created Sproc (MERGE).
108    */
109    As
110    Begin
111      Declare @RC int = 0;
112            Begin Try
113         -- ETL Processing Code --
114         Merge Into Patients as t
115          Using vETLNewPatient as s -- For Merge to work with SCD tables, I need to
116    insert a new row when the following is not true:
117          On  t.FName = s.FName
118          And t.LName = s.LName
119          And t.Email = s.Email
120        When Not Matched -- At least one column value does not match add a new
121    row:
122         Then
123          Insert (FName, LName, Email, Address, City, State, ZipCode)
124           Values (s.FName
125                  ,s.LName
126                  ,s.Email
127                  ,s.Address
128                          ,s.City
129                          ,s.State
130                  ,s.ZipCode)
131         When Matched
132             AND (t.Address <> s.Address OR t.City <> s.City OR t.State <>
```

*Created view for staging purposes.*

*Use MERGE statement for incremental loading into the database table*

```
133   s.State OR t.ZipCode <> s.ZipCode)
134            -- If there is a row in the target (dim) table that is no longer in
135   the source table
136         Then -- indicate that row is no longer current
137          Update
138           Set t.Address = s.Address
139             ,t.City = s.City
140                       ,t.State = s.State
141                       ,t.ZipCode = s.ZipCode
142      ;
143
144      -- ETL Logging Code --
145              Exec pInsETLLog
146              @ETLAction = 'pETLSyncPatients'
147             ,@ETLLogMessage = 'Patients synced';
148              Set @RC = +1
149         End Try
150
151         -- Error Handling Code --
152         Begin Catch
153             Declare @ErrorMessage nvarchar(1000) = Error_Message
154
155      -- ETL Logging Code --
156             Exec pInsETLLog
157            @ETLAction = 'pETLSyncPatients'
158           ,@ETLLogMessage = @ErrorMessage;
159             Set @RC = -1
160         End Catch
161         Return @RC;
162   End
163   go
164   /* Testing Code:
165    Declare @Status int;
166    Exec @Status = pETLSyncDimCustomers;
167    Print @Status;
168   */
169   go
170
171   /*
172   go
173   Declare @Status int = 0;
174   Exec @Status = pETLSyncPatients;
175   Select [Object] = 'pETLSyncPatients', [Status] = @Status;
176   select * from dbo.ETLLog
177
178   select * from Patients.dbo.Patients
179   --WHERE fname = 'Michael'
180   */
```

*Logging the ETL processes into the ETLLog table.*

*You can always **uncomment** these code for **debugging and testing** purposes.*

### 6.2.2.   DWETLObjects.sql (Milestone 2)

**🞣   Note:**
- The SQL Script is very long. Thus, we only pick up one of these tables as an example.

```
1    --**********************************************************************--
2    -- Title: DWFinal
3    -- Author: Yuanlong Zhang
4    -- Desc: This file tests your knowledge on how to create an Incremental ETL
5    process with SQL code
6    -- Change Log: When,Who,What
7    -- 2021-01-17,RRoot,Created File
8    -- 2022-01-24,Yuanlong Zhang, Completed File
9    -- 2022-03-07,Yuanlong Zhang, Updated the file for Final Projects
10
```

