# ETL Technical Manual

*Version: 1.0*
*Date: 03/25/2025*
*Author: Girum Obse*
*UW ETL Processing Final*

## Table of Contents

# 1. Introduction

## 1.1 Purpose of the Document

This document serves as a comprehensive technical guide for understanding and implementing ETL (Extract, Transform, Load) processes, of a small medical clinic data analytics project. The project involves automating the data ingestion, transformation, and loading processes from CSV files to a structured Data Warehouse. The ETL solution prioritizes scalability, efficiency, and maintainability using SQL Server, SSIS, and Python. This manual provides step-by-step instructions on ETL development, best practices, and troubleshooting techniques.

Additionally, it outlines the four milestones of the project, covering file imports, data warehouse integration, non-SQL ETL processes, and automation. Security, compliance, and performance optimization strategies are also discussed to ensure smooth data flow between systems. Whether building an ETL pipeline for the first time or optimizing an existing one, this document serves as a valuable reference for designing scalable and efficient ETL solutions.

## Current Situation

Currently, individual clinics send daily CSV files to the corporate office, where they are manually reviewed, cleaned, and added to one of two databases. This manual ETL process is time-consuming, prone to human error, and lacks scalability, especially as the number of clinics increases. The process delays data availability, affecting decision-making and efficiency.

# 2. Project Overview

The ETL project automates data movement from daily CSV uploads into structured databases, providing analytical insights. The existing system is entirely manual, and this solution introduces automation to streamline data integration.

- **Source Systems**: CSV files from clinics

- **Destination Systems**: Patients and DoctorsSchedules databases, DWClinicReportData (Data Warehouse), and Excel reports

- **ETL Tools**: SQL Server, SSIS, Python

- **Deployment**: Visual Studio-based structured solution

## 2.1 Current Workflow over view

**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 <u>current</u> ETL process is **entirely manual**.
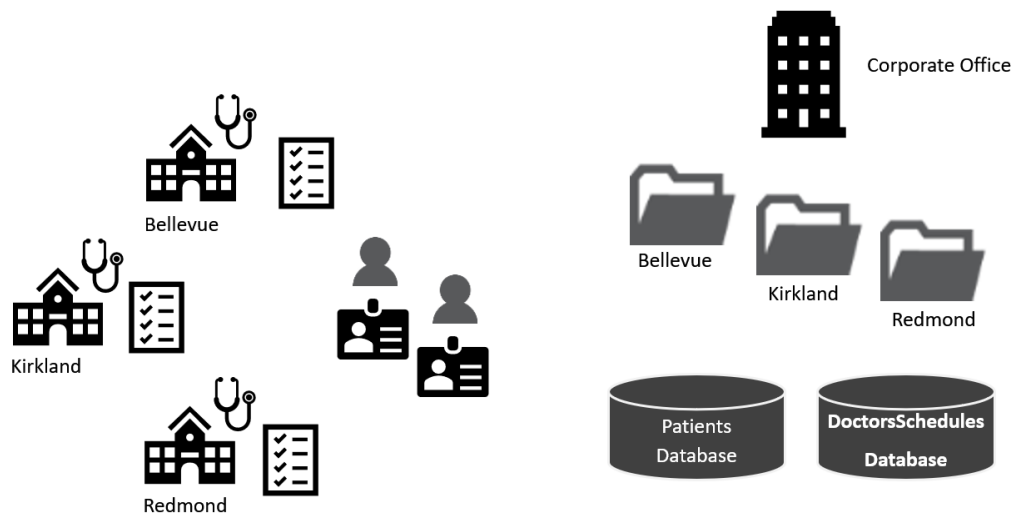
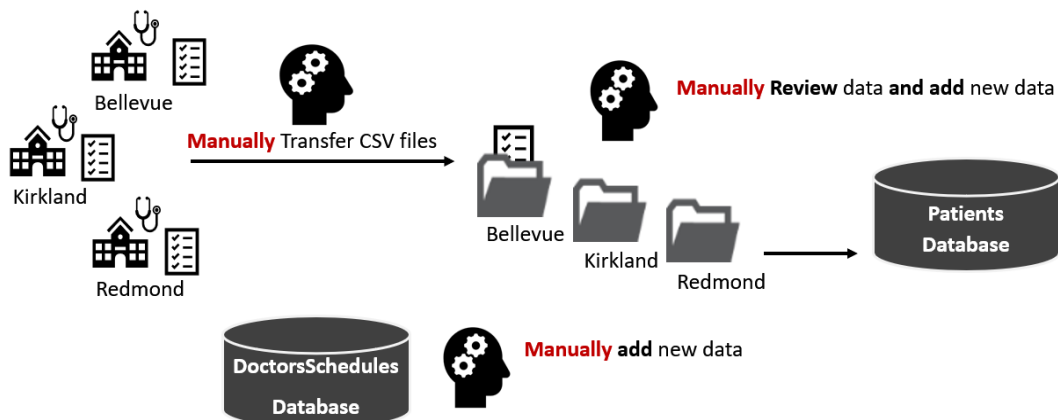

Figure 1. Topology of the current design.



