

My First Attempt at Building a Data Warehouse

Contents:

- 1 Intro**
- 2 Project Resources & Plan Development**
- 3 Project Requirements**
- 4 Designing the Data Architecture**
 - a. Choosing the Data Management Approach
 - b. Define Architecture Layers
 - c. Build High-Level Architecture Diagram
- 5 Project Initialization**
 - a. Define Naming Conventions
 - b. Set up Git Repo
 - c. Create Schemas
- 6 Bronze Layer Build**
 - a. Analyzing: Source Systems
 - b. Coding: Data Ingestion
 - c. Validating: Data Completeness & Schema Checks
 - d. Document: Draw Data Flow Diagram
 - e. Commit Code to Git
- 7 Silver Layer Build**
 - a. Analyzing: Explore & Understand the Data
 - b. Document: High-Level Data Integration Diagram
 - c. Coding & Validating: Data Cleansing & Quality Checks
 - d. Document & Commit Code to Git
- 8 Gold Layer Build**
 - a. Analyzing: Explore Business Objects
 - b. Coding & Validating: Data Integration & Quality Checks
 - c. Document & Commit Code to Git
- 9 Exploratory SQL Data Analysis**
- 10 SQL Advanced Analytics**

1. INTRO

This project represents my journey in learning SQL and building a data warehouse in SQL Server through a guided project delivered by Baraa Khatib Salkini. The project was completed as part of Baraa's course '[The Complete SQL Bootcamp](#)'.

I strongly recommend the course for anyone seeking to learn SQL and the basics of data warehouse design.

Project Overview

I built a three-layer Medallion Architecture in SQL Server, ingesting and integrating data from fictitious CRM and ERP source systems and applying extensive data-quality checks, cleansing rules, and key-mapping logic.

In the Silver and Gold layers, I standardized and joined cross-system customer, product, and sales data, and designed a Star Schema with dimension and fact views to support analytics.

Throughout the project, I used SQL for modeling and transformation, Git for version control, and Draw.io for architectural and lineage documentation.

The final warehouse delivers a fully validated, analytics-ready dataset and demonstrates practical competency in SQL, data modeling, and data product development.

2. Project Resources & Plan Development

Project Resources:

Data

- o Cust_info.csv
- o Prd_info.csv
- o sales_details.csv
- o CUST_AZ12.csv
- o LOC_A101csv
- o PX_CAT_G1V2.csv

Data / Code Repository: [GitHub](#)

Project Management Tool: [Notion](#)

Diagram Development: [draw.io](#)

What I did to get started:

1. Obtained the project data from the instructor and saved locally on my machine
2. Created a Github account.
3. Downloaded and installed [draw.io](#) for desktop. [Draw.io](#) was used for the architecture diagram for this project.
4. Created a Notion account and set up high-level project plan with tasks for the project. Used the collection of Epics and tasks provided by the instructor, but made some tweaks as the project progressed..

Project Plan As Built in Notion

Aa. Name	Q. Progress	DWH Tasks
Requirements Analysis	100.00%	<ul style="list-style-type: none"> Analyze & Understand the Requirements
Design Data Architecture	100.00%	<ul style="list-style-type: none"> Choose Data Management Approach Brainstorm & Design the Layers Draw the Data Architecture (Draw.io)
Project Initialization	100.00%	<ul style="list-style-type: none"> Prepare the Git Repo & Prepare the Structure Create DB & Schemas Create Detailed Project Tasks (Notion) Define Project Naming Conventions
Build Bronze Layer	100.00%	<ul style="list-style-type: none"> Analyzing: Source Systems Coding: Data Ingestion Validating: Data Completeness & Schema Checks Document: Draw Data Flow (Draw.io) Commit Code in Git Repo
Build Silver Layer	100.00%	<ul style="list-style-type: none"> Analyzing: Explore & Understand Data Coding: Data Cleansing Validating: Data Correctness Checks Documenting & Versioning in GIT Commit Code in GIT Repo Document: Draw Data Integration (Draw.io)
Build Gold Layer	100.00%	<ul style="list-style-type: none"> Analyzing: Explore Business Objects Coding: Data Integration Validating: Data Integration Checks Document: Draw Data Model of Star Schema (Draw.io) Document: Create Data Catalog
Exploratory SQL Data Analysis	0.00%	<ul style="list-style-type: none"> Change Over Time Analysis Dimensions vs Measures Database Exploration Dimensions Exploration Date Exploration Measure Exploration Magnitude Analysis Ranking Analysis
SQL Advanced Analytics	0.00%	<ul style="list-style-type: none"> Change Over Time Analysis Cumulative Analysis Performance Analysis Part to Whole Analysis Data Segmentation Build Customers Report Build Products Report Document in GIT

3. Project Requirements

Objective

Develop a modern data warehouse using SQL Server to consolidate sales data, enabling analytical reporting and informed decision-making.

Specifications

- **Data Sources:** Import data from two source systems (ERP and CRM) provided as CSV files.
- **Data Quality:** Cleanse and resolve data quality issues prior to analysis.
- **Integration:** Combine both sources into a single, user-friendly data model designed for analytical queries.
- **Scope:** Focus on the latest dataset only; historization of data is not required.
- **Documentation:** Provide clear documentation of the data model to support both business stakeholders and analytics teams.

BI: Analytics & Reporting (Data Analysis)

Objective

Develop SQL-based analytics to deliver detailed insights into:

- Customer Behavior
- Product Performance
- Sales Trends

These insights empower stakeholders with key business metrics, enabling strategic decision-making.

4. Designing the Data Architecture

a. Choosing the Data Management Approach

In this task, I evaluated the requirements of the project and selected Data Warehouse as the appropriate approach from the following data architectures:

- Data Warehouse
- Data Lake
- Data Lakehouse
- Data Mesh

Next I learned about the types of data warehouse approaches, evaluated each against the requirements and selected the Medallion Architecture approach as it is most appropriate for this effort

- **Inmon:** Start with staging, next layer data is modeled in an EDW in 3NF, then broken into a third layer of data marts, then data BI tools are connected to the data marts
- **Kimbal:** Similar to Inmon, however, there is no EDW and data is moved from staging directly to the data marts as the EDW is time intensive. Kimbal is a faster approach than Inmon, but it can lead to chaos in the data mart with duplicate efforts
- **Data Vault:** Similar to Inmon however, the EDW layer is split into a Raw Vault (original data) and Business Vault (includes the business rules and transformations which prepare data for data marts)
- **Medallion Architecture:** Comprised of three layers (bronze, silver, and gold). Bronze Layer: Original data as-is. Enables traceability and issue identification; Silver: Transformed and cleansed data. Gold Layer: Similar to Data Mart in which business-ready objects are built and ready to share as products (machine learning, AI, reports) with business rules / logic applied.

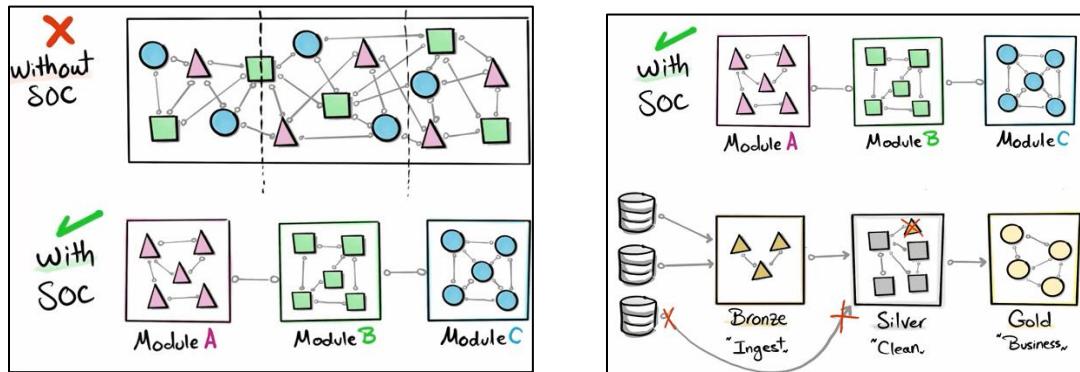
b. Define Architecture Layers

Summary:

Here I documented the definition and objective for each layer (Bronze, Silver, and Gold)

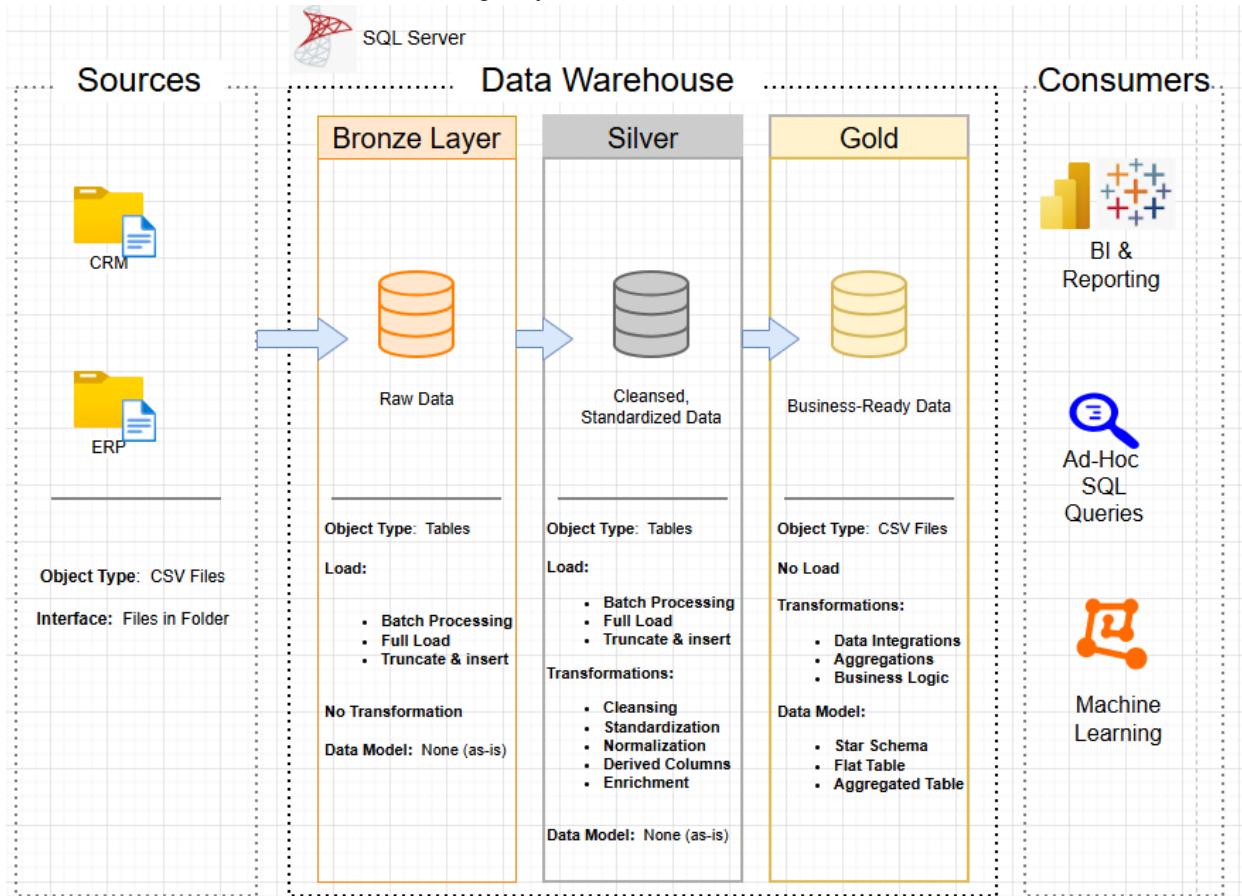
	Bronze	Silver	Gold Layer
Definition	Raw, unprocessed data as-is from sources	Cleansed and standardized data	Business-Ready Data
Objective	Traceability & Debugging	Intermediate Layer Prepare Data for Analysis	Provide data to be consumed for reporting & analytics
Object Type	Tables	Tables	Views
Load Method	Full Load (Truncate & Insert)	Full Load (Truncate & Insert)	None (not required as there are no tables)
Data Transformation	None (as-is)	<ul style="list-style-type: none">• Data Cleansing• Data Standardization• Data Normalization• Derived Columns• Data Enrichment	<ul style="list-style-type: none">• Data Integration• Data Aggregation• Business Logic & Rules
Data Modeling	None (as-is)	None (as-is)	<ul style="list-style-type: none">• Star Schema• Aggregated Objects• Flat Tables
Target Audience	Data Engineers	Data Analysts Data Engineers	Data Analysts Business Users

DOCUMENTING DETAIL IN THE TABLE ABOVE ENABLED SEPARATION OF CONCERNS (SOCs): Breaking the system into separate components with specific purposes with no overlap



c. Build High-Level Architecture Diagram

The objective here was to draw the data architecture using [Draw.io](#) according to the requirements, medallion architecture, and defined design layers within the architecture



5. Project Initialization

a. Define Naming Conventions

Here I defined the set of rules for naming everything in the project to ensure standardization so all members of the project are using the same naming conventions. This is done for:

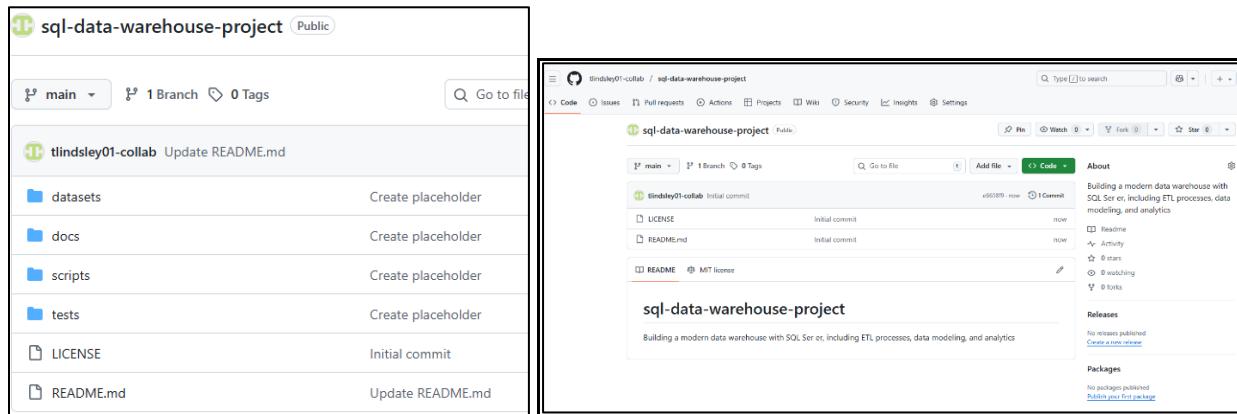
- The Database
- Schema
- Tables
- Stored Procedures, etc

b. Create Git Repo

For this task, I created a Github account, and learned how to create a repo structure for a project. A Git repo (repository) is a safe location for storing and tracking code used for a project and collaborating with the team. Additionally you can share your repository as part of your portfolio to show that you know how to build and maintain a well-documented git repository for your project.

I created my first repository in Git hub through this effort and gave it the name ‘sql-data-warehouse-project’. I then began creating folders and placeholders required for the project.

The project content structure in Git



c. Create Schemas

Below are the steps I took to create the database:

1. Open SQL Server Management Studio
2. Access the database master (system database) by creating a new query and entering: “USE MASTER”
3. Created the database with the following script in the same query window: “CREATE DATABASE DataWarehouse”;
4. I then went to the object explorer on the left hand side of SQL Server Management Studio, right clicked on ‘Databases’ and selected ‘Refresh’. The new database ‘DataWarehouse’ now appears in the collection of databases.

5. I now need to begin building the database. To do so, I entered 'USE DataWarehouse' in the same query window. The screenshot below shows the queries and outputs of steps 1-5

The screenshot shows the SQL Server Management Studio interface. The Object Explorer on the left shows a connection to 'PF33Z2Q8\SQLEXPRESS (SQL Server 16.0.1150 - ANSWERTHINK)'. Under 'Databases', a new database named 'DataWarehouse' is selected. The 'SQLQuery16...lindsay (61)*' query window contains the following SQL code:

```

1 --CREATE Database 'Data Warehouse'
2 USE MASTER; --Entered this to access the database master (system database) in SQL Server
3
4 CREATE DATABASE DataWarehouse; --Once in the system database, I entered and ran this to create my new database
5
6 USE DataWarehouse; --Once the database was created, I entered and ran this to begin using the database for build-out of the db
7

```

The 'Messages' pane at the bottom right shows 'Commands completed successfully.' and a completion time of '2025-11-28T09:29:9980095-08:00'.

6. Next I executed the following query to create the schema for each layer of the medallion architecture (Bronze, Silver, and Gold)

Schema Creation SQL

The screenshot shows the SQL Server Management Studio interface. The Object Explorer on the left shows a connection to 'PF33Z2Q8\SQLEXPRESS (SQL Server 16.0.1150 - ANSWERTHINK)'. Under 'Schemas', several schemas are listed: 'bronze', 'db_datareader', 'db_datawriter', 'db_ddladmin', 'db_denydatareader', 'db_denydatawriter', 'db_owner', 'db_securityadmin', 'dbo', 'gold', 'guest', 'INFORMATION_SCHEMA', 'silver', and 'sys'. The 'SQLQuery16...lindsay (61)*' query window contains the following SQL code:

```

1 --CREATE Database 'Data Warehouse'
2 USE MASTER; --Entered this to access the database master (system database) in SQL Server
3
4 CREATE DATABASE DataWarehouse; --Once in the system database, I entered and ran this to create my new database
5
6 USE DataWarehouse; --Once the database was created, I entered and ran this to begin using the database for build-out of the db
7
8
9 --Creates the schemas for each layer of the medallion architecture. They appear in
10 -- 'DataWarehouse' > 'Security' > 'Schemas'
11 CREATE SCHEMA bronze;
12 GO --This is the separator command that tells SQL to complete each command separately, then move to the next
13 CREATE SCHEMA silver;
14 GO --This is the separator command that tells SQL to complete each command separately, then move to the next
15 CREATE SCHEMA gold;
16 GO --This is the separator command that tells SQL to complete each command separately, then move to the next

```

The 'Messages' pane at the bottom right shows 'Commands completed successfully.' and a completion time of '2025-11-28T09:40:2450134-08:00'. A status bar at the bottom right indicates 'PF33Z2Q8\SQLEXPRESS (16.0 RTM) | ANSWERTHINK\tom.lindsay.'

7. Once Schemas were built for each layer, I then returned to Github to commit the code to my repo using the steps listed below. Completion of these steps also resulted in the completion of all tasks within the 'Project Initialization' EPIC in my project plan:

- a. Opened the 'scripts' folder in my sql-datawarehouse-project
- b. Selected 'Add file'
- c. Selected 'Create New File'
- d. Assigned file name 'init_database.sql'
- e. Pasted the code from SQL Server to document the steps taken to execute steps 1-6
- f. Added cautionary and informative commentary to the file to provide users of the script with an understanding of what the code will do and considerations before executing.
- g. Finally, I committed the code by selecting 'Commit Changes'

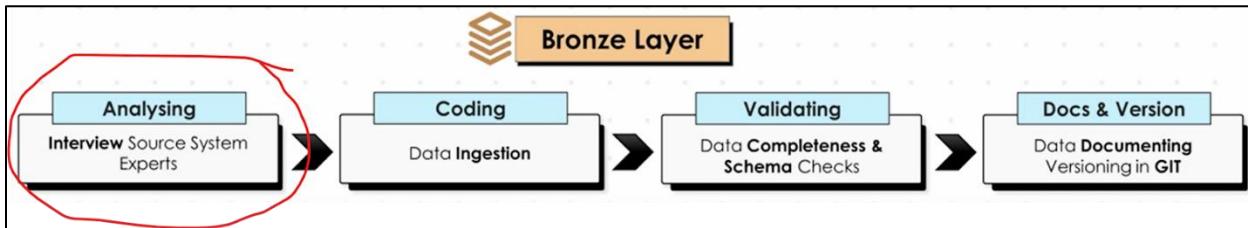
Code Commit in Git

The screenshot shows a GitHub repository interface. The top navigation bar includes 'Code', 'Issues', 'Pull requests', 'Actions', 'Projects', 'Wiki', 'Security', 'Insights', and 'Settings'. The main area is titled 'Code Commit in Git' and shows the file structure of 'sql-data-warehouse-project'. The 'Files' sidebar lists 'main', 'datasets', 'docs', 'scripts' (which contains 'init_database.sql' and 'placeholder'), 'tests', 'LICENSE', and 'README.md'. The right panel displays the contents of 'init_database.sql'. The code is as follows:

```
1  /*
2  =====
3  Create Database and Schemas
4  =====
5  Script Purpose:
6
7  This script creates a new database named 'DataWarehouse' after checking if it already exists.
8  If the database exists, it is dropped and recreated. Additionally, the script sets up three schemas
9  within the database: 'bronze', 'silver', and 'gold'.
10
11 --WARNING:
12 --Running this script will drop the entire 'DataWarehouse' database if it exists.
13 --All data in the database will be permanently deleted. Proceed with caution
14 --and ensure you have proper backups before running this script.
15 */
16 USE MASTER; --Entered this to access the database master (system database) in SQL Server
17 GO
18
19 --ONLY USE THE FOLLOWING IF RECREATING THE DATABASE. IF NOT...CONTINUE TO 'CREATE DATABASE SECTION BELOW.
20 --This will drop and recreate the 'DataWarehouse' database
21 IF EXISTS (SELECT 1 FROM sys.databases WHERE name = 'DataWarehouse')
22 BEGIN
23     ALTER DATABASE DataWarehouse SET SINGLE_USER WITH ROLLBACK IMMEDIATE;
24     DROP DATABASE DataWarehouse;
25 END;
26 GO
```

6. Bronze Layer Build

a. Analyzing Source Systems



The first step in building the layer is to understand the source system. To do so in a normal project, I would conduct interviews with the source system owners / SMEs to understand the nature of the system to be connected to the warehouse. As this is a fictional project with fictional systems....no interviews were conducted

Conducting the interview allows you to understand:

- Business Context & Ownership
- Who owns the data
- What business processes it supports
- System and data documentation
- Data model and data catalog for the source system

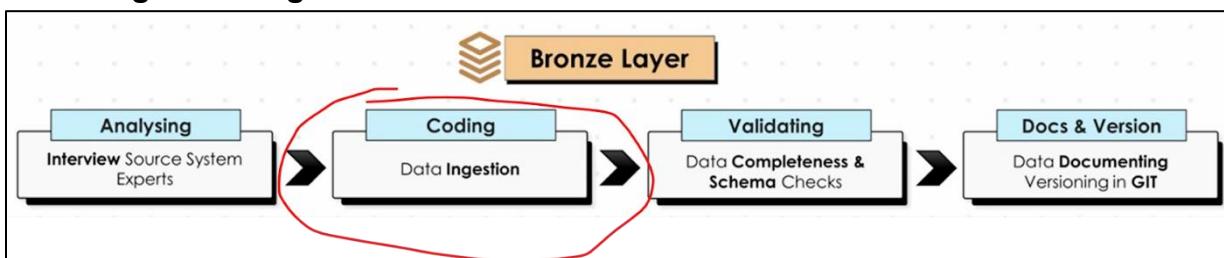
Architecture & Technology Stack

- How is the data stored? (On prem in SQL Server or Oracle. In the cloud on Azure or AWS)
- What are the integration capabilities of the source system? Is it through APIs, Kafka, File Extract, Direct DB connection?)

Extract & Load

- Can we do incremental loads vs full loads?
- Data scope & historical needs?
- Expected extract size (gigabytes / terabytes). Important for understanding whether we have the right tools and platforms to connect and extract source system data.
- Are there data volume limitations of the source system
- How do we avoid impacting the source system's performance in the connection and extraction
- What are the authentication and authorization requirements of the source system that need to be adhered to / made available to access and extract the data (tokens, SSH keys, VPN, IP whitelisting....)

b. Coding: Data Ingestion



The first step in this task is to understand the meta data, the structure, the schema of the incoming data before defining the structure of the tables for the Bronze layer in DDL (data definition language). This can be accomplished through the interviews with the source system SMEs / owners **and** through data

profiling / exploring the data on your own to identify column names and data types. If exploring on your own first, ensure that your assumptions are validated by the system owners / SMEs before building the bronze layer.

The following steps were used in this project to gain this understanding of the data sets provided for this project.