```
11   --**************************************************************--
12   IF NOT EXISTS (SELECT 1 from sys.servers where name =
13   'continuumsql.westus2.cloudapp.azure.com')
14   BEGIN
15     EXEC sp_addlinkedserver @server = 'continuumsql.westus2.cloudapp.azure.com'
16     EXEC sp_addlinkedsrvlogin 'continuumsql.westus2.cloudapp.azure.com'
17                               ,'false'
18                               ,NULL
19                               ,'BICert'
20                               ,'BICert'
21   END
22
23   IF NOT EXISTS (SELECT 1 from sys.servers where name = 'is-
24   root01.ischool.uw.edu\BI')
25   BEGIN
26     EXEC sp_addlinkedserver @server = 'is-root01.ischool.uw.edu\BI'
27     EXEC sp_addlinkedsrvlogin 'is-root01.ischool.uw.edu\BI'
28                               ,'false'
29                               ,NULL
30                               ,'BICert'
31                               ,'BICert'
32   END
33
34
35   USE [DWClinicReportDataYuanlongZhang];
36   go
37   SET NoCount ON;
38
39   --  Setup Logging Objects ----------------------------------------------
40
41                        Script for ETLLog table is omitted
42
43   --**************************************************************--
44   -- A) Drop the FOREIGN KEY CONSTRAINTS and Clear the tables
45    -- NOT NEEDED FOR INCREMENTAL LOADING:
46   --**************************************************************--
47
48
49   --**************************************************************--
50   -- B) Synchronize the Tables
51   --**************************************************************--
52                        Script for DimDate table is omitted
53
54
55   /****** [dbo].[DimDoctors] ******/
56   go
57   Create or Alter View vETLDimDoctors
58   /* Author: Yuanlong Zhang
59   ** Desc: Extracts and transforms data for DimDoctors
60   ** Change Log: When,Who,What
61   ** 2022-01-24,Yuanlong Zhang,Created Sproc (MERGE).
62   ** 2022-03-07,Yuanlong Zhang,Created Sproc (MERGE) for final project.
63   */
64   As
65     Select
66              [DoctorID] = d.DoctorID
67             ,[DoctorFullName] = d.FirstName + ' ' + d.LastName
68             ,[DoctorEmailAddress] = d.EmailAddress
69             ,[DoctorCity] = TRIM(d.City)
70             ,[DoctorState] = CASE WHEN LEN(TRIM(d.State)) > 2 THEN '(ERROR) ' +
71   TRIM(d.State) ELSE TRIM(d.State) END
72             ,[DoctorZip] = d.Zip
73     From
74   [continuumsql.westus2.cloudapp.azure.com].[DoctorsSchedules].[dbo].[Doctors] d
75     UNION
76     Select -1, 'No Doctor', 'N/A', 'N/A', 'N/A', 000000
77   go
78   /* Testing Code:
79    Select * From vETLDimDoctors;
80   */
```

*Define external server connections*

*Created view for ETL purposes.*

*To avoid violation of none null value in doctor, union a "no doctor" column.*

```
81
82   go
83   Create or Alter Procedure pETLSyncDimDoctors
84   /* Author: Yuanlong Zhang
85   ** Desc: Updates data in DimDoctors using the vETLDimDoctors view
86   ** Change Log: When,Who,What
87   ** 2022-01-24,Yuanlong Zhang, Created Sproc (MERGE).
88   ** 2022-03-07,Yuanlong Zhang, Created Sproc (MERGE) for the final project.
89   */
90   AS
91   Begin
92          Declare @RC int = 0;
93     Begin Try
94       -- ETL Processing Code --
95       Merge Into DimDoctors as t
96        Using vETLDimDoctors as s -- For Merge to work with SCD tables, I need to
97   insert a new row when the following is not true:
98        On  t.DoctorID = s.DoctorID
99       When Not Matched -- At least one column value does not match add a new
100  row:
101       Then
102        Insert (DoctorID, DoctorFullName,
103  DoctorEmailAddress,DoctorCity,DoctorState, DoctorZip)
104        Values (s.DoctorID
105                       ,s.DoctorFullName
106            ,s.DoctorEmailAddress
107            ,s.DoctorCity
108            ,s.DoctorState
109                       ,s.DoctorZip)
110       When Matched
111          AND t.DoctorEmailAddress <> s.DoctorEmailAddress
112          OR t.DoctorCity <> s.DoctorCity
113          OR t.DoctorState <> s.DoctorState
114          OR t.DoctorZip <> s.DoctorZip-- If there is a row in the target
115  (dim) table that is no longer in the source table
116        Then -- indicate that row is no longer current
117         Update
118         Set t.DoctorEmailAddress = s.DoctorEmailAddress
119            ,t.DoctorCity = s.DoctorCity
120                       ,t.DoctorState = s.DoctorState
121                       ,t.DoctorZip = s.DoctorZip
122       When Not Matched by Source
123       Then
124        Update
125         Set t.DoctorFullName =
126  iif(patindex('%(Deleted)%',[DoctorFullName]) > 0, [DoctorFullName],
127  [DoctorFullName] + ' (Deleted)')
128      ;
129
130     -- ETL Logging Code --
131             Exec pInsETLLog
132             @ETLAction = 'pETLSyncDimDoctors'
133            ,@ETLLogMessage = 'DimDoctors synced';
134             Set @RC = +1
135       End Try
136       Begin Catch
137             Declare @ErrorMessage nvarchar(1000) = Error_Message();
138             Exec pInsETLLog
139           @ETLAction = 'pETLSyncDimDoctors'
140          ,@ETLLogMessage = @ErrorMessage;
141             Set @RC = -1
142       End Catch
143       Return @RC;
144  End
145  go
146  /* Testing Code:
147   Declare @Status int;
148   Exec @Status = pETLSyncDimDoctors;
149   Print @Status;
150   Select * From DimDoctors
```

*MERGE Statement for
incremental loading.
We use doctor ID to
identify the unique record*

*To avoid violate the FK
constraints, we add
"(deleted)" when we
detected the column had
been deleted*