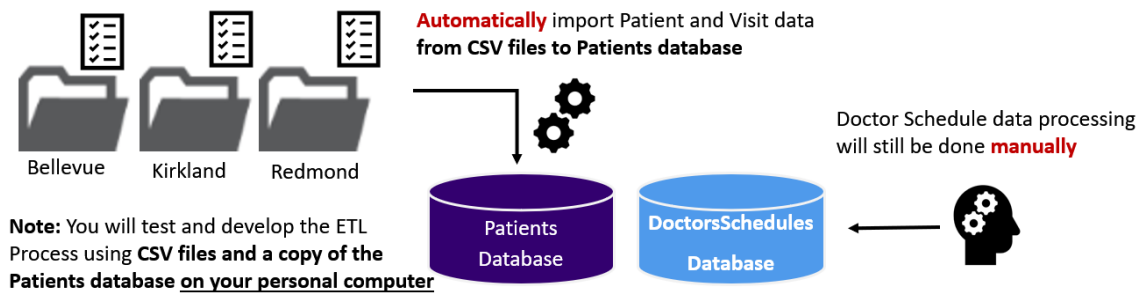Figure 2. Visualization of the manual ETL file process

Figure 3. Visualization of the manual ETL file process

## 2.2,Proposed ETL process

A new ETL process will automate parts of this workflow, with potential for full automation in the future. Below, you'll find detailed descriptions of the entire process and steps, along with sample code snippets and screenshots for clarity.

## 2.3 Audience

This manual is intended for a wide range of professionals involved in data management and analytics, including:
- ETL Developers – Those responsible for designing and implementing ETL pipelines.
- Data Engineers – Professionals who manage data pipelines and ensure data flows correctly.
- Database Administrators – Responsible for maintaining data integrity and optimizing database performance.
- Business Intelligence Analysts – Those who analyze and report on business data.
- IT Managers – Decision-makers overseeing data-related projects.
This document assumes some familiarity with databases, SQL, and basic scripting but does not require advanced ETL experience.

## 2.4 ETL Overview

### 2.4.1 What is ETL?
ETL stands for Extract, Transform, and Load—a three-step process used for integrating data from multiple sources into a centralized data warehouse or analytics system. The Extract

phase retrieves data from various sources such as databases, web APIs, and flat files. The Transform phase applies business rules to clean, standardize, and enrich the data. The Load phase ensures the data is stored in the target system, ready for reporting and analysis. ETL is widely used in data warehousing, analytics, and business intelligence, ensuring that organizations work with high-quality, consistent, and structured data.

## 2.4.1 Why ETL is Important?

ETL plays a crucial role in modern data-driven organizations by enabling them to consolidate data from multiple sources and prepare it for analysis. Without ETL, businesses would struggle with inconsistent, incomplete, and redundant data. ETL ensures that data is cleaned, standardized, and formatted properly before it reaches reporting tools. Additionally, ETL processes support historical data storage, allowing companies to perform trend analysis, forecasting, and strategic decision-making. ETL also enhances data governance and compliance, ensuring that sensitive information is processed securely according to industry regulations such as GDPR and HIPAA.

## 3.Environment Setup

Before beginning ETL development, ensure the following components are installed and configured:

- SQL Server with Management Studio

- Visual Studio with SSIS extensions

- Python 3.9 with required libraries

- Access to database backup files (.BAK) for restoration

- CSV files for testing in the C:\_BISolutions\ClinicDailyData directory

---

## 4. Milestone 1: File-Based ETL

**Data Review**

- Assess CSV file structures for inconsistencies (e.g., column names, data types, missing values).

**Setting Up the Development Environment**

- Create a Visual Studio solution with organized folders.

**Database Restoration**

Execute the following SQL script to restore databases:

ALTER DATABASE [Patients] SET SINGLE_USER WITH ROLLBACK IMMEDIATE;

RESTORE DATABASE [Patients]

 FROM DISK = N'C:/_BISolutions/Databases/Patients.bak'

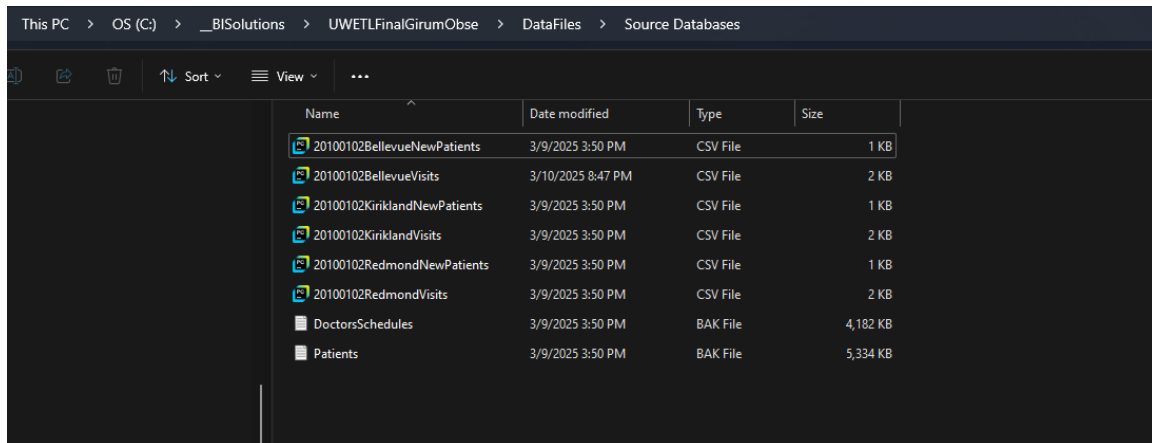 WITH RECOVERY, REPLACE;