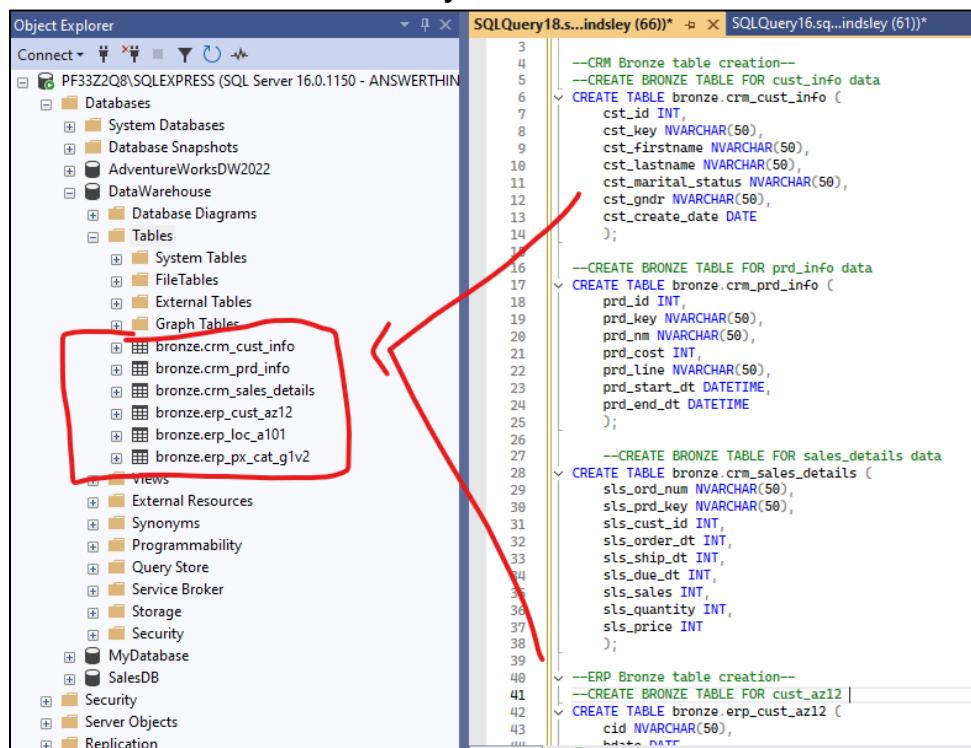
1. Open the 'cust_info.csv' data file for CRM and identify the column names, review the data and determine the data types
2. Open a new query in 'DataWarehouse' database and enter the following command to create the table in the bronze layer of the database (adhering to naming conventions defined previously in this project):

```
CREATE TABLE bronze.crm_cust_info
```
3. In the same query, create the following columns and assign their data types based on what is understood from the provided data set. The names should be identical to how they are defined in the data set

```
cst_id INT,
cst_key NVARCHAR(50),
cst_firstname NVARCHAR(50),
cst_lastname NVARCHAR(50),
cst_marital_status NVARCHAR(50),
cst_gndr NVARCHAR(50),
cst_create_date DATE
```

4. Execute the query, then refresh the tables in the database to confirm creation.
5. Execute steps 1 - 3 for all remaining files for both source systems

Bronze Layer Table Creation



The screenshot shows the SQL Server Management Studio interface. On the left, the Object Explorer pane displays a tree structure of databases and objects. A red box highlights the 'Tables' node under the 'bronze' database. On the right, the SQL Query Editor pane contains several 'CREATE TABLE' statements for different bronze tables:

```

3  --CRM Bronze table creation--
4  --CREATE BRONZE TABLE FOR cust_info data
5
6  CREATE TABLE bronze.crm_cust_info (
7      cst_id INT,
8      cst_key NVARCHAR(50),
9      cst_firstname NVARCHAR(50),
10     cst_lastname NVARCHAR(50),
11     cst_marital_status NVARCHAR(50),
12     cst_gndr NVARCHAR(50),
13     cst_create_date DATE
14 );
15
16  --CREATE BRONZE TABLE FOR prd_info data
17  CREATE TABLE bronze.crm_prd_info (
18      prd_id INT,
19      prd_key NVARCHAR(50),
20      prd_nm NVARCHAR(50),
21      prd_cost INT,
22      prd_line NVARCHAR(50),
23      prd_start_dt DATETIME,
24      prd_end_dt DATETIME
25 );
26
27  --CREATE BRONZE TABLE FOR sales_details data
28  CREATE TABLE bronze.crm_sales_details (
29      sls_ord_num NVARCHAR(50),
30      sls_prd_key NVARCHAR(50),
31      sls_cust_id INT,
32      sls_order_dt INT,
33      sls_ship_dt INT,
34      sls_due_dt INT,
35      sls_sales INT,
36      sls_quantity INT,
37      sls_price INT
38 );
39
40  --ERP Bronze table creation--
41  --CREATE BRONZE TABLE FOR cust_az12 |
42  CREATE TABLE bronze.erpcust_az12 (
43      cid NVARCHAR(50),
44      bdate DATE
45 );

```

6. Loading the data: We are doing a bulk insert directly from files into the database. A single operation in one go. I wrote the script to load the cust_info data file into the bronze.crm_cust_info' table.

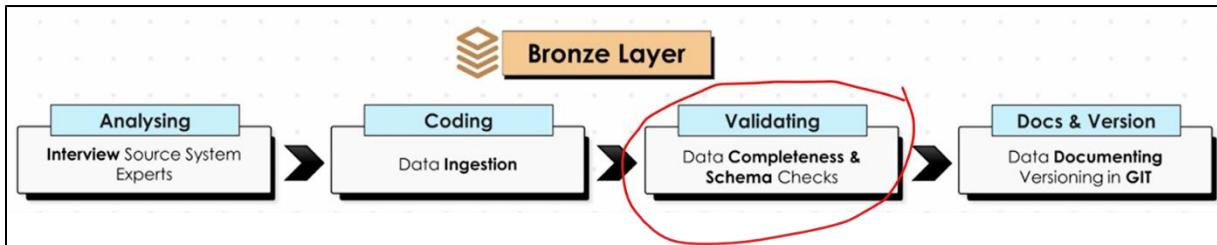
NOTE: I learned that this script will encounter an error if the file is open by a user. I had the file open when I attempted to execute the script and encountered just that. Closed the file and executed again successfully

--Script to load the cust_info data file into the bronze.crm_cust_info' table

```
BULK INSERT bronze.crm_cust_info
FROM 'C:\sql-data-warehouse-project\datasets\source_crm\cust_info.csv'
WITH (
    FIRSTROW = 2, --Tells SQL Server that the data starts on row 2 in the file
    FIELDTERMINATOR = ',', -- Tells SQL Server that commas are used as the separator between data for each column
    TABLOCK --Option to lock the table to improve performance while data is being loaded
);
```

c. Validating: Data Completeness & Schema Checks

This was essentially a continuation of the previous tasks conducted in 'Coding'



a. The next step was to run a query to view the data in the table to confirm that the data was pulled into the correct columns. I ran the following: SELECT * FROM bronze.crm_cust_info

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date
1	AW00011000	Jon	Yang	M	M	2025-10-06
2	AW00011001	Eugene	Huang	S	M	2025-10-06
3	AW00011002	Ruben	Tores	M	M	2025-10-06
4	AW00011003	Christy	Zhu	S	F	2025-10-06
5	AW00011004	Elizabeth	Johnson	S	F	2025-10-06
6	AW00011005	Julio	Ruiz	S	M	2025-10-06
7	AW00011006	Janet	Alvarez	S	F	2025-10-06
8	AW00011007	Marco	Mehta	M	M	2025-10-06
9	AW00011008	Rob	Verhoff	S	F	2025-10-06
10	AW00011009	Shannon	Carlson	S	M	2025-10-06
11	AW00011010	Jacquelyn	Suarez	S	F	2025-10-06
12	AW00011011	Curtis	Lu	M	M	2025-10-06
13	AW00011012	Lauren	Walker	M	F	2025-10-06
14	AW00011013	Ian	Jankine	M	M	2025-10-06

b. The previous steps were only a portion of the full code that needed to be written. The command needed to include a 'TRUNCATE' command to enable bulk load if the data is refreshed in the source file

without adding to the existing data in the table. Truncating will replace existing data with the new / refreshed data source. So....the command was rewritten and run again

```
--Script to load the cust_info data file into the bronze.crm_cust_info' table

TRUNCATE TABLE bronze.crm_cust_info; --Empties existing data in the table, then loads with refreshed data from scratch.
--This is a full load.

BULK INSERT bronze.crm_cust_info
FROM 'C:\Users\tom.lindsay\OneDrive - The Hackett Group, Inc\Desktop\Data Knowledge Development\SQL Bootcamp Files\sql-data-warehouse-project\datasets\source_crm\cust_info.csv'
WITH (
    FIRSTROW = 2, --Tells SQL Server that the data starts on row 2 in the file
    FIELDTERMINATOR = ',', -- Tells SQL Server that commas are used as the separator between data for each column
    TABLOCK --Option to lock the table to improve performance while data is being loaded
);

-- QUERY TO CHECK THE QUALITY OF THE DATA LOAD AND ENSURE DATA WENT TO THE CORRECT COLUMNS
SELECT * FROM bronze.crm_cust_info -- QUERY TO CHECK THE QUALITY OF THE DATA LOAD AND ENSURE DATA WENT TO THE
CORRECT COLUMNS

SELECT COUNT(*) FROM bronze.crm_cust_info ---- QUERY TO CHECK THAT THE TABLE ROW COUNT MATCHES THE CSV FILE
```

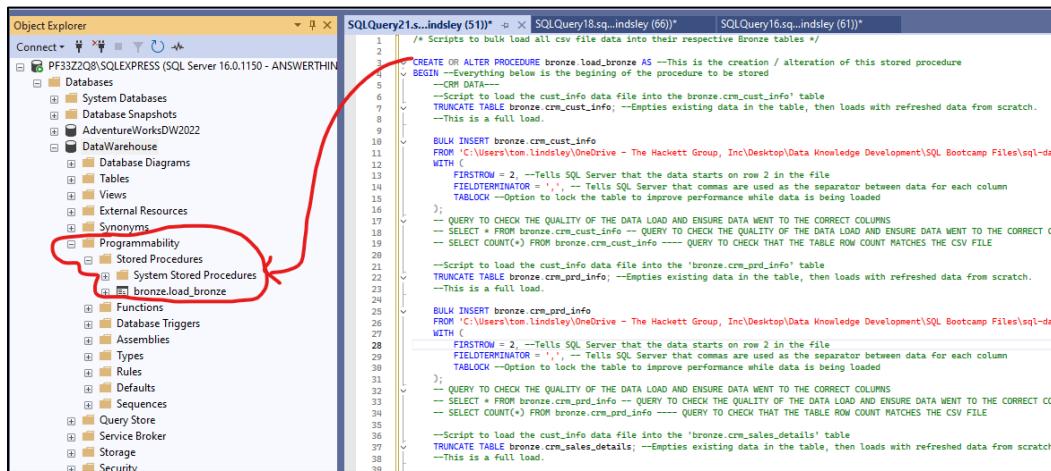
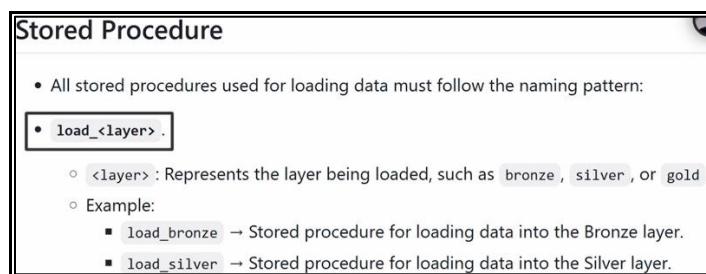
c. The same script was written and run for each data file for CRM and ERP data tables.

d. After running each script, I checked the results of the load against the source files to confirm that the data was placed into the correct columns and that the tables and the data sources had the same number of rows.

e. With the checks complete, I was good to move on to creating a stored procedure as this script would need to be run on a recurring basis.

f. Creating a stored procedure:

- I wrote a stored procedure for the loading of the data into the database as it is intended to be a recurring action to bring in the data.



- Once executing the stored procedure and refreshing the DB to visually confirm its existence, I then created a new query to test the execution of the procedure. The following SQL statement was used: EXEC bronze.load_bronze
- The message output in SQL Server Studio is fairly vague. I edited the script in the procedure to enhance the information produced in the message output for greater clarity / confirmation of successful load to the intended tables. The result of the enhancement:

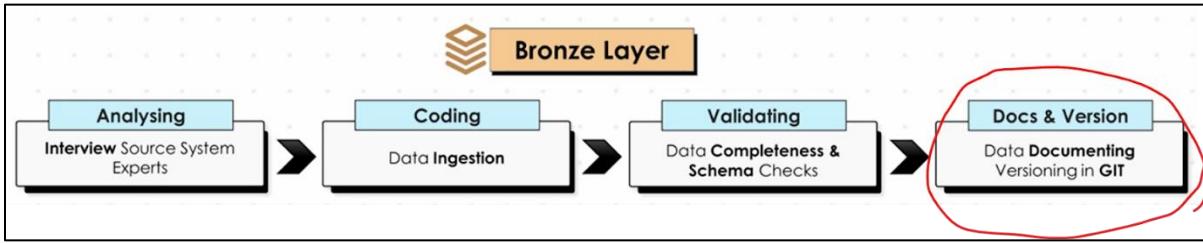
```

Messages
=====
Loading Bronze Layer
-----
Loading CRM Tables
-----
>>Truncating Table bronze.crm_cust_info
>>Inserting Data Into: bronze.crm_cust_info
(18493 rows affected)
>>Truncating Table bronze.crm_prd_info
>>Inserting Data Into: bronze.crm_prd_info
(397 rows affected)
>>Truncating Table bronze.crm_sales_details
>>Inserting Data Into: crm_sales_details
(60398 rows affected)
-----
Loading ERP Tables
-----
>>Truncating Table bronze.erp_cust_az12
>>Inserting Data Into: erp_cust_az12
(18483 rows affected)
>>Truncating Table bronze.erp_loc_a101
>>Inserting Data Into: erp_loc_a101
(18484 rows affected)
>>Truncating Table bronze.erp_px_cat_g1v2
>>Inserting Data Into: erp_px_cat_g1v2

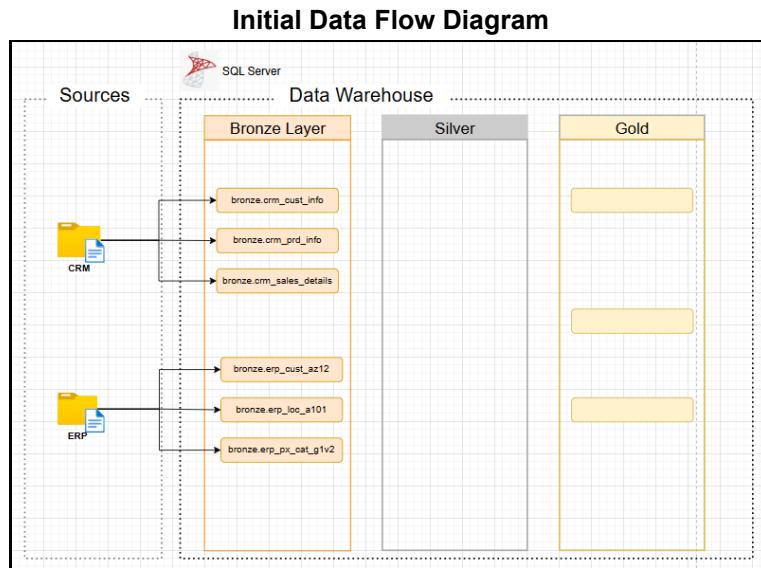
```