```
151  */
152
153
154                    Script for other table is omitted
155
156
157  /****** [dbo].[DimPatients] ******/
158  go
159  Create or Alter View vETLDimPatients
160  /* Author: Yuanlong Zhang
161  ** Desc: Extracts and transforms data for DimPatients - SCD-2
162  ** Change Log: When,Who,What
163  ** 2022-01-25,Yuanlong Zhang,Created Sproc (MERGE).
164  ** 2022-03-07,Yuanlong Zhang,Created Sproc (MERGE) for final project.
165  */
166  As
167    Select [PatientID] = p.ID
168          ,[PatientFullName] = Cast(p.FName+ ' ' + p.LName as nvarchar(100))
169          ,[PatientCity] = p.City
170          ,[PatientState] = p.State
171          ,[PatientZipCode] = p.ZipCode
172      FROM [is-root01.ischool.uw.edu\BI].[Patients].[dbo].[Patients] p
173  go
174  /* Testing Code:
175   Select * From vETLDimPatients;
176  */
177
178  go
179  Create or Alter Procedure pETLSyncDimPatients
180  /* Author: Yuanlong Zhang
181  ** Desc: Inserts data into DimPatients  - Type 2 SCD
182  ** Change Log: When,Who,What
183  ** 2022-01-24,Yuanlong Zhang, Created Sproc (MERGE).
184  ** 2022-03-07,Yuanlong Zhang, Created Sproc (MERGE) for the final project.
185  */
186  As
187  Begin
188    Declare @RC int = 0;
189          Begin Try
190      -- ETL Processing Code --
191      Merge Into DimPatients as t
192       Using vETLDimPatients as s -- For Merge to work with SCD tables, I need
193  to insert a new row when the following is not true:
194        On  t.PatientID = s.PatientID
195        And t.PatientFullName = s.PatientFullName
196        And t.PatientCity = s.PatientCity
197        And t.PatientState = s.PatientState
198            And t.PatientZipCode = s.PatientZipCode
199      When Not Matched -- At least one column value does not match add a new
200  row:
201        Then
202          Insert (PatientID, PatientFullName, PatientCity, PatientState,
203  PatientZipCode,
204                  StartDate, EndDate, IsCurrent)
205          Values (s.PatientID
206                  ,s.PatientFullName
207                  ,s.PatientCity
208                  ,s.PatientState
209                         ,s.PatientZipCode
210                  ,Cast(Convert(nvarchar(100), GetDate(), 112) as date) -- Smart
211  Key can be joined to the DimDate
212                  ,Null
213                  ,1)
214      When Not Matched By Source -- If there is a row in the target (dim)
215  table that is no longer in the source table
216         Then -- indicate that row is no longer current
217          Update
218           Set t.EndDate = Cast(Convert(nvarchar(100), GetDate(), 112) as date)
219  -- Smart Key can be joined to the DimDate
220               ,t.IsCurrent = 0
```

*Patient table use SCD Type-2.*
*We added 3 columns for StartDate, EndDate and IsCurrent.*
*The date will get system time when executed.*

*When deletion, we update the EndDate and marked iscurrent as 0.*