ALTER DATABASE [Patients] SET MULTI_USER;

GO

## Data Transformation Documentation

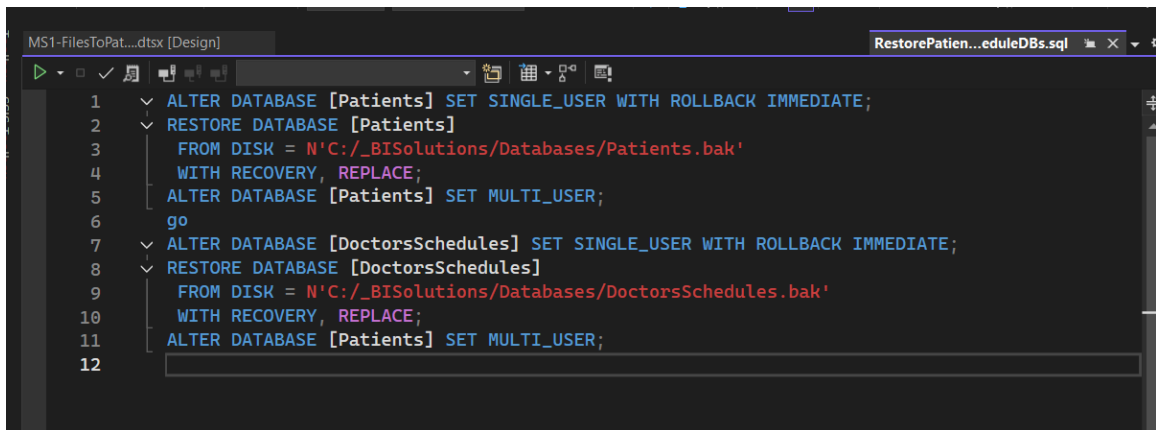- Maintain an Excel metadata sheet tracking transformations.

## SSIS Package Implementation

- Develop an SSIS package (ETLFilesToDatabases.dtsx) to automate CSV imports.



Figure 4 Excel source file for different clinics

```
MS1-FilesToPat....dtsx [Design]                                    RestorePatien...eduleDBs.sql

1    ∨ ALTER DATABASE [Patients] SET SINGLE_USER WITH ROLLBACK IMMEDIATE;
2    ∨ RESTORE DATABASE [Patients]
3        FROM DISK = N'C:/_BISolutions/Databases/Patients.bak'
4        WITH RECOVERY, REPLACE;
5      ALTER DATABASE [Patients] SET MULTI_USER;
6      go
7    ∨ ALTER DATABASE [DoctorsSchedules] SET SINGLE_USER WITH ROLLBACK IMMEDIATE;
8    ∨ RESTORE DATABASE [DoctorsSchedules]
9        FROM DISK = N'C:/_BISolutions/Databases/DoctorsSchedules.bak'
10       WITH RECOVERY, REPLACE;
11     ALTER DATABASE [Patients] SET MULTI_USER;
12
```

Figure5, sql code script showing restoring the patients and doctorschedule database



Fiqure6, showing doctorschedule and patients database tables

```sql
317        Create or Alter Proc pETLInsertVisitsData
318          (@Date date)
319          As
320          /*****************************************************************************************************
321          Desc: Inserts Data from Staging Tables into Visits
322          Dev:GObse
323          Date:03/05/2025
324          *****************************************************************************************************/
325      Begin
326          Declare @RC int =0;
327          Begin Try
328          Begin Tran;
329          Insert Into [Patients].[dbo].[Visits]
330          ([Date],[Clinic],[Patient],[Doctor],[Procedure],[Charge])
331          Exec pETLSelectVistsData @Date = @Date --using my select sproc
332          Commit Tran;
333          Set @RC = 1;
334          End Try
335          Begin Catch
336              Rollback Tran;
337              Print 'Error Inserting Data from Staging Tables into Visits'
338              Print ERROR_MESSAGE()
339              Set @RC = -1;
340          End Catch
341          Return @RC;
342      End
343      go
344      Exec pETLSelectVistsData @Date = '20100102'
345      go
346
347      --Test Code--
348      USE [master]
349      go
350      -- reset the database
351      ALTER DATABASE [patients] SET SINGLE_USER WITH ROLLBACK IMMEDIATE;
352      RESTORE DATABASE [Patients]
353        FROM DISK = N'C:/_BISOLUTIONS/Databases/Patients.bak'
354        WITH RECOVERY, REPLACE;
355      ALTER DATABASE [Patients] SET MULTI_user;
356      go
357      Select Count(*) From [Patients].[dbo].[Visits] --40150
358      go
359      Use tempdb;
360      go
361      Exec pETLCreateStagingTables
362      go
363      Exec pETLImportFileDataToStagingTables
364      go
365      Exec pETLTransformVistsData
366      go
367      Exec pETLInsertVisitsData @Date = '20100102'
368      go
369      Select Count(*) From [Patients].[dbo].[Visits] --40304
370      go
371
```

Figure7: part of  sql  script showing how data loaded to databases

Figure 8, ssis package for dataloaded from excel source to visits table in patients database



Figure9, showing screenshot of scheduled job to load data to patients database tables

Figure 10, showing data in patients database tables

# 5. Milestone 2: Data Warehouse ETL

**Process Steps**

**Creating the Data Warehouse**

**- Run the provided SQL script (`\Scripts\CreateDWClinicReportData.sql`) to create the **DWClinicReportData** database.**

Figure11, sql code script of the datawarehouse

- Verify that all required tables, schemas, and constraints are properly created.