- The next step was to further enhance the stored procedure by adding a TRY CATCH command to handle errors that might occur in execution
- I also added script to include the duration of the loading effort for each table. To do this I declared two variables just before the 'BEGIN TRY' command for start time and end time. These variables will be used to output the execution duration that will be included in the message detail

d. Document: Draw Data Flow (Draw.io)

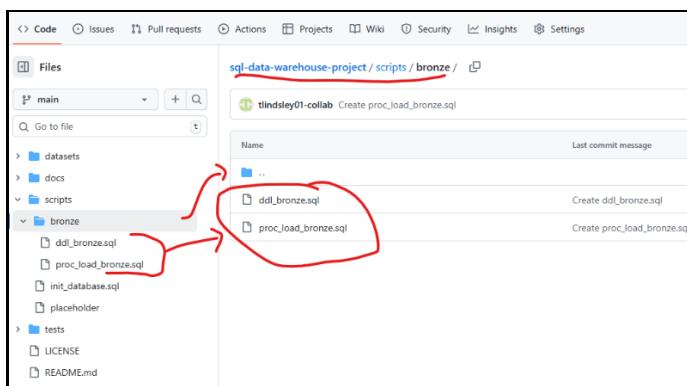


1. For the first step in this task, I drew a simple diagram in [Draw.io](#) of the flow of data from sources to the tables in the bronze layer to document data lineage. The output completed for the bronze layer is depicted below.

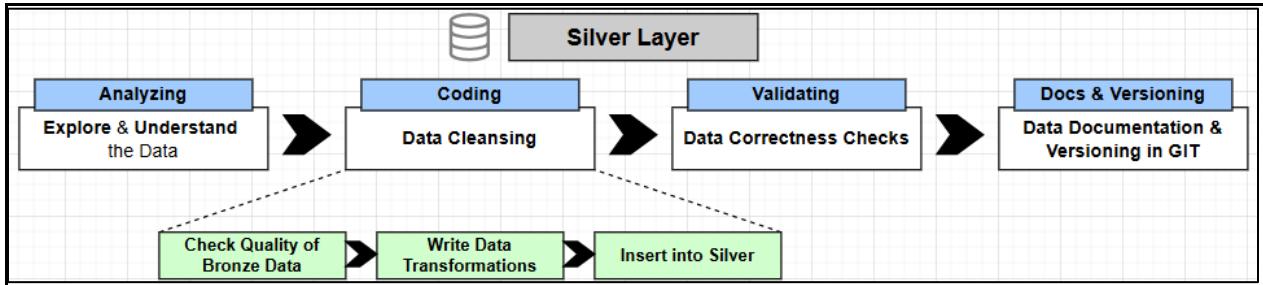


d. Commit Code to Git

The last task in this Epic is to commit my code into my Git repository. All of the code generated for table creation, and table loading for the bronze layer were added as separate files in the repository in a folder specific to the bronze layer.



7. Silver Layer Build

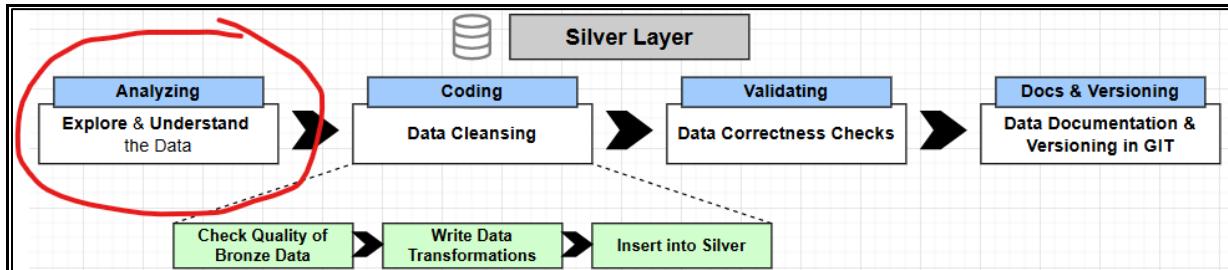


It is now time to build the Silver layer of the data warehouse now that the bronze layer has been built, documented, and its supporting code has committed to my Git repo.

This Epic is comprised of the following key tasks that will be completed in the order in which they are displayed below:

- Analyzing: Explore & Understand Data
- Coding: Data Cleansing
 - Check Quality of Bronze Data
 - Write Data Transformations
 - Insert into Silver Layer
- Validating: Data Correctness Checks
- Documenting & Versioning in GIT
- Commit Code in GIT Repo

a. Analyzing: Explore & Understand Data



Exploring and understanding the data is key to building out the Silver layer, the relationships, etc. To do this I explored each table one-by-one in the bronze layer.

For each table, I selected the top 1000 rows to get an initial understanding. This can be done by creating and executing the query below, or (in SQL Server) by right clicking on the desired table and selecting 'Select Top 1000 Rows'.

A key step in the data exploration task is to build an integration model to document what I observe / assume about how the tables are related. This is important so that I would not later forget what I observed

CRM Data Tables

- **bronze.crm_cust_info:** This table contains customer info with standard customer detail such as first name, last name, customer creation date. It includes a customer id (cust_id) as the primary key as well as a customer key (cust_key) which seems duplicative. An initial view of this table indicates that both the cust_id and the cust_key could be used to define relationships with other tables.
- **bronze.crm_prd_info:** This table contains basic product information along with historical information (historical product pricing). It includes a product id (prd_id) that appears to serve as the primary key, but it also includes a product key (prd_key). Like the cust_info table, an initial view of this table indicates that two columns (prd_id and prd_key) could be used to define relationships with other tables. I will have to determine which is most appropriate as I review the other tables.
- **bronze.crm_sales_details:** This table contains historical transactional records about sales and orders. Each record includes an order number (sls_ord_num), a product key for the product in the order (sls_prd_key), the id of the customer that placed the order (sls_cust_id), as well as order date, ship date, total order sales, sales quantity, and sales price. It appears that the sls_ord_num is the primary key for this table.

It appears that the sls_prd_key can be used to define the relationship with the crm_prd_info table but the value in the sls_prd_key for each record appears to only be a subset of the value listed for each record in the prd_key field of the crm_prd_info table.

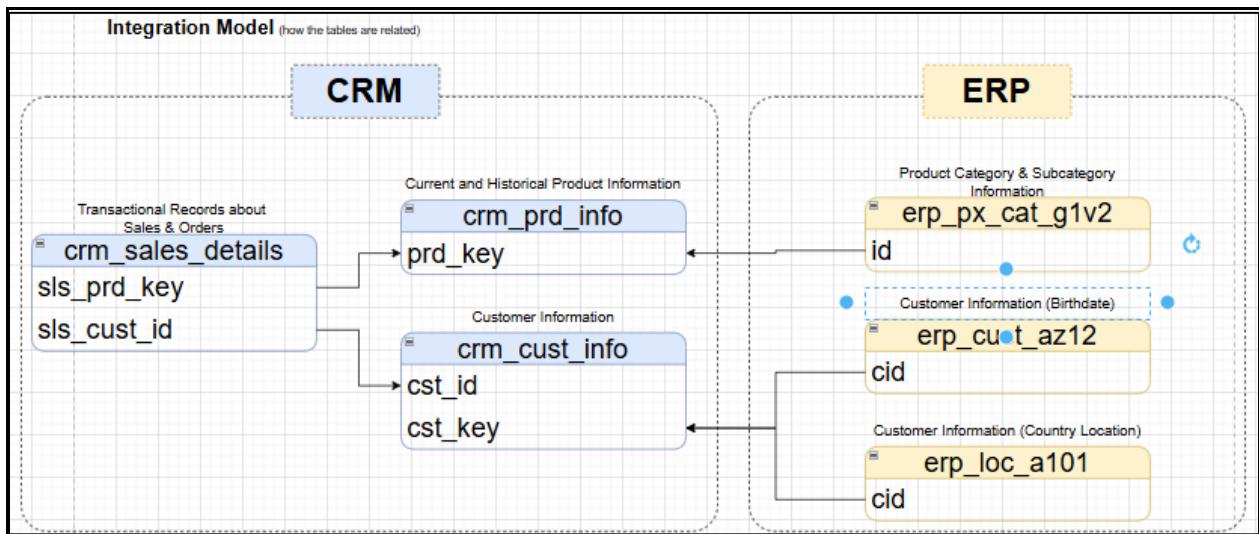
The data type and character count for records in the sls_cust_id field appear to be a match for the data type and character count for the cust_id field in the crm_cust_info table. Therefore I will use the sls_cust_id field to define the relationship with the crm_cust_info table matching to the cst_id field.

ERP Data Tables

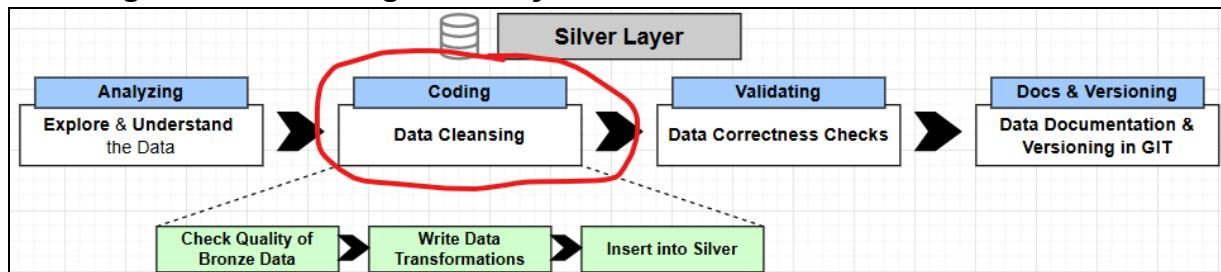
- **bronze.erp_cust_az12:** This table contains customer birth dates and gender. Its primary key appears to be the customer id field (cid). The values in this field for each record appear to be partial match to the customer key (cst_key) in the CRM table crm_cust_info, but it includes additional leading characters. No other field in this table could be used to create a relationship with any other table. So the customer id field in this table will be used to create a relationship with the crm_cust_info table by matching on its cst_key
- **bronze.erp_loc_a101:** This table contains country for each customer. Similar to the erp_cust_az12 table, its primary key appears to be the customer id field (cid) with partial match to the customer key (cst_key) in the CRM table crm_cust_info, but it includes “-” as separators for characters in this field. No other field in this table could be used to create a relationship with any other table. So the customer id (cid) field in this table will be used to create a relationship with the crm_cust_info table by matching on its cst_key
- **bronze.erp_px_cat_g1v2:** This table contains product categories, subcategories, and an identifier for “maintenance”. Not sure what the “maintenance” field is used for, but it only includes “Yes” or “No” as values for each record. The primary key appears to be the “id” field and is comprised of five (5) characters →Example: ‘AC_BR’. I assume this is the “id” for each product subcategory. The values in this field appear to be a match for the first five (5) characters in the prd_key field in the crm_prd_info table. This allows me to define a relationship between the two tables by matching on these fields. Doing so will require some transformation / manipulation first.

b. Document: High-Level Data Integration Diagram

As I explored the data in each table, I built the model below to help me visually understand the relationships that I will need to define.



c. Coding: Data Cleansing & Quality Checks



The first step in this task is to build the structure of the silver layer. It will be identical to the structure used for the Bronze layer. To accomplish this, I open the ddl script used for the bronze layer that I have saved in Git, paste it into a new query in SQL Server, then replace all references to 'bronze' with 'silver'.