```
221        ;
222
223        -- ETL Logging Code --
224                Exec pInsETLLog
225                @ETLAction = 'pETLSyncDimPatients'
226               ,@ETLLogMessage = 'DimPatients synced';
227                Set @RC = +1
228            End Try
229
230            -- Error Handling Code --
231            Begin Catch
232                Declare @ErrorMessage nvarchar(1000) = Error_Message();
233
234        -- ETL Logging Code --
235                Exec pInsETLLog
236               @ETLAction = 'pETLSyncDimPatients'
237              ,@ETLLogMessage = @ErrorMessage;
238                Set @RC = -1
239            End Catch
240            Return @RC;
241    End
242    go
243    /* Testing Code:
244     Declare @Status int;
245     Exec @Status = pETLSyncDimPatients;
246     Print @Status;
247     select * from dimpatients
248    */
249    go
250
                        Script for other table is omitted
```

### 6.2.3.  DWMongoDBView.sql (Milestone 3)

The two views have a similar structure. Thus, we use vRptDoctorShifts as an example:

```
1   --**********************************************************************--
2   -- Title: DWFinal-DWMongoDBView
3   -- Author: Yuanlong Zhang
4   -- Desc: This file used to create the view to feed MongoDB
5   -- Change Log: When,Who,What
6   -- 2021-01-17,RRoot,Created File
7   -- 2022-01-24,Yuanlong Zhang, Completed File
8   -- 2022-03-10,Yuanlong Zhang, Updated the file for Final Projects
9
10  --**********************************************************************--
11
12  USE DWClinicReportDataYuanlongZhang;
13  GO
14
15  IF (Object_ID('vRptDoctorShifts') is not null) Drop View vRptDoctorShifts;
16  GO
17  Create View vRptDoctorShifts
18  AS
19  Select
20  [ShiftDate] = Cast(Cast([FullDate] as date) as varchar(100))
21  ,[ClinicName] = dc.ClinicName
22  ,[ClinicCity] = dc.ClinicCity
23  ,[ClinicState] = dc.ClinicState
24  ,[ShiftID] = ds.ShiftID
25  ,[ShiftStart] = ds.ShiftStart
26  ,[ShiftEnd] = ds.ShiftEnd
27  ,[DoctorFullName] = ddo.DoctorFullName
28  ,[HoursWorked] = fds.HoursWorked
29  FROM
30  dbo.FactDoctorShifts fds
31  JOIN dbo.DimDates dd
```

*To avoid dirty read, drop the view every time.*

*Convert all the data into string to avoid data type issue*

```
32   ON fds.ShiftDateKey = dd.DateKey
33   JOIN dbo.DimClinics dc
34   ON fds.ClinicKey = dc.ClinicID
35   JOIN dbo.DimShifts ds
36   ON fds.ShiftKey = ds.ShiftID
37   JOIN dbo.DimDoctors ddo
38   ON fds.DoctorKey = ddo.DoctorKey
39
```

### 6.2.4. ETLJob.sql (Milestone 4)