- Check for any errors or missing objects in the database after execution.

Reviewing the Database Design

- Identify the **dimension tables** (e.g., DimPatients, DimDoctors, DimShifts) and fact tables(e.g., FactDoctorShifts and FactVisits).

- Understand relationships between tables, primary and foreign keys.

**SQL ETL and SSIS Implementation**

- **Create SQL Views and Stored Procedures**:

  - Develop **views** to extract necessary data from the **DoctorsSchedules** and **Patients** OLTP databases.

  - Write **stored procedures** to transform and load the data into the **DWClinicReportData** warehouse.

  - Implement **data cleansing** to remove duplicates, handle null values, and ensure consistency.

Figure12, showing ETL script of data load from OLTP database to the OLAP database example of Dimdoctor table



- **Develop an SSIS Package (DWClinicReportDataETL.dtsx)**:

- Create an **SSIS package** to automate the ETL process.

- Configure **Data Flow tasks** to extract data from the source OLTP databases.

- Use **Staging Tables** for incremental loads and ensure proper error handling.

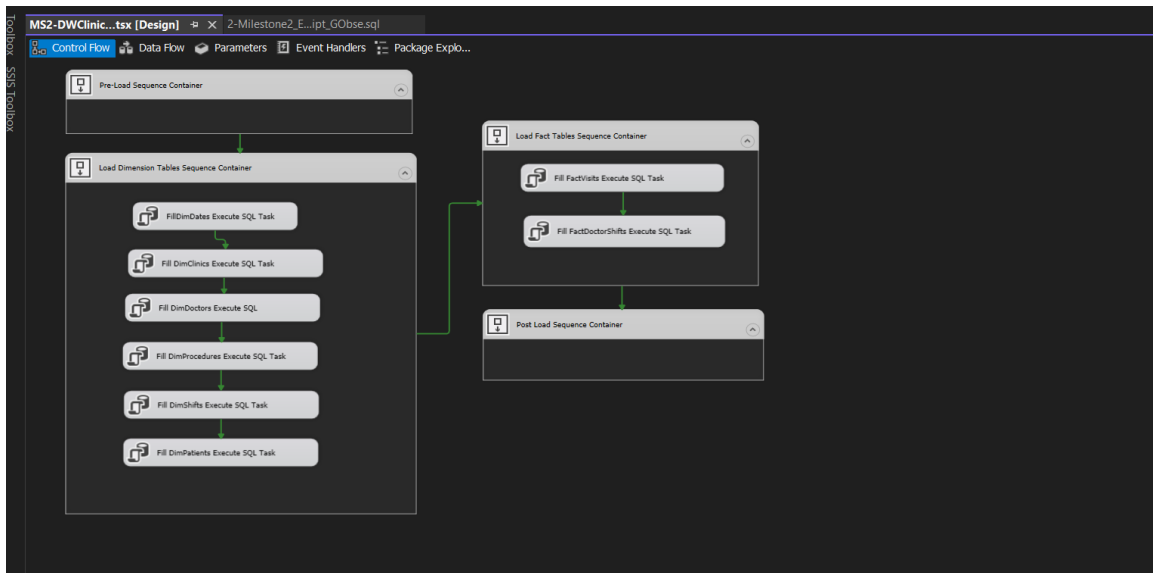- Implement **Logging and Notifications** to track process execution and failures.

**Figure13, showing ssis package of the datawarehouse**

- **Schedule the SSIS package for execution via **SQL Server Agent** or manual triggers.**