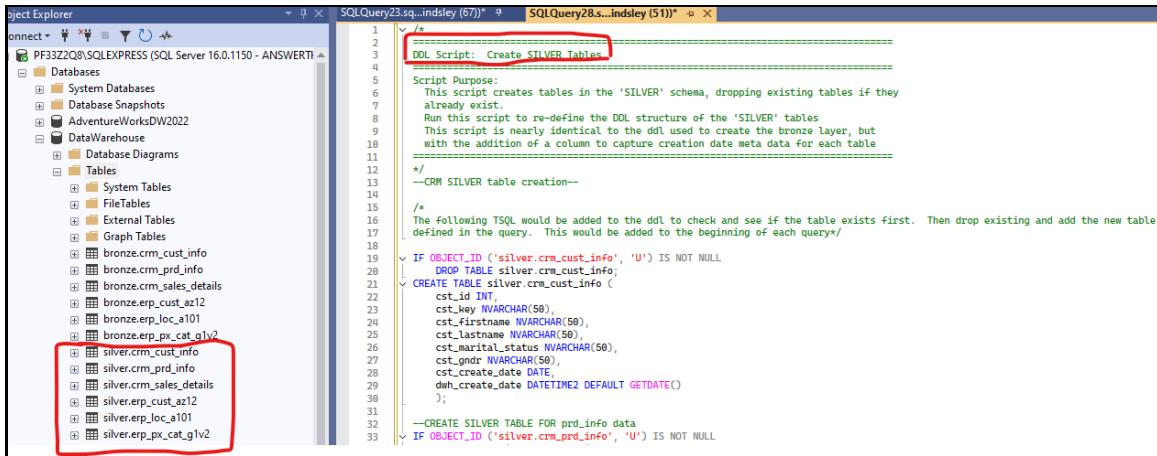
In addition to replicating the ddl to build the Silver layer, I also have to add metadata columns to each table to capture the following data. This information is needed to enable tracking and root cause analysis of corrupt data errors, data gaps in loads, etc:

- **create_date:** The record's load timestamp. Added "dwh_create_date DATETIME2 DEFAULT GETDATE()" in each CREATE TABLE command as an additional column for every table

These columns were not created in this step for each table. Will need to determine if these will be added later.

- **update_date:** The records last update timestamp
- **source_system:** The origin system of the record
- **file_location:** The file source of the record (where it is stored and made available for loading into the db)

Creation and execution of the ddl script to build the silver layer successfully complete resulted in the silver layer tables being added to the Data Warehouse as depicted below:



```

Object Explorer
PF3Z2Q\SQLEXPRESS (SQL Server 16.0.1150 - ANSWERTI)
Databases
Tables
System Tables
FileTables
External Tables
Graph Tables
bronze.crm_cust_info
bronze.crm_prd_info
bronze.crm_sales_details
bronze.emp_cust_az12
bronze.emp_loc_a101
bronze.emp_px_cat_g1v2
silver.crm_cust_info
silver.crm_prd_info
silver.crm_sales_details
silver.emp_cust_az12
silver.emp_loc_a101
silver.emp_px_cat_g1v2

SQLQuery23.sql...indsley (67)* SQLQuery28.s...indsley (51)*
=====
DDL Script: Create SILVER Tables
=====
Script Purpose:
This script creates tables in the 'SILVER' schema, dropping existing tables if they already exist.
Run this script to re-define the DDL structure of the 'SILVER' tables
This script is nearly identical to the ddl used to create the bronze layer, but with the addition of a column to capture creation date meta data for each table
=====
/*
--CRM SILVER table creation

The following TSQL would be added to the ddl to check and see if the table exists first. Then drop existing and add the new table defined in the query. This would be added to the beginning of each query*/
IF OBJECT_ID ('silver.crm_cust_info', 'U') IS NOT NULL
DROP TABLE silver.crm_cust_info;
CREATE TABLE silver.crm_cust_info (
    cst_id INT,
    cst_key NVARCHAR(50),
    cst_firstname NVARCHAR(50),
    cst_lastname NVARCHAR(50),
    cst_marital_status NVARCHAR(50),
    cst_gndr NVARCHAR(50),
    cst_create_date DATE,
    dim_create_date DATETIME2 DEFAULT GETDATE()
);

--CREATE SILVER TABLE FOR prd_info data
IF OBJECT_ID ('silver.crm_prd_info', 'U') IS NOT NULL

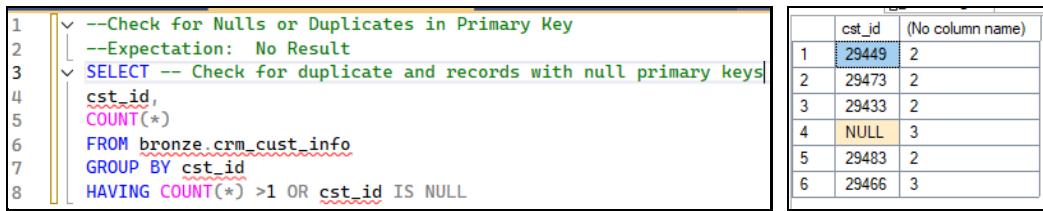
```

With the Silver layer tables built, the next step is to identify the data quality issues in each data set. This must be done before writing any transformation scripts. This exploration is done by evaluating each table in the bronze layer.

CRM → Customer Info Table

Checking for and Correcting Records with Duplicate and Null Primary Keys

Primary Key Duplicate and Null For each table, the primary key was evaluated to determine whether or not there were NULL values or duplicates using the following script. The result was the identification of 6 customer ids with two or more duplicates. Three (3) of those records had null values in the cst_id field (no customer id)



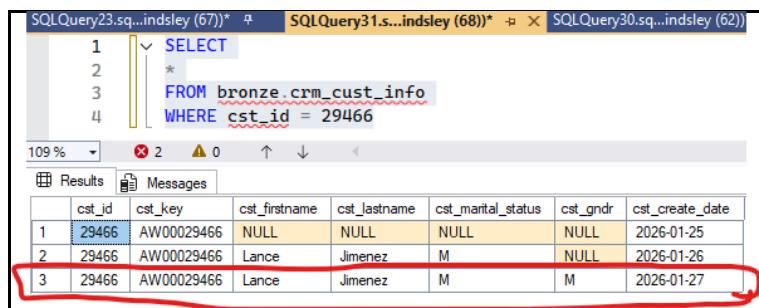
```

1 --Check for Nulls or Duplicates in Primary Key
2 ---Expectation: No Result
3 SELECT -- Check for duplicate and records with null primary keys
4 cst_id,
5 COUNT(*)
6 FROM bronze.crm_cust_info
7 GROUP BY cst_id
8 HAVING COUNT(*) >1 OR cst_id IS NULL

```

cst_id	(No column name)
29449	2
29473	2
29433	2
NULL	3
29483	2
29466	3

Next I needed to create a new query to clean the data to correct / account for the duplicates and null values. Before starting the transformation query, I needed to explore the data a bit more. I started by focusing on one of the IDs (29466) that appears more than once and running a query to view the associated records. This review allowed me to see that cst_id 29466 appears for three different records in the data set. Each record has a different creation date. The most recent creation date was selected as the primary record in the dataset.



```

SELECT *
FROM bronze.crm_cust_info
WHERE cst_id = 29466

```

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date
1 29466	AW00029466	NULL	NULL	NULL	NULL	2026-01-25
2 29466	AW00029466	Lance	Jimenez	M	NULL	2026-01-26
3 29466	AW00029466	Lance	Jimenez	M	M	2026-01-27

I assumed that all other duplicates likely had different creation dates. So, to set the most recently added records for duplicate cst_ids, I needed to rank them according to their creation date and assign a ranking number. The following query was used to rank all records in order by their creation date with the assignment of "1" for the most recent. Records with the assignment of "1" will be used as the primary records from the data set.

```

1  SELECT
2  *
3  ROW_NUMBER() OVER (PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_last
4  FROM bronze.crm_cust_info
    
```

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date	flag_last
1	NULL	PO25	NULL	NULL	NULL	NULL	1
2	NULL	SF566	NULL	NULL	NULL	NULL	2
3	NULL	13451235	NULL	NULL	NULL	NULL	3
4	11000	AW00011000	Jon	Yang	M	2025-10-06	1
5	11001	AW00011001	Eugene	Huang	S	2025-10-06	1
6	11002	AW00011002	Ruben	Torres	M	2025-10-06	1
7	11003	AW00011003	Christy	Zhu	S	2025-10-06	1
8	11004	AW00011004	Elizabeth	Johnson	S	2025-10-06	1
9	11005	AW00011005	Julio	Ruiz	S	2025-10-06	1
10	11006	AW00011006	Janet	Alvarez	S	2025-10-06	1
11	11007	AW00011007	Marco	Mehta	M	2025-10-06	1
12	11008	AW00011008	Rob	Verhoff	S	2025-10-06	1
13	11009	AW00011009	Shannon	Carlson	S	2025-10-06	1
14	11010	AW00011010	Jacquelyn	Suarez	C	2025-10-06	1

Creating and running the query below allowed me to see all records that are duplicates (or nulls) with ranking identifiers not equal to "1". These will not be used in my data set for this table.

```

1  SELECT
2  *
3  FROM (
4  SELECT
5  *
6  ROW_NUMBER() OVER (PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_last
7  FROM bronze.crm_cust_info
8  )t WHERE flag_last != 1
    
```

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date	flag_last
1	NULL	SF566	NULL	NULL	NULL	NULL	2
2	NULL	13451235	NULL	NULL	NULL	NULL	3
3	29433	AW00029433	NULL	NULL	M	2026-01-25	2
4	29449	AW00029449	NULL	Chen	S	2026-01-25	2
5	29465	AW00029465	Lance	Jimenez	M	2026-01-26	2
6	29466	AW00029466	NULL	NULL	NULL	2026-01-25	3
7	29473	AW00029473	Carmen	NULL	NULL	2026-01-25	2
8	29483	AW00029483	NULL	Navarro	NULL	2026-01-25	2

Creating and running the query below allowed me to see all records with ranking identifiers equal to "1". These records will be used in my data set for this table.

```

1  /*Order all records and their cst_ids based on their creation date starting with
the most recent first and assign an order number to each.
Example: Every customer record with a cst_id that appears once, will be assigned a "1", while
every customer record with a cst_id that appears more than once will be ordered by the most recent
creation date and assigned a rank of 1, 2, 3, etc for each occurrence of that cst_id
*/
2
3  SELECT
4  *
5  FROM (
6  SELECT
7  *
8  ROW_NUMBER() OVER (PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_last
9  FROM bronze.crm_cust_info
10 )t WHERE flag_last = 1
    
```

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date	flag_last
1	NULL	PO25	NULL	NULL	NULL	NULL	1
2	11000	AW00011000	Jon	Yang	M	2025-10-06	1
3	11001	AW00011001	Eugene	Huang	S	2025-10-06	1
4	11002	AW00011002	Ruben	Torres	M	2025-10-06	1
5	11003	AW00011003	Christy	Zhu	S	2025-10-06	1
6	11004	AW00011004	Elizabeth	Johnson	S	2025-10-06	1
7	11005	AW00011005	Julio	Ruiz	S	2025-10-06	1
8	11006	AW00011006	Janet	Alvarez	S	2025-10-06	1
9	11007	AW00011007	Marco	Mehta	M	2025-10-06	1
10	11008	AW00011008	Rob	Verhoff	S	2025-10-06	1
11	11009	AW00011009	Shannon	Carlson	S	2025-10-06	1
12	11010	AW00011010	Jacquelyn	Suarez	C	2025-10-06	1
13	11011	AW00011011	Curtis	Lu	M	2025-10-06	1
14	11012	AW00011012	Lauren	Walker	M	2025-10-06	1

Checking for Unwanted Spaces in String Values

Each field with strings as data types (NVARCHAR) needs to be checked for unwanted spaces and have unwanted spaces removed (leading or trailing spaces). The following columns were identified as being strings and needed to be checked and corrected if applicable.

```
cst_key NVARCHAR(50),
cst_firstname NVARCHAR(50),
cst_lastname NVARCHAR(50),
cst_marital_status NVARCHAR(50),
cst_gndr NVARCHAR(50),
```

The following query was written and executed on the cst_firstname field to test the identification of records with unwanted spaces as depicted below. The same was done for each string value field. Only the 'cst_firstname' and 'cst_lastname' fields had unwanted spaces that required correction.

A screenshot of the SQL Server Management Studio interface. The top pane shows a query window with the following code:

```
1 /* Check string fields for unwanted spaces (leading or trailing)
2 */
3 /*
4  SELECT cst_firstname
5   FROM bronze.crm_cust_info
6  WHERE cst_firstname != TRIM(cst_firstname)
7  --identifies records where the original field and trimmed field are not equal. Those are the records
8  --with unwanted spaces.
```

The bottom pane shows the results of the query, which lists 15 rows of data from the 'cst_firstname' column. A red bracket is drawn under the first two rows, 'Jon' and 'Elizabeth', indicating they are the records identified as having unwanted spaces.

cst_firstname
1 Jon
2 Elizabeth
3 Lauren
4 Ian
5 Chloe
6 Destiny
7 Angela
8 Caleb
9 Wille
10 Ruben
11 Javier
12 Nicole
13 Maria
14 Allison
15 Adrián

The following query was written and executed to produce a table that included all fields from cust_info table, but corrects the wanted spaces in the cst_firstname and cst_lastname fields, and includes only unique records records (non-duplicates) with a ranking of "1" based on their creation date. The ranking identifier is not included in the output.