```
1    --**************************************************************************--
2    -- Title: Create the DW ETL Job
3    -- Desc: This file will drop and create a SQL Agent Job
4    -- Change Log: When,Who,What
5    -- 2020-01-01,RRoot,Created File
6    -- 2022-03-14,Yuanlong Zhang, Updated file for Final Project
7    --**************************************************************************--
8
9
10   USE [master]
11   GO
12   Begin Try
13   -- Access to the Server
14   CREATE LOGIN [DESKTOP-AO8N3T1\i_ecn]
15    FROM WINDOWS
16     WITH DEFAULT_DATABASE=[master], DEFAULT_LANGUAGE=[us_english]
17
18   ALTER SERVER ROLE [sysadmin] ADD MEMBER [DESKTOP-AO8N3T1\i_ecn]
19   End Try
20   Begin Catch
21         Print Error_Message()
22   End Catch
23   GO
24
25   -- Abstaction layer to the Login
26   Begin Try
27   CREATE CREDENTIAL [CredentialForETLAutomations]
28    WITH IDENTITY = N'DESKTOP-AO8N3T1\i_ecn'
29   End Try
30   Begin Catch
31         Print Error_Message()
32   End Catch
33   GO
34
35   -- Connection to an account with enough permission
36   -- using the abstraction layer credentials
37   Begin Try
38         EXEC msdb.dbo.sp_add_proxy
39          @proxy_name=N'SSIS Proxy'
40         ,@credential_name=N'CredentialForETLAutomations'
41         ,@enabled=1
42
43         -- Map to SSIS subsystems
44         EXEC msdb.dbo.sp_grant_proxy_to_subsystem
45          @proxy_name=N'SSIS Proxy'
46         ,@subsystem_id=11  -- SSIS Package
47   End Try
48   Begin Catch
49         Print Error_Message()
50   End Catch
51   GO
52
53
54   USE [msdb]
55   GO
56   BEGIN TRY
57     IF Exists (Select * from SysJobs Where Name = 'ETLDWClinicReportData')
```

*Create Login for local account.*
*You have to update the account name on your computer.*

*Create a credential for he identity we set up before.*

*Create a SSIS Proxy for SQL Agent Jobs;*
*Use Credential which we created before*

```
58        Begin
59          Exec sp_delete_job @job_name = ETLDWClinicReportData
60        End
61
62    /****** Object:   Job [DWClinicReportDataYuanlongZhang]      Script Date:
63  8/21/2021 3:46:14 PM ******/
64    BEGIN TRANSACTION
65    DECLARE @ReturnCode INT
66    SELECT @ReturnCode = 0
67    /****** Object:   JobCategory [[Uncategorized (Local)]]      Script Date:
68  8/21/2021 3:46:14 PM ******/
69    IF NOT EXISTS (SELECT name FROM msdb.dbo.syscategories WHERE
70  name=N'[Uncategorized (Local)]' AND category_class=1)
71    BEGIN
72    EXEC @ReturnCode = msdb.dbo.sp_add_category @class=N'JOB', @type=N'LOCAL',
73  @name=N'[Uncategorized (Local)]'
74    IF (@@ERROR <> 0 OR @ReturnCode <> 0) GOTO QuitWithRollback
75
76    END
77
78    DECLARE @jobId BINARY(16)
79    EXEC @ReturnCode =  msdb.dbo.sp_add_job @job_name=N'ETLDWClinicReportData',
80                      @enabled=1,
81                      @notify_level_eventlog=0,
82                      @notify_level_email=0,
83                      @notify_level_netsend=0,
84                      @notify_level_page=0,
85                      @delete_level=0,
86                      @description=N'Performs ETL tasks for
87  ETLDWClinicReportData',
88                      @category_name=N'[Uncategorized (Local)]',
89                      @owner_login_name=N'sa', @job_id = @jobId OUTPUT
90    IF (@@ERROR <> 0 OR @ReturnCode <> 0) GOTO QuitWithRollback
91    /****** Object:   Step [Run DWIndependentBookSellersETLpackage.dtsx]
92  Script Date: 8/21/2021 3:46:14 PM ******/
93    EXEC @ReturnCode = msdb.dbo.sp_add_jobstep @job_id=@jobId, @step_name=N'Run
94  ETLJob.dtsx',
95                      @step_id=1,
96                      @cmdexec_success_code=0,
97                      @on_success_action=1,
98                      @on_success_step_id=0,
99                      @on_fail_action=2,
100                     @on_fail_step_id=0,
101                     @retry_attempts=0,
102                     @retry_interval=0,
103                     @os_run_priority=0, @subsystem=N'SSIS',
104                     @command=N'/FILE
105 "C:\_BISolutions\ETLFinal_YuanlongZhang\ETLPackages\ETLJob.dtsx"
106 /CHECKPOINTING OFF /REPORTING E',
107                     @database_name=N'master',
108                     @flags=0,
109                     @proxy_name=N'SSIS Proxy'
110   IF (@@ERROR <> 0 OR @ReturnCode <> 0) GOTO QuitWithRollback
111   EXEC @ReturnCode = msdb.dbo.sp_update_job @job_id = @jobId, @start_step_id =
112 1
113   IF (@@ERROR <> 0 OR @ReturnCode <> 0) GOTO QuitWithRollback
114   EXEC @ReturnCode = msdb.dbo.sp_add_jobschedule @job_id=@jobId,
115 @name=N'EachNight',
116                     @enabled=1,
117                     @freq_type=4,
118                     @freq_interval=1,
119                     @freq_subday_type=1,
120                     @freq_subday_interval=0,
121                     @freq_relative_interval=0,
122                     @freq_recurrence_factor=0,
123                     @active_start_date=20210821,
124                     @active_end_date=99991231,
125                     @active_start_time=10000,
126                     @active_end_time=235959,
127                     @schedule_uid=N'ac7412ed-e42f-46a0-a8bb-d16ccf0310fb'
```

*Deleted SQL Agent Job if existed*

*Created a new job which execute the ETLJob.dtsx*

*Scheduled the job to run each night.*