**Figure 14, showing the ssis scheduled job and job history.**

Figure15, demonstrating the data in DataClinicReport datawarehouse tables



Figure16, demonstrating the metadata worksheets

## ETL Objects

### Database ETL Objects

| Object Name | Type | Description | Location |
|---|---|---|---|
| pETLCreateStagingTables | Stored Procedure | Creates staging tables in a temporary database to store raw data before transformation and loi | temp database |
| pETLImportFileDataToStagingTables | Stored Procedure | Imports data from source files into staging tables in the temporary database for processing. | temp database |
| pETLTransformVistsData | Stored Procedure | Transforms and cleans visit-related data in the staging tables before loading | temp database |
| pETLSelectVistsData | Stored Procedure | Retrieves transformed visit data from the staging tables for validation or further processing. | temp database |
| pETLInsertVisitsData | Stored Procedure | Inserts transformed visit data from the staging tables into the visits table. | Patients |
| pETLDropForeignKeyConstraints | Stored Procedure | Drops the foreign keys before truncation | DWClinicReportDataGirum database |
| pETLTruncateTables | Stored Procedure | Truncates (empties) all relevant tables to prepare for new data loading | DWClinicReportDataGirum database |
| pETLFillDimDates | Stored Procedure | Populates the date dimension table with a range of dates used for analysis | DWClinicReportDataGirum database |
| pETLDimClinics | Stored Procedure | Loads or updates data for the clinics dimension, which stores clinic details | DWClinicReportDataGirum database |
| pETLDimDoctors | Stored Procedure | Loads or updates data for the doctors dimension, containing doctor-specific details | DWClinicReportDataGirum database |
| pETLDimProcedures | Stored Procedure | Loads or updates data for the procedures dimension, which lists medical procedures | DWClinicReportDataGirum database |
| pETLDimShifts | Stored Procedure | Loads or updates data for the shifts dimension, tracking doctor and clinic shifts. | DWClinicReportDataGirum database |
| pETLDimPatients | Stored Procedure | Loads or updates data for the patients dimension, storing patient details. | DWClinicReportDataGirum database |
| pETLFactVisits | Stored Procedure | Populates the fact table for patient visits, which records interactions between patients and clir | DWClinicReportDataGirum database |
| pFactDoctorShifts | Stored Procedure | Loads or updates the fact table that tracks doctor shifts. | DWClinicReportDataGirum database |
| pETLAddForeignKeyConstraints | Stored Procedure | Re-adds the foreign key constraints after the ETL process is complete to enforce data integrity | DWClinicReportDataGirum database |

### SSIS ETL Objects

| Object Name | Type | Description | Location |
|---|---|---|---|
| MS1-FilesToPatientsDB.dtsx | SSIS Package | Set of task that move data from files to staging tables | C:\_BISolutions\UWETLFinalGirumObse\BIETLFinalSSISPackagees\MS1-FilesToPatientsDB.dtsx |
| MS2-DWClinicReportDataETL.dstx | SSIS Package | Set of task that move data from staging tables to the datawarehouse | C:\_BISolutions\UWETLFinalGirumObse\BIETLFinalSSISPackagees\MS2-DWClinicReportDataETL.dtsx |
| MS3-ETLClincReportsDocumentData.dstx | SSIS Package | Execute task that move data from Datawarehouse view tables to Excell reports us | C:\_BISolutions\UWETLFinalGirumObse\BIETLFinalSSISPackagees\MS3-ETLClincReportsDocumentData.dtsx |

### Non-SQL ETL Objects

| Object Name | Type | Description | Location |
|---|---|---|---|
| ETL_Clinic_Reports_To_Excel.py | Python script | Script that moves CSV report data to an Excel spreadsheet | C:\_BISolutions\UWETLFinalGirumObse\PythonETL\PythonApplication1\ETL_Clinic_Reports_To_Excel.py |

Team Members | ETL Transformations | ETL Objects | +

## 6. Milestone 3: Non-SQL ETL

**Process Steps**

**Creating SQL Views**

- **Develop two SQL views in the DWClinicReportData database:**
    - **One for doctor shifts data.**
    - **One for patient visits data.**
- **Ensure the views contain relevant columns for analytical reporting.**
- **Validate the views by querying sample data.**

**Sql code for the views..**



```sql
1
2       If (OBJECT_ID('vRrtDoctorShifts') is not null) Drop View vRrtDoctorShifts;
3       go
4
5     CREATE OR ALTER VIEW vRrtDoctorShifts
6       AS
7       SELECT
8           CAST(CAST(d.FullDate AS DATE) AS VARCHAR(100)) AS ShiftDate,
9           c.ClinicName,
10          c.ClinicCity,
11          c.ClinicState,
12          s.ShiftID,
13          s.ShiftStart,
14          s.ShiftEnd,
15          doc.DoctorFullName,
16          ABS(fds.HoursWorked) AS HoursWorked
17      FROM FactDoctorShifts AS fds
18      JOIN DimDates AS d ON fds.ShiftDateKey = d.DateKey
19      JOIN DimClinics AS c ON fds.ClinicKey = c.ClinicKey
20      JOIN DimShifts AS s ON fds.ShiftKey = s.ShiftKey
21      JOIN DimDoctors AS doc ON fds.DoctorKey = doc.DoctorKey;
22
23
24      --select * from  vRrtDoctorShifts
```

```sql
1    If (OBJECT_ID('vRrtPatientVisits') is not null) Drop View vRrtPatientVisits;
2    go
3
4    CREATE OR ALTER VIEW vRrtPatientVisits AS
5      SELECT
6          CAST(CAST(d.FullDate AS DATE) AS VARCHAR(100)) AS VisitDate,
7          c.ClinicName,
8          c.ClinicCity,
9          c.ClinicState,
10         p.PatientFullName,
11         doc.DoctorFullName,
12         pr.ProcedureName,
13         ABS(fv.ProcedureVistCharge) AS ProcedureVistCharge
14     FROM FactVisits AS fv
15     JOIN DimDates AS d ON fv.DateKey = d.DateKey
16     JOIN DimClinics AS c ON fv.ClinicKey = c.ClinicKey
17     JOIN DimPatients AS p ON fv.PatientKey = p.PatientKey
18     JOIN DimDoctors AS doc ON fv.DoctorKey = doc.DoctorKey
19     JOIN DimProcedures AS pr ON fv.ProcedureKey = pr.ProcedureKey;
20
21     --select * from vRrtPatientVisits
22
23
```