A screenshot of the SQL Server Management Studio interface. The top pane shows a query window with the following code:

```
1 /* Produce a table from cust_info that corrects unwanted spaces in
2  cst_firstname and cst_lastname, and excludes duplicate records while
3  keeping only the most recently added record for any duplicates
4 */
5 /*
6  SELECT
7   cst_id,
8   cst_key,
9   TRIM(cst_firstname) AS cst_firstname,
10  TRIM(cst_lastname) AS cst_lastname,
11  cst_marital_status,
12  cst_gndr,
13  cst_create_date
14
15  FROM (
16   SELECT
17   *,
18   ROW_NUMBER() OVER (PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_last
19   FROM bronze.crm_cust_info
20   WHERE cst_id IS NOT NULL
21 )t WHERE flag_last = 1
```

The bottom pane shows the results of the query, which is a table with 13 rows of corrected data. The columns are cst_id, cst_key, cst_firstname, cst_lastname, cst_marital_status, cst_gndr, and cst_create_date.

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date
1	11000	AW00011000	Jon	Yang	M	2025-10-05
2	11001	AW00011001	Eugene	Huang	S	2025-10-06
3	11002	AW00011002	Ruben	Torres	M	2025-10-06
4	11003	AW00011003	Christy	Zhu	S	2025-10-06
5	11004	AW00011004	Elizabeth	Johnson	S	2025-10-06
6	11005	AW00011005	Julio	Ruiz	S	2025-10-06
7	11006	AW00011006	Janet	Alvarez	S	2025-10-06
8	11007	AW00011007	Marco	Mehta	M	2025-10-06
9	11008	AW00011008	Rob	Verhoff	S	2025-10-06
10	11009	AW00011009	Shannon	Carlson	S	2025-10-05
11	11010	AW00011010	Jacquelyn	Suarez	S	2025-10-06
12	11011	AW00011011	Curtis	Lu	M	2025-10-06
13	11012	AW00011012	Lauren	Walker	F	2025-10-06

Checking for consistency of values in low cardinality columns (columns used for categorical data)

The columns capturing customer marital status and gender have low quantity of possible variations. The data in these columns needs to be assessed for consistency. Example: The cst_gndr column should only have two possible values → “M” for male and “F” for female. The values in this column should be “M” or “F” only (or maybe NULL). It should be list “Male”, “Female”, “Q”, “Duck”, or anything else. So I run the following query to identify inconsistencies.

The following query was run and produced results for three (3) possible values: M, F, and NULL as depicted below.

cst_gndr
NULL
F
M

Now...to make this data business / reader friendly, we will adjust our primary transformation query for this table so that M is replaced with “Male”, “F” is replaced with “Female” and “NULL” is replaced with “n/a”. The query will also be written so that any unwanted spaces are trimmed and we account for any possible lower case occurrences of “M” or “F”.

The same approach was applied to the marital status field (cst_marital_status) to replace “S”, “M”, and “NULL” with “Single”, “Married”, and “n/a” respectively.

```

4
5
6 /*/
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
    */
SELECT
    cst_id,
    cst_key,
    TRIM(cst_firstname) AS cst_firstname,
    TRIM(cst_lastname) AS cst_lastname,
    CASE WHEN UPPER(TRIM(cst_marital_status)) = 'S' THEN 'Single'
        WHEN UPPER(TRIM(cst_marital_status)) = 'M' THEN 'Married'
        ELSE 'n/a'
    END cst_marital_status,
    CASE WHEN UPPER(TRIM(cst_gndr)) = 'F' THEN 'Female'
        WHEN UPPER(TRIM(cst_gndr)) = 'M' THEN 'Male'
        ELSE 'n/a'
    END cst_gndr,
    cst_create_date
FROM (
    SELECT
        *,
        ROW_NUMBER() OVER (PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_last
    FROM bronze.crm_cust_info
    WHERE cst_id IS NOT NULL
) t WHERE flag_last = 1

```

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date	
1	11000	AW00011000	Jon	Yang	Mamed	Male	2025-10-06
2	11001	AW00011001	Eugene	Huang	Single	Male	2025-10-06
3	11002	AW00011002	Ruben	Tones	Mamed	Male	2025-10-06
4	11003	AW00011003	Christy	Zhu	Single	Female	2025-10-06
5	11004	AW00011004	Elizabeth	Johnson	Single	Female	2025-10-06
6	11005	AW00011005	Julio	Ruiz	Single	Male	2025-10-06
7	11006	AW00011006	Janet	Alvarez	Single	Female	2025-10-06
8	11007	AW00011007	Marco	Mehta	Mamed	Male	2025-10-06
9	11008	AW00011008	Rob	Verhoff	Single	Female	2025-10-06
10	11009	AW00011009	Shannon	Cadeen	Single	Male	2025-10-06

Inserting Customer Info data into Silver layer table

At this point, I am now ready to insert the cleansed data into the Silver layer table as the data quality checks and cleansing complete for the table `crm_cust_info`. An ‘`INSERT INTO`’ command was added to the data quality script so it would execute the load after the data cleansing script is complete.

```

SQLQuery28.sql...indsley (51)*
6   v [INSERT INTO silver.crm_cust_info (
7     cst_id,
8     cst_key,
9     cst_firstname,
10    cst_lastname,
11    cst_marital_status,
12    cst_gndr,
13    cst_create_date)
14
15   SELECT
16     cst_id,
17     cst_key,
18     TRIM(cst_firstname) AS cst_firstname,
19     TRIM(cst_lastname) AS cst_lastname,
20     CASE WHEN UPPER(TRIM(cst_marital_status)) = 'S' THEN 'Single'
21       WHEN UPPER(TRIM(cst_marital_status)) = 'M' THEN 'Married'
22       ELSE 'n/a'
23     END cst_marital_status,
24     CASE WHEN UPPER(TRIM(cst_gndr)) = 'F' THEN 'Female'
25       WHEN UPPER(TRIM(cst_gndr)) = 'M' THEN 'Male'
26       ELSE 'n/a'
27     END cst_gndr,
28     cst_create_date
29
30   FROM
31   *
32   ROW_NUMBER() OVER (PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_last
33   FROM bronze.crm_cust_info
34   WHERE cst_id IS NOT NULL
35   )t WHERE flag_last = 1
  
```

Messages

(18484 rows affected)

Compilation time: 2023-11-26T12:48:19.6286900-03:00

I now check that the data quality corrections and standardizations are correct by running a query for each correction / standardizationWith the load complete. The screenshot below depicts the results of the query to confirm values listed for `cst_gndr` are correctly converted from M, F, and NULL to “Male”, “Female”, and “n/a”....and that no other value types exist

```

SQLQuery28.sql...indsley (51)*
1  --Check for Nulls or Duplicates in Primary Key
2  --Expectation: No Result
3  v SELECT -- Check for duplicate and records with null primary keys
4    cst_id,
5    COUNT(*)
6  FROM silver.crm_cust_info
7  GROUP BY cst_id
8  HAVING COUNT(*) >1 OR cst_id IS NULL
9
10 --Check for cst_key Unwanted Spaces
11 --Expectation: No Result
12 v SELECT -- Check for duplicate and records with null primary keys
13   cst_key
14  FROM silver.crm_cust_info
15  WHERE cst_key != TRIM(cst_key)
16
17 --
18 v SELECT
19   cst_firstname
20  FROM silver.crm_cust_info
21  WHERE cst_firstname != TRIM(cst_firstname)
22
23 --
24 v SELECT
25   cst_lastname
26  FROM silver.crm_cust_info
27  WHERE cst_lastname != TRIM(cst_lastname)
28
29 --Check for data standardization
30 --Expectation: No Result
31 v SELECT DISTINCT
32   cst_marital_status
33  FROM silver.crm_cust_info
34
35 v SELECT DISTINCT
36   cst_gndr
37  FROM silver.crm_cust_info
  
```

Results

cst_gndr
1 n/a
2 Male
3 Female

CRM → Product Info Table

Checking for and Correcting Records with Duplicate and Null Primary Keys

Primary key was evaluated to determine whether or not there were NULL values or duplicates using the following script. This resulted in zero (0) duplicates or null values for the prd_id field

```

13  --> SELECT -- Check for duplicate and records with 'null primary keys'
14      prd_id,
15      COUNT(*)
16  FROM bronze.crm_prd_info
17  GROUP BY prd_id
18  HAVING COUNT(*) >1 OR prd_id IS NULL

```

Deriving New Columns to Enable Mapping

During the data exploration, I found that the prd_key field is a concatenation of the category id 'id' field from the 'erp_px_cat_g1v2' table and additional characters native to this table. The first five characters of each value in this field are the category id and the remaining are just....a bunch of other characters. So, this column needs to separate the first five characters into a new column so the 'erp_px_cat_g1v2' can be mapped to it.

prd_id	prd_key	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt
210	CO-RF-FR-R92B-58	HL Road Frame - Black- 58	NULL	R	2003-07-01 00:00:00.000	NULL

CAT ID

Additionally, I have to correct for the difference between the characters as the "id" in the 'erp_px_cat_g1v2' are separated by a "_", whereas the characters in the prd_key from the prd_info table are separated by a "-".

I also need to create a separate column for the prd_key without the category id characters. So...some standardization is needed here. This is accomplished by the script (with results below) in the following image:

```

2  --> SELECT
3      prd_id,
4      prd_key,
5      PARSENAME(SUBSTRING(prd_key,1,5), '-1', '_') AS cat_id, --Extracts the category id, replaces the "-" with "_",
6      --creates new column so that this table can be mapped to the category id ("id") field in the category table
7      --(erp_px_cat_g1v2)
8      SUBSTRING(prd_key,7,LEN(prd_key)) AS prd_key, -- Removes the first five characters that are the category id from 'erp_px_cat_g1v2'
9      prd_nm,
10     prd_cost,
11     prd_line,
12     prd_start_dt,
13     prd_end_dt
14  FROM bronze.crm_prd_info
15
16  --Check bronze.crm_prd_info for existence of duplicate prd_id

```

prd_id	prd_key	cat_id	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt		
1	CO-RF-FR-R92B-58	CO_RF	FR-R92B-58	HL Read Frame - Black	58	NULL	R	2003-07-01 00:00:00.000	NULL
2	CO-RF-FR-R92B-58	CO_RF	FR-R92B-58	HL Read Frame - Red	58	NULL	R	2003-07-01 00:00:00.000	NULL
3	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet Red	12	S	2011-07-01 00:00:00.000	2007-12-28 00:00:00.000	
4	AC-HE-HL-U509-B	AC_HE	HL-U509-B	Sport-100 Helmet Black	14	S	2011-07-01 00:00:00.000	2007-12-27 00:00:00.000	
5	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet Red	13	S	2011-07-01 00:00:00.000	NULL	
6	AC-HE-HL-U509-B	AC_HE	HL-U509-B	Sport-100 Helmet Black	12	S	2011-07-01 00:00:00.000	2007-12-28 00:00:00.000	
7	AC-HE-HL-U509	AC_HE	HL-U509	Sport-100 Helmet Black	14	S	2012-07-01 00:00:00.000	2008-12-27 00:00:00.000	
8	AC-HE-HL-U509	AC_HE	HL-U509	Sport-100 Helmet Black	13	S	2013-07-01 00:00:00.000	NULL	
9	CL-SO-SO-8909-M	CL_SO	SO-8909-M	Mountain Bike Socks- M	3	M	2011-07-01 00:00:00.000	2007-12-28 00:00:00.000	
10	CL-SO-SO-8909-L	CL_SO	SO-8909-L	Mountain Bike Socks- L	3	M	2011-07-01 00:00:00.000	2007-12-28 00:00:00.000	
11	AC-HE-HL-U509-B	AC_HE	HL-U509-B	Sport-100 Helmet Blue	12	S	2011-07-01 00:00:00.000	2007-12-28 00:00:00.000	
12	AC-HE-HL-U509-B	AC_HE	HL-U509-B	Sport-100 Helmet Blue	14	S	2012-07-01 00:00:00.000	2008-12-27 00:00:00.000	
13	AC-HE-HL-U509-B	AC_HE	HL-U509-B	Sport-100 Helmet Blue	13	S	2013-07-01 00:00:00.000	NULL	

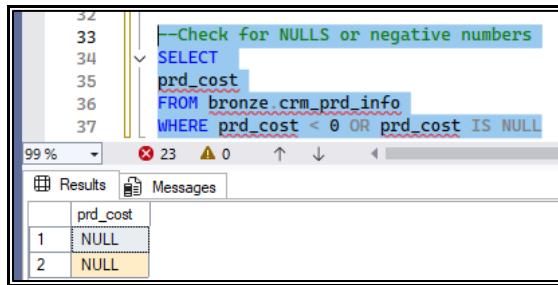
Checking for unwanted spaces

I checked the prd_nm field for unwanted spaces and found none using the following query

```
SELECT  
prd_nm  
FROM bronze.crm_prd_info  
WHERE prd_nm != TRIM(prd_nm)
```

Checking for product costs <0 and Null values

I checked the prd_cost field to identify any products with negative values or nulls using the following query. No products had negative costs, but there are two products with null values for cost. In a business context, the business stakeholder would need to determine whether to retain the null or replace it with a "0" or some other value.