```
128    IF (@@ERROR <> 0 OR @ReturnCode <> 0) GOTO QuitWithRollback
129    EXEC @ReturnCode = msdb.dbo.sp_add_jobserver @job_id = @jobId, @server_name
130 = N'(local)'
131    IF (@@ERROR <> 0 OR @ReturnCode <> 0) GOTO QuitWithRollback
132    COMMIT TRANSACTION
133    GOTO EndSave
134    QuitWithRollback:
135        IF (@@TRANCOUNT > 0) ROLLBACK TRANSACTION
136    EndSave:
137
138 END TRY
139 BEGIN CATCH
140    IF (@@TRANCOUNT > 0) ROLLBACK TRANSACTION
141    Print Error_Message()
142 END CATCH
143
144 GO
145
```

### 6.2.5.  ETLViews.sql (Milestone 4)

```
1  --***************************************************************--
2  -- Title: Testing the Reporting Views
3  -- Author: RRoot, Yuanlong Zhang
4  -- Desc: This file creates several ETL views used in Admin reports
5  -- Change Log: When,Who,What
6  -- 2018-02-07,RRoot,Created File
7  -- 2022-03-15, Yuanlong Zhang, Updated for the final project
8  --***************************************************************--
9  Use DWClinicReportDataYuanlongZhang;
10
11
12 --Create views for SSIS Job
13 Create or Alter View vDWClinicReportDataYuanlongZhangETLJobHistory
14 As
15 Select Top 100000
16  [JobName] = j.name
17 ,[StepName] = h.step_name
18 ,[RunDateTime] = msdb.dbo.agent_datetime(run_date, run_time)
19 ,[RunDurationSeconds] = h.run_duration
20 ,[RunStatus] = iif(h.run_status = 1, 'Success', 'Failure')
21 From msdb.dbo.sysjobs as j
22    Inner Join msdb.dbo.sysjobhistory as h
23      ON j.job_id = h.job_id
24 --Where j.enabled = 1 And j.name = 'ETLDWClinicReportData'
25 Order by JobName, RunDateTime desc;
26
27
28 --Create a view for row count reports
29 Create or Alter View DWClinicReportDataRowCounts
30 As
31 With [RowCounts] -- Using a CTE to access the Top Command for the Order By
32 statement in the view
33 As(
34 Select [SortCol] = 1, [TableName] = 'DimDates', [CurrentNumberOfRows] =
35 Count(*) From [DimDates]
36 Union
37 Select [SortCol] = 2, [TableName] = 'DimClinics', [CurrentNumberOfRows] =
38 Count(*) From [DimClinics]
39 Union
40 Select [SortCol] = 3, [TableName] = 'DimDoctors', [CurrentNumberOfRows] =
41 Count(*) From [DimDoctors]
42 Union
43 Select [SortCol] = 4, [TableName] = 'DimPatients', [CurrentNumberOfRows] =
44 Count(*) From [DimPatients]
45 Union
46 Select [SortCol] = 5, [TableName] = 'DimProcedures', [CurrentNumberOfRows] =
47 Count(*) From [DimProcedures]
48 Union
```

*Created a SQL Agent Job log view from MSDB.*

*Convert the running status from 1,0 to success and failure.*

*Create rowcount table by union all the individual row count stats*