**Extract, Transform, and Load (ETL) using Python**

- **Develop a Python script to perform ETL operations:**

  - **Extract data from the SQL views.**

  - **Transform the data as needed (e.g., date formatting, column renaming).**

  - **Load the data into an Excel spreadsheet named ClinicReportsData_.xlsx.**

- **Store the Excel file in the solution's Reports folder.**

- **Validate that the data is correctly loaded into the two worksheets (one per view).**

**Screernshot of the python script for the ETL for Clinic Reports**

```python
import os
import pyodbc
from datetime import datetime
from openpyxl import Workbook

# Database connection parameters
server = 'localhost\GMSSQLSERVER_DEV'  # Your SQL Server name
database = 'DWClinicReportDataGirum'  # Your database name

# Define the folder path where the Excel file will be saved
reports_folder = r"C:\__BISolutions\UWETLFinalGirumObse\Reports"

# Get the current date in the format YYYY-MM-DD
current_date = datetime.now().strftime("%Y-%m-%d")

# Construct the full file path for the Excel file
file_path = os.path.join(reports_folder, f"ClinicReportsData_{current_date}.xlsx")

# Establish the database connection using Windows authentication
conn = pyodbc.connect(f'DRIVER={{SQL Server}};SERVER={server};DATABASE={database};Trusted_Connection=yes;')
cursor = conn.cursor()

# Query to get Doctor Shift data from the vRrtDoctorShifts view
doctor_shifts_query = """
SELECT ShiftDate, ClinicName, ClinicCity, ClinicState, ShiftID, ShiftStart, ShiftEnd, DoctorFullName, HoursWorked
FROM vRrtDoctorShifts
"""

# Query to get Patient Visit data from the vRrtPatientVisits view
patient_visits_query = """
SELECT VisitDate, ClinicName, ClinicCity, ClinicState, PatientFullName, DoctorFullName, ProcedureName, ProcedureVistCharge
FROM vRrtPatientVisits
"""

# Execute the queries
cursor.execute(doctor_shifts_query)
doctor_shifts_data = cursor.fetchall()

cursor.execute(patient_visits_query)
patient_visits_data = cursor.fetchall()

# Create a new Excel workbook
workbook = Workbook()

# Create the first sheet for Doctor Shifts
sheet1 = workbook.active
sheet1.title = "DoctorShifts"

# Set column headers for DoctorShifts sheet
sheet1["A1"] = "ShiftDate"
sheet1["B1"] = "ClinicName"
sheet1["C1"] = "ClinicCity"
```
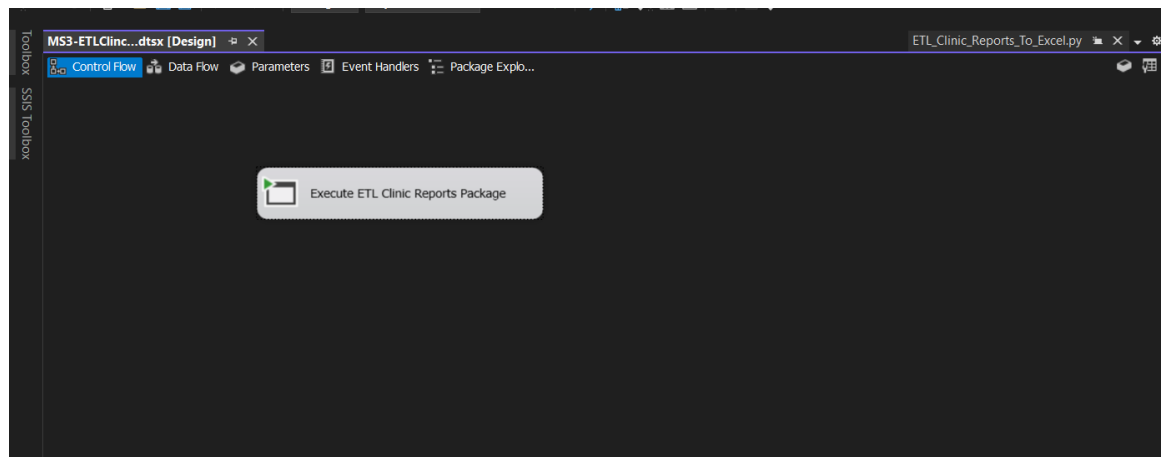
**SSIS Package Integration**