The screenshot shows a SQL query window with the following code:

```
--Check for NULLS or negative numbers  
SELECT  
prd_cost  
FROM bronze.crm_prd_info  
WHERE prd_cost < 0 OR prd_cost IS NULL
```

The results grid shows two rows with the column 'prd_cost' containing 'NULL' for both rows.

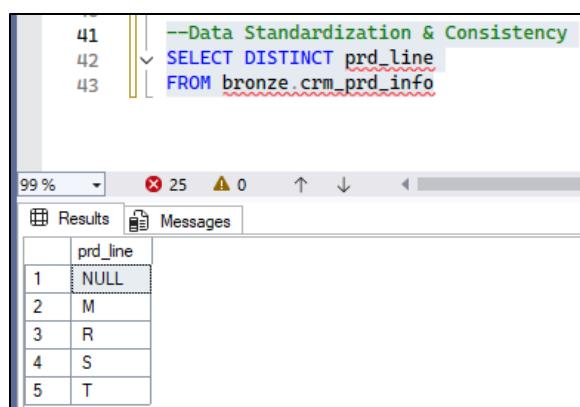
prd_cost
1 NULL
2 NULL

As this is a fictional situation without actual stakeholders, I chose to replace the NULL values with a "0". I did this by adding the following to the transformation script for the prd_info table:

```
ISNULL(prd_cost,0) AS prd_cost, --Replaces null values in cost with 0
```

Checking 'prd_line' field for consistency of values (similar to the check conducted for marital status and gender in the 'crm_cust_info' table)

The 'prd_line' column has a low quantity of possible variations. The data in this column needs to be assessed for consistency. I ran 'SELECT DISTINCT' query on the prd_info table and found the following possible values to be: "NULL", "M", "R", "S", and "T".



The screenshot shows a SQL query window with the following code:

```
--Data Standardization & Consistency  
SELECT DISTINCT prd_line  
FROM bronze.crm_prd_info
```

The results grid shows five distinct values in the 'prd_line' column: 'NULL', 'M', 'R', 'S', and 'T'.

prd_line
1 NULL
2 M
3 R
4 S
5 T

In a real-world situation, I would ask the data owner / business stakeholder what those values represent so that I could convert them to a user / business friendly value. There are no stakeholders in this case as this is a fictional project. So...I added a CASE WHEN with UPPER(TRIM) to the transformation script to replace the values with the following per the instructor for this course:

```
CASE WHEN UPPER(TRIM(prd_line)) = 'M' THEN 'Mountain'
WHEN UPPER(TRIM(prd_line)) = 'R' THEN 'Road'
WHEN UPPER(TRIM(prd_line)) = 'S' THEN 'Other Sales'
WHEN UPPER(TRIM(prd_line)) = 'T' THEN 'Touring'
ELSE 'n/a'
END AS prd_line,
```

Checking the quality of Product Start and End Date fields and correcting quality issues.

The first part of this effort is to identify records that have a product end date that is earlier than the product start date. To do this, I ran the following query and found 201 records in which this occurs.

```
SELECT * FROM bronze.crm_prd_info
WHERE prd_end_dt < prd_start_dt
```

I assumed that the start and end dates for these records were simply reversed and in the incorrect field. So, the start date is actually the end date and the end date is actually the start date. Assuming this is the case, I also found that several records had product end dates that overlap with the start date of the product with the same product key but next prd_id sequence. So....I had some corrections to make.

In a real-world scenario, these observations would need to be reviewed and verified with the data owner before making any decisions on transformations / clean-up. In this case, I pretend that I reviewed the issues with a data owner / business stakeholder and had them confirm the dates were reversed, and that the end dates should not overlap but end the day prior to the start date of the next product id having the same value for 'prd_key'

The below image summarizes observations against a sample set of records with approach for correction:

INITIAL REVIEW OF THE FIELDS									OBSERVATION / NOTE
prd_id	prd_key	cat_id	prd_key	prd_name	prd_cost	prd_line	prd_start_dt	prd_end_dt	
212	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	12	Other Sales	7/1/2011 0:00	12/28/2007 0:00	Prod end dates is earlier than start. Date values are in reverse column order.
213	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	14	Other Sales	7/1/2012 0:00	12/27/2008 0:00	
214	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	13	Other Sales	7/1/2013 0:00	NULL	

PASS AT TRANSPOSING DATES - OVERLAP IDENTIFIED									OBSERVATION / NOTE
prd_id	prd_key	cat_id	prd_key	prd_name	prd_cost	prd_line	prd_start_dt	prd_end_dt	
212	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	12	Other Sales	12/28/2007 0:00	7/1/2011 0:00	Transposed start and end dates so start date is earlier than end date. This reveals an overlap between the end date for prd_id 212 and the start date for prd_id 213 (of the same prod key). Need to set end date to be day prior to start date of the next product id
213	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	14	Other Sales	12/27/2008 0:00	7/1/2012 0:00	
214	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	13	Other Sales	7/1/2013 0:00	NULL	

CORRECTING OVERLAP OF END DATE AND START DATES									OBSERVATION / NOTE
prd_id	prd_key	cat_id	prd_key	prd_name	prd_cost	prd_line	prd_start_dt	prd_end_dt	
212	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	12	Other Sales	7/1/2011 0:00	6/30/2012 0:00	Transposed start and end dates so start date is earlier than end date. Corrected overlap between the end date for prd_id 212 and the start date for prd_id 213 (of the same prod key) by setting end date = start date of next record - 1
213	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	14	Other Sales	7/1/2012 0:00	6/30/2013 0:00	
214	AC-HE-HL-U509-R	AC_HE	HL-U509-R	Sport-100 Helmet- Red	13	Other Sales	7/1/2013 0:00	NULL	

The following commands were added to the script to correct the start and end dates, and to remove the time values from the date field as they were "00:00:00". Below is a screenshot of the full script and results.

```
CAST(prd_start_dt AS DATE) as prd_start_dt,--"CAST" removes time data
CAST(LEAD(prd_start_dt) OVER (PARTITION BY prd_key ORDER BY prd_start_dt)-1 AS DATE)
AS prd_end_dt
```

prd_id	prd_key	cat_id	prd_key	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt	
7	405	AC-FE-FE-6654	AC_FE	FE-6654	Fender Set - Mountain	8	Mountain	2013-07-01	NULL
8	218	AC-HE-HL-U509	AC_HE	HL-U509	Sport-100 Helmet: Black	12	Other Sales	2011-07-01	2012-06-30
9	216	AC-HE-HL-U509	AC_HE	HL-U509	Sport-100 Helmet: Black	14	Other Sales	2012-07-01	2013-06-30
10	217	AC-HE-HL-U509	AC_HE	HL-U509	Sport-100 Helmet: Black	13	Other Sales	2013-07-01	NULL
11	220	AC-HE-HL-U509-B	AC_HE	HL-U509	Sport-100 Helmet: Blue	12	Other Sales	2011-07-01	2012-06-30
12	221	AC-HE-HL-U509-B	AC_HE	HL-U509	Sport-100 Helmet: Blue	14	Other Sales	2012-07-01	2013-06-30
13	222	AC-HE-HL-U509-R	AC_HE	HL-U509	Sport-100 Helmet: Blue	13	Other Sales	2013-07-01	NULL
14	212	AC-HE-HL-U509-R	AC_HE	HL-U509	Sport-100 Helmet: Red	12	Other Sales	2011-07-01	2012-06-30
15	213	AC-HE-HL-U509-R	AC_HE	HL-U509	Sport-100 Helmet: Red	14	Other Sales	2012-07-01	2013-06-30
16	214	AC-HE-HL-U509-R	AC_HE	HL-U509	Sport-100 Helmet: Red	13	Other Sales	2013-07-01	NULL
17	487	AC-HP-HY-1023-70	AC_HP	HY-102...	Hydration Pack - 70 oz.	21	Other Sales	2013-07-01	NULL
18	451	AC-ULT-H902	AC_UL	LT-H902	Headlights - Dual-Beam	14	Road	2012-07-01	NULL
19	452	AC-ULT-H903	AC_UL	LT-H903	Headlights - Weatherp...	19	Road	2012-07-01	NULL
20	453	AC-ULT-H904	AC_UL	LT-H904	Headlights - Weatherp...	17	Road	2012-07-01	2013-06-30

I needed to update the ddl for the prd_info table in the Silver layer to account for the removal of time values from the start and end date columns in the transformation script. I also added the cat_id field to the ddl as that is being created in the transformation script. Below is the updated ddl to account for this.

```
--CREATE SILVER TABLE FOR prd_info data
IF OBJECT_ID ('silver.crm_prd_info', 'U') IS NOT NULL
DROP TABLE silver.crm_prd_info;
CREATE TABLE silver.crm_prd_info (
    prd_id INT,
    cat_id NVARCHAR(50), ←
    prd_key NVARCHAR(50),
    prd_nm NVARCHAR(50),
    prd_cost INT,
    prd_line NVARCHAR(50),
    prd_start_dt DATE, ←
    prd_end_dt DATE, ←
    dwh_create_date DATETIME2 DEFAULT GETDATE()
);
```

Inserting and validating the transformation query results into the Silver layer prd_info table

It is now time to insert the transformed data now that the data quality corrections are complete for the prd_info table and the ddl has been updated.

To insert the corrected data, I wrote and executed an INSERT INTO script ahead of the transformation query for the prd_info table and executed checks against the output. Below is the INSERT INTO with the transformation script:

```

4   v INSERT INTO silver.crm_prd_info (
5     prd_id,
6     cat_id,
7     prd_key,
8     prd_nm,
9     prd_cost,
10    prd_line,
11    prd_start_dt,
12    prd_end_dt
13  )
14
15  SELECT
16    prd_id,
17    REPLACE(SUBSTRING(prd_key, 1, 5), '-' , '_') AS cat_id, --Extracts the category id, replaces the "-" with "_",
18    --creates new column so that this table can be mapped to the category id ("id") field in the category table
19    --(erp_px_cat_glv2)
20    SUBSTRING(prd_key, 7,LEN(prd_key)) AS prd_key, -- Removes the first five characters that are the category id from 'erp_px_cat_glv2'
21    --and creates new column of just the product key characters so we can map to the crm_sales_details table
22    prd_nm,
23    ISNULL(prd_cost,0) AS prd_cost, --Replaces null values in cost with 0
24    CASE UPPER(TRIM(prd_line)) --Replaces single letter values with business-friendly names
25    and removes unwanted spaces /*
26      WHEN 'M' THEN 'Mountain'
27      WHEN 'R' THEN 'Road'
28      WHEN 'O' THEN 'Other Sales'
29      WHEN 'T' THEN 'Touring'
30      ELSE 'n/a'
31    END AS prd_line,
32    CAST(prd_start_dt AS DATE) as prd_start_dt,--"CAST" removes time data
33    CAST(LEAD(prd_start_dt) OVER (PARTITION BY prd_key ORDER BY prd_start_dt)-1 AS DATE) AS prd_end_dt /* create new column
34    in which the end date is equal to original start date -1 to correct for the transposed data in these columns.*/
35  FROM bronze.crm_prd_info

```

The next step was to run checks against the transformed data and derived fields in the prd_info table in the silver table. Each of the following scripts were run and executed. All were successful.

```

50  --Check silver.crm_prd_info for existence of duplicate prd_id
51  v SELECT -- Check for duplicate and records with null primary keys
52    prd_id,
53    COUNT(*)
54  FROM silver.crm_prd_info
55  GROUP BY prd_id
56  HAVING COUNT(*) > 1 OR prd_id IS NULL
57
58  --Check for cst_key Unwanted Spaces
59  --Expectation: No Result
60  v SELECT
61    prd_nm
62    FROM silver.crm_prd_info
63    WHERE prd_nm != TRIM(prd_nm)
64
65  --Check for NULLs or negative numbers
66  --Expectation: No Result
67  v SELECT
68    prd_cost
69    FROM silver.crm_prd_info
70    WHERE prd_cost < 0 OR prd_cost IS NULL
71
72  --Data Standardization & Consistency
73  v SELECT DISTINCT prd_line
74    FROM silver.crm_prd_info
75
76  --Checking for invalid product start and end dates
77  v SELECT
78    *
79    FROM silver.crm_prd_info
80    WHERE prd_end_dt < prd_start_dt --Checks for records with end dates that are earlier than the prd start date
81
82  v SELECT
83    *
84    FROM silver.crm_prd_info --View full data set in the silver table to get a warm fuzzy|

```

Results

prd_id	cat_id	prd_key	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt	dwh_create_date	
1	478	AC_BC	BC-M005	Mountain Bottle Cage	4	Mountain	2013-07-01	NULL	2025-11-27 09:56:50.0500000
2	479	AC_BC	BC-R205	Road Bottle Cage	3	Road	2013-07-01	NULL	2025-11-27 09:56:50.0500000
3	477	AC_BC	WB-H098	Water Bottle - 30 oz.	2	Other Sales	2013-07-01	NULL	2025-11-27 09:56:50.0500000
4	483	AC_BR	RA-H123	Hitch Rack - 4-Bike	45	Other Sales	2013-07-01	NULL	2025-11-27 09:56:50.0500000
5	486	AC_BS	ST-1401	All-Purpose Bike Stand	59	Mountain	2013-07-01	NULL	2025-11-27 09:56:50.0500000
6	484	AC_CL	CL-9009	Bike Wash - Dissolver	3	Other Sales	2013-07-01	NULL	2025-11-27 09:56:50.0500000

DISCLAIMER: Now....I completely understand that this is best case scenario as I was following step-by-step with the instructor. I likely would have dorked this up completely and not known how to correct any issues if running this on my own. It would have been a disaster in all likelihood. But I'm learning....baby steps. This is fun.....so far.