```
49  Select [SortCol] = 6, [TableName] = 'DimShifts', [CurrentNumberOfRows] =
50  Count(*) From [DimShifts]
51  Union
52  Select [SortCol] = 7, [TableName] = 'FactDoctorShifts', [CurrentNumberOfRows]
53  = Count(*) From [FactDoctorShifts]
54  Union
55  Select [SortCol] = 8, [TableName] = 'FactVisits', [CurrentNumberOfRows] =
56  Count(*) From [FactVisits]
57  Union
58  Select [SortCol] = 9, [TableName] = 'ETLLog', [CurrentNumberOfRows] = Count(*)
59  From [ETLLog]
60  )
61  Select Top 100000 [SortCol],[TableName],[CurrentNumberOfRows]
62    From [RowCounts]
63    Order By [SortCol] asc; -- Use a sort column, so it does not sort by table
64  name.
65  go
66
```

*Sorted by SortCol rather than by table name.*

# *7. Summary*

This manual provides an overview of our ETL process.

You must carefully follow <mark>section 2 Checklist</mark> and ensure the zip has been <mark>unzipped into</mark> <mark>C:\_BISolutions\ETLFinal_YuanlongZhang</mark> before starting.
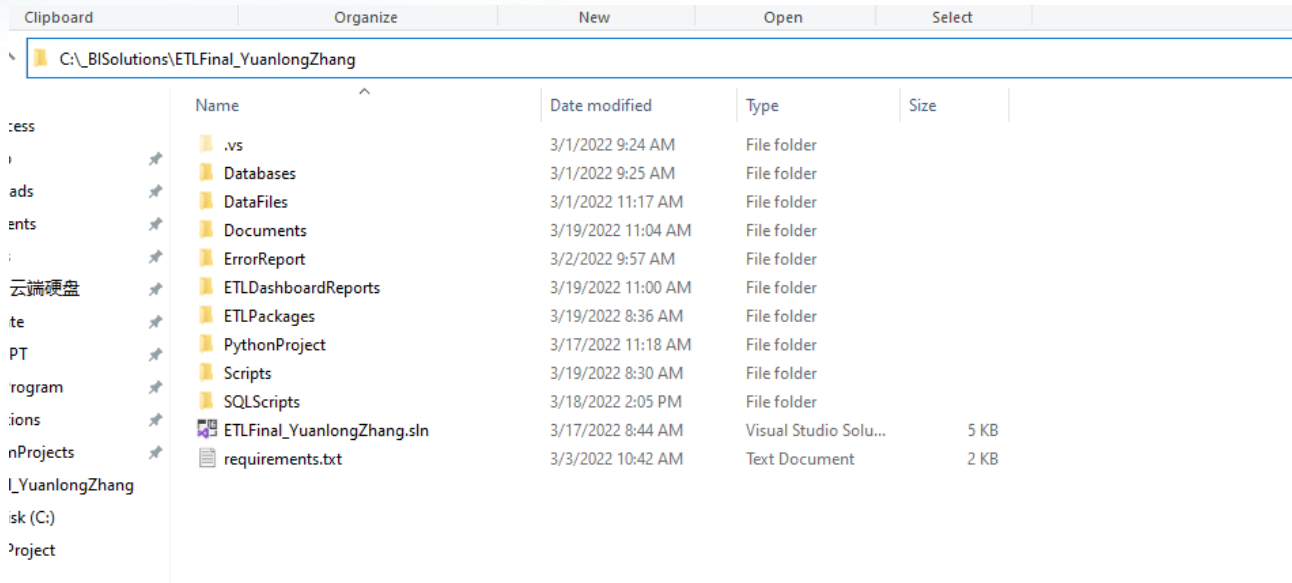


*Figure 57 ETL Solution Folder*

You can search the keywords to locate the topic you are interested in quickly.

For more information, don't hesitate to contact Yuanlong Zhang via email: yuanlong@uw.edu.