- **Develop an SSIS package (ETLClinicReportsDocumentData.dtsx):**

    o **Automate execution of the Python script.**

    o **Handle errors and log execution results.**

    o **Ensure proper folder and file management for storing output reports.**

**Updating ETL Process Metadata**

- **Update the Excel metadata spreadsheet:**

    o **Document source, destination, and transformation rules.**

    o **Ensure all tabs are updated as per the ETL workflow.**

- **Verify the metadata consistency with the implemented ETL process.**
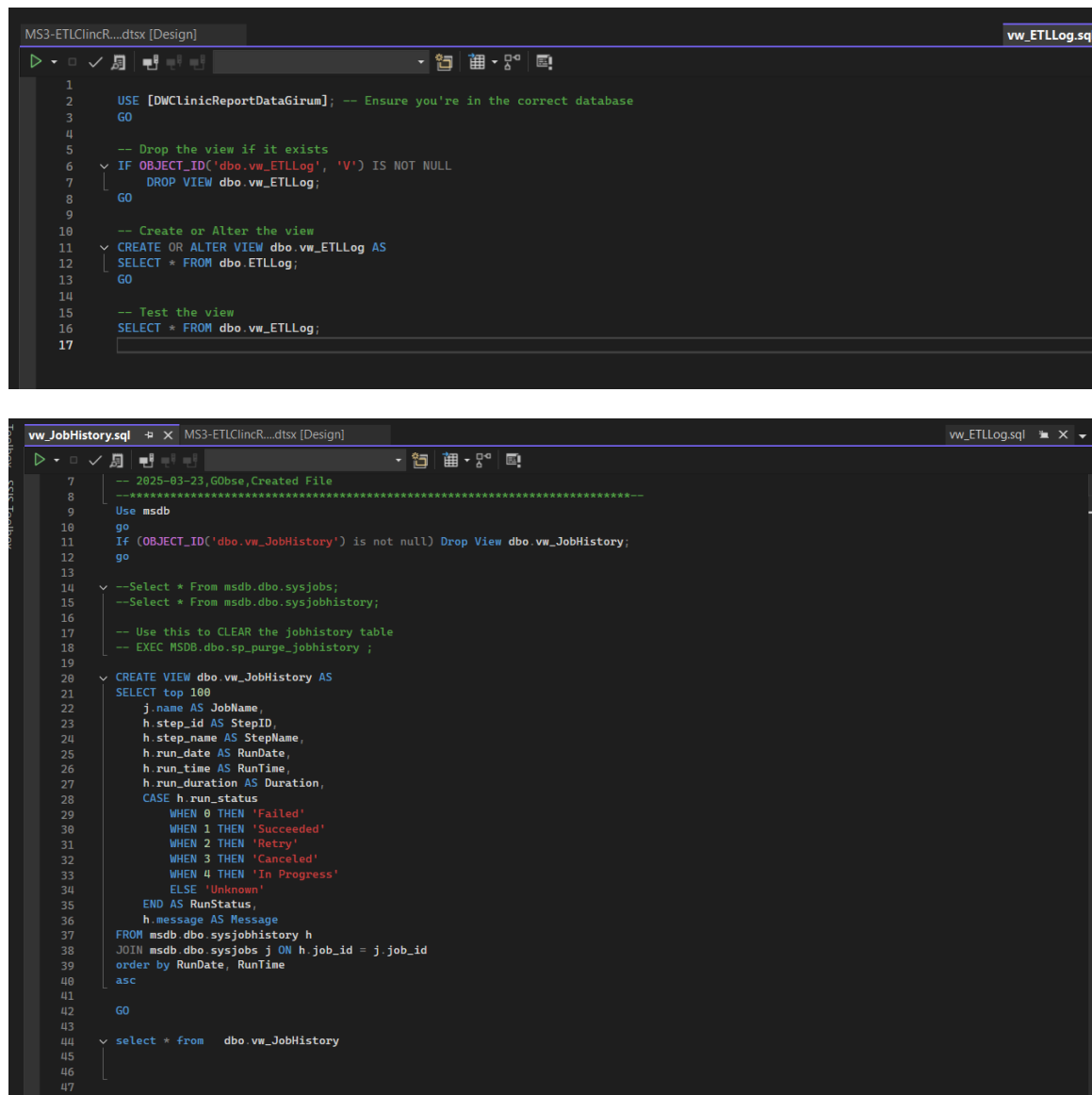
**Ssis package the executes the python ETL script**


**Excel worksheet of the clinic reports**

| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 1 | VisitDate | ClinicNam | ClinicCity | ClinicState | PatientFu | DoctorFul | Procedure | ProcedureVistCharge | |
| 2 | 2005-01- | Kirkland | Kirkland | WA | Thierry Bl | Ben Golds | Abdomina | 342 | |
| 3 | 2005-01- | Kirkland | Kirkland | WA | Joy Ceregl | Ben Golds | Skel xray- | 388 | |
| 4 | 2005-01- | Redmond | Redmond | WA | Janet Core | Lee Casts | Ther ultra | 438 | |
| 5 | 2005-01- | Redmond | Redmond | WA | Keith Jord | Ben Golds | Interview | 188 | |
| 6 | 2005-01- | Redmond | Redmond | WA | Christian ( | Ben Golds | Comprehe | 288 | |
| 7 | 2005-01- | Bellevue | Bellevue | WA | Margaret S | George Al | Skel xray- | 452 | |
| 8 | 2005-01- | Redmond | Redmond | WA | Sandeep E | Ben Golds | Dx ultraso | 388 | |
| 9 | 2005-01- | Kirkland | Kirkland | WA | Anton Des | Ben Golds | IVUS intra | 568 | |
| 10 | 2005-01- | Kirkland | Kirkland | WA | Yvonne Le | Ben Golds | Ther ult p | 435 | |
| 11 | 2005-01- | Bellevue | Bellevue | WA | Margaret I | Michael Jc | Contrast a | 238 | |
| 12 | 2005-01- | Redmond | Redmond | WA | Christoph | Ben Golds | Dx ultraso | 588 | |
| 13 | 2005-01- | Bellevue | Bellevue | WA | Michael H | Michael Jc | Other ske | 788 | |
| 14 | 2005-01- | Bellevue | Bellevue | WA | Thomas V | Michael Jc | not other | 488 | |
| 15 | 2005-01- | Kirkland | Kirkland | WA | Dominic C | Ben Golds | Skel xray- | 345 | |
| 16 | 2005-01- | Redmond | Redmond | WA | Jared Krar | Laura Gerr | Dx ultraso | 588 | |
| 17 | 2005-01- | Kirkland | Kirkland | WA | Matthias ( | Kevin Mit | Interview | 188 | |
| 18 | 2005-01- | Bellevue | Bellevue | WA | Kay Suffin | Kevin Mit | Pelvimetr | 488 | |
| 19 | 2005-01- | Kirkland | Kirkland | WA | Alan Haug | Kevin Mit | Contrast a | 238 | |
| 20 | 2005-01- | Kirkland | Kirkland | WA | Frank Gag | Kevin Mit | IVUS coro | 488 | |
| 21 | 2005-01- | Redmond | Redmond | WA | Samuel Ha | Laura Gerr | Dx ultraso | 348 | |
| 22 | 2005-01- | Redmond | Redmond | WA | Tomas All | Laura Gerr | IVUS intra | 568 | |
| 23 | 2005-01- | Kirkland | Kirkland | WA | Jacob Cha | Kevin Mit | Skel xray- | 388 | |
| 24 | 2005-01- | Redmond | Redmond | WA | Raquel Gia | Lee Casts | IVUS extra | 448 | |
| 25 | 2005-01- | Kirkland | Kirkland | WA | Grant Johr | Kevin Mit | Retroperit | 234 | |
| 26 | 2005-01- | Kirkland | Kirkland | WA | Julie Beck | Kevin Mit | Skel xray- | 388 | |
| 27 | 2005-01- | Kirkland | Kirkland | WA | Kay Lisboa | Kevin Mit | Skel xray- | 452 | |
| 28 | 2005-01- | Kirkland | Kirkland | WA | Dylan McC | Elizabeth | Brief inter | 188 | |
| 29 | 2005-01- | Bellevue | Bellevue | WA | Margaret ( | Kevin Mit | IVUS extra | 448 | |
| 30 | 2005-01- | Bellevue | Bellevue | WA | Stanley M | Elizabeth | Upper lim | 348 | |
| 31 | 2005-01- | Bellevue | Bellevue | WA | Vanessa C | Elizabeth | IVUS extra | 448 | |
| 32 | 2005-01- | Kirkland | Kirkland | WA | Nieves Kir | Elizabeth | Limited cc | 288 | |
| 33 | 2005-01- | Bellevue | Bellevue | WA | Rose Cave | Elizabeth | IVUS coro | 488 | |
| 34 | 2005-01- | Redmond | Redmond | WA | Clarence E | Laura Gerr | Skel xray- | 538 | |
| 35 | 2005-01- | Bellevue | Bellevue | WA | Betty Micl | Kevin Mit | Pelvimetr | 488 | |
| 36 | 2005-01- | Bellevue | Bellevue | WA | Keith Harι | Kevin Mit | IVUS perip | 488 | |
| 37 | 2005-01- | Kirkland | Kirkland | WA | Dylan Byh | Elizabeth | Skeletal s | 343 | |
| 38 | 2005-01- | Redmond | Redmond | WA | Gary Alexa | Laura Gerr | X-ray NEC | 188 | |
| 39 | 2005-01- | Kirkland | Kirkland | WA | Dominic C | Kevin Mit | Consultati | 288 | |

‹   ›        DoctorShifts    **PatientVisits**    +

## 7. Milestone 4: Automation, Reports, and Documentation

Process Steps

## Creating SQL Reporting Views

• Develop ETL report views to extract data from the ETL log and MSDB jobs tables.
• Ensure the views include relevant columns for tracking ETL execution status, errors, and processing times.
• Validate the views by running queries to confirm correct data extraction and filtering.
• Optimize performance by adding appropriate indexing where necessary.





Figure20, sql script for the views vw_ETL Log and vw_JobHistory

## Creating ETL Dashboard Reports with SSRS

Create an ETL report that shows the contents of the ETL log and MSDB jobs table.

- The report should display execution details, success/failure statuses, error messages, and timestamps.

- Use appropriate filters and sorting options to ensure the report provides useful insights.

- Ensure that your Data Sources and DataSets are properly configured and named according to their respective database objects.

Figure22, screenshot of the SSRS report for MSDB jobs and ETL Log Report

## Summery

The ETL Technical Manual outlines the implementation of an automated ETL solution for a small medical clinic, replacing the current manual CSV-based process with a structured system using SQL Server, SSIS, and Python.

The project involves automating data ingestion, transformation, and loading into a Data Warehouse, improving efficiency and accuracy. Key milestones include file-based ETL for CSV imports, data warehouse ETL for structured reporting, Python-based ETL for Excel-based analytics, and automation using SSIS.

The solution also integrates SQL views, stored procedures, and SSRS dashboards to track ETL performance and errors. Designed for scalability and reliability, this manual serves as a guide for data engineers, ETL developers, and  BI analysts involved in data integration projects.