CRM → Sales Details Table

Checking primary key (sls_ord_num) for unwanted spaces:

Executed the following query with no results for records with unwanted spaces:

```

Cust_Info_Fla...lindsay (68) | cust_info_tran...lindsay (60))
1  SELECT
2   sls_ord_num,
3   sls_prd_key,
4   sls_cust_id,
5   sls_order_dt,
6   sls_ship_dt,
7   sls_due_dt,
8   sls_sales,
9   sls_quantity,
10  sls_price
11  FROM bronze.crm_sales_details
12  WHERE sls_ord_num != TRIM(sls_ord_num)

```

Results

sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_quantity	sls_price

Checking the sls_prd_key to determine whether or not there are matches in the prd_key of the the crm_prd_info table: This is necessary as I determined earlier in the data diagram that I would need to using a portion of the values in this column to match to the prd_key field in the crm_prd_info table. The following was written and executed...resulting in zero (0) matches.

```

Cust_Info_Fla...lindsay (68) | cust_info_tran...lindsay (60))
1  SELECT
2   sls_ord_num,
3   sls_prd_key,
4   sls_cust_id,
5   sls_order_dt,
6   sls_ship_dt,
7   sls_due_dt,
8   sls_sales,
9   sls_quantity,
10  sls_price
11  FROM bronze.crm_sales_details
12  WHERE sls_prd_key NOT IN (SELECT prd_key FROM silver.crm_prd_info)

```

Results

sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_quantity	sls_price

Checking the sls_cust_id to determine whether or not there are matches in the cst_id field of the crm_cust_info table: This is necessary as I determined earlier in the data diagram that I would need to using a portion of the values in this column to match to the cst_id field in the crm_cust_info table. The following was written and executed...resulting in zero (0) matches.

```

Cust_Info_Fla...lindsay (68) | cust_info_tran...lindsay (60))
1  SELECT
2   sls_ord_num,
3   sls_prd_key,
4   sls_cust_id,
5   sls_order_dt,
6   sls_ship_dt,
7   sls_due_dt,
8   sls_sales,
9   sls_quantity,
10  sls_price
11  FROM bronze.crm_sales_details
12  WHERE sls_cust_id NOT IN (SELECT cst_id FROM silver.crm_cust_info)

```

Results

sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_quantity	sls_price

Check values in the following fields to identify non-date values: sls_order_dt; sls_ship_dt; and sls_due_dt. I ran the following query and it returned 17 records with a “0” in the sls_order_dt field

```
SELECT  
    sls_order_dt  
FROM bronze.crm_sales_details  
WHERE sls_order_dt <=0
```

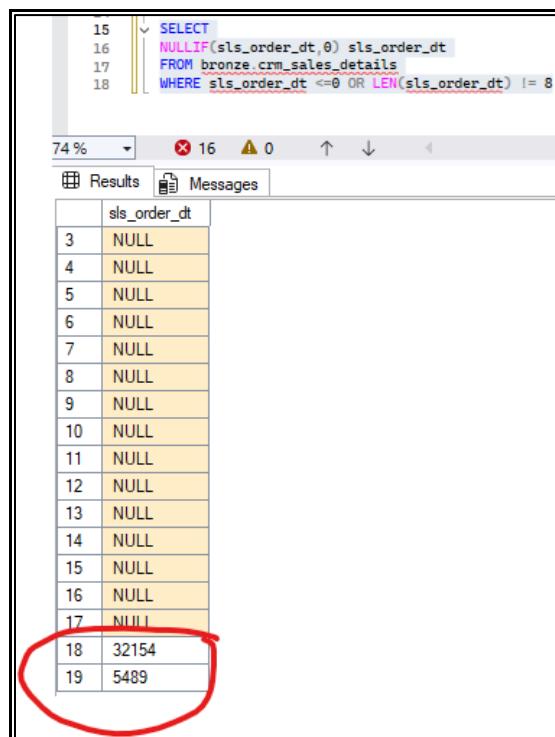
I needed to convert the “0” for these records to NULL. The following was written to correct for this:

```
SELECT  
    NULLIF(sls_order_dt,0) sls_order_dt  
FROM bronze.crm_sales_details  
WHERE sls_order_dt <=0
```

A further check of the sls_order_dt field was done to identify:

- Any additional non-date numerical values that can not be converted to a date
- Any dates beyond 20500101 (fictional upper boundary chosen for this exercise)
- Any dates earlier than 19990101 (fictional lower boundary chosen for this exercise).

This query was run and returned all NULLs as well as two records with non-date numerical values.



	sls_order_dt
3	NULL
4	NULL
5	NULL
6	NULL
7	NULL
8	NULL
9	NULL
10	NULL
11	NULL
12	NULL
13	NULL
14	NULL
15	NULL
16	NULL
17	NULL
18	32154
19	5489

Converting the values in the following fields from INT to dates: sls_order_dt; sls_ship_dt; and sls_due_dt. This is necessary because the values should be dates but are written as "20101229" for example. They should be in this format "2010-12-29"

I added the following to the transformation query to convert non-date numerical values to NULL and set all others to dates:

```

Cust_Info_Fla...lindsay (68) | cust_info_tran...lindsay (60)

1 | SELECT
2 |     sls_ord_num,
3 |     sls_prd_key,
4 |     sls_cust_id,
5 |     CASE WHEN sls_order_dt = 0 OR LEN(sls_order_dt) != 8 THEN NULL
6 |     ELSE CAST(CAST(sls_order_dt AS VARCHAR) AS DATE) --Order here is casting the integer values to VARCHAR then
7 |     --casting the VARCHAR to a date
8 |     END AS sls_order_dt,
9 |     sls_ship_dt,
10 |     sls_due_dt,
11 |     sls_sales,
12 |     sls_quantity,
13 |     sls_price
14 |     FROM bronze.crm_sales_details
15 |
16 | 74 % 18 ▲ 0 ↑ ↓
17 | Results Messages
18 |
```

	sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_quantity	sls_price
4...	SO69210	HL-U509	12426	2013-10-26	20131102	20131107	35	1	35
4...	SO69211	TI-M823	12432	2013-10-26	20131102	20131107	35	1	35
4...	SO69211	FE-6654	12432	2013-10-26	20131102	20131107	22	1	22
4...	SO69211	HL-U509-R	12432	2013-10-26	20131102	20131107	35	1	35
4...	SO69212	WB-H098	11495	2013-10-26	20131102	20131107	5	1	5
4...	SO69212	BC-R205	11495	2013-10-26	20131102	20131107	9	1	9
4...	SO69213	FE-6654	13368	2013-10-26	20131102	20131107	22	1	22
4...	SO69213	LJ-0192-S	13368	2013-10-26	20131102	20131107	50	1	50
4...	SO69214	BK-M68B-38	18595	2013-10-26	20131102	20131107	2295	1	2295
4...	SO69214	TI-M823	18595	2013-10-26	20131102	20131107	35	1	35
4...	SO69215	BK-M68B-46	16864	2013-10-26	20131102	20131107	2295	1	2295
4...	SO69215	TI-M823	16864	NULL	20131102	20131107	-35	1	35
4...	SO69215	TT-M928	16864	NULL	20131102	20131107	5	1	5
4...	SO69215	HL-U509	16864	2013-10-26	20131102	20131107	35	1	35

The same series of checks and corrections done on the sls_order_dt field were conducted for the remaining date fields and added to the transformation query. No quality issues were identified, but I added the same CASE WHEN transformation script to each date field to ensure similar potential quality issues that might arise in future data sets are automatically corrected.

```

Cust_Info_Fla...lindsay (68) | cust_info_tran...lindsay (60)

1 | SELECT
2 |     sls_ord_num,
3 |     sls_prd_key,
4 |     sls_cust_id,
5 |     sls_order_dt,
6 |     CASE WHEN sls_order_dt = 0 OR LEN(sls_order_dt) != 8 THEN NULL
7 |     ELSE CAST(CAST(sls_order_dt AS VARCHAR) AS DATE) --Order here is casting the integer values to VARCHAR then
8 |     --casting the VARCHAR to a date
9 |     END AS sls_order_dt,
10 |     sls_ship_dt,
11 |     CASE WHEN sls_ship_dt = 0 OR LEN(sls_ship_dt) != 8 THEN NULL
12 |     ELSE CAST(CAST(sls_ship_dt AS VARCHAR) AS DATE) --Order here is casting the integer values to VARCHAR then
13 |     --casting the VARCHAR to a date
14 |     END AS sls_ship_dt,
15 |     sls_due_dt,
16 |     CASE WHEN sls_due_dt = 0 OR LEN(sls_due_dt) != 8 THEN NULL
17 |     ELSE CAST(CAST(sls_due_dt AS VARCHAR) AS DATE) --Order here is casting the integer values to VARCHAR then
18 |     --casting the VARCHAR to a date
19 |     END AS sls_due_dt,
20 |     sls_sales,
21 |     sls_quantity,
22 |     sls_price
23 |     FROM bronze.crm_sales_details
24 |
25 | 74 % 23 ▲ 0 ↑ ↓
26 | Results Messages
27 |
```

	sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_quantity	sls_price
1	SO43697	BK-R93R-62	21768	2010-12-29	2011-01-05	2011-01-10	3578	1	3578
2	SO43698	BK-M825-44	28389	2010-12-29	2011-01-05	2011-01-10	3400	1	3400
3	SO43699	BK-M825-44	25863	2010-12-29	2011-01-05	2011-01-10	3400	1	3400
4	SO43700	BK-R50B-62	14501	2010-12-29	2011-01-05	2011-01-11	699	1	699
5	SO43701	BK-M825-44	11003	2010-12-29	2011-01-05	2011-01-10	3400	1	3400
6	SO43702	BK-R93R-44	27645	2010-12-30	2011-01-06	2011-01-11	3578	1	3578
7	SO43703	BK-R93R-62	16624	2010-12-30	2011-01-06	2011-01-11	3578	1	3578
8	SO43704	BK-M82B-48	11005	2010-12-30	2011-01-06	2011-01-11	3375	1	3375
9	SO43705	BK-M825-38	11011	2010-12-30	2011-01-06	2011-01-11	3400	1	3400
10	SO43706	BK-R93R-48	27621	2010-12-31	2011-01-07	2011-01-12	3578	1	3578
11	SO43707	BK-R93R-48	27616	2010-12-31	2011-01-07	2011-01-11	3578	1	3578
12	SO43708	BKR50R-52	20042	2010-12-31	2011-01-07	2011-01-12	699	1	699
13	SO43709	BK-R93R-52	16351	2010-12-31	2011-01-07	2011-01-12	3578	1	3578
14	SO43710	BK-R93R-56	16517	2010-12-31	2011-01-07	2011-01-12	3578	1	3578
15	SO43711	BK-R93R-56	27606	2011-01-01	2011-01-07	2011-01-11	3578	1	3578
16	SO43712	BK-R93B-44	13513	2011-01-01	2011-01-08	2011-01-13	3578	1	3578

Checking for invalid dates. The purpose of this check is to ensure that the ship dates and due dates are not earlier than the order dates. The following query was run and returned no results - confirming that order dates are not later than ship or due dates.

```

SELECT
*
FROM bronze.crm_sales_details
WHERE sls_order_dt > sls_ship_dt OR sls_order_dt > sls_due_dt

```

Checking quality of data for sales, quantity and price. In this fictional project, the following business rules have been defined for these fields:

- Sales must be equal to the product of quantity and price ($\text{sales} = \text{quantity} * \text{price}$).
- Negative values, zeros, and Nulls are not allowed

To check for data consistency, I ran the following query that returned results with records in all three fields that required correction.

	sls_sales	sls_quantity	sls_price
1	0	1	10
2	16	2	NULL
3	NULL	1	9
4	NULL	1	35
5	35	1	NULL
6	30	1	-30
7	70	2	NULL
8	10	2	NULL
9	NULL	1	24
10	21	1	-21
11	40	1	2
12	-18	1	9
13	NULL	1	10
14	100	10	NULL

In a real-world situation, I would need to meet with the data owners / business stakeholders to review the data quality issues and gain alignment on a decision to resolve them. The decision could be for them to fix the data in the source system or to correct them through transformations in the data warehouse until they are able to correct them in the source system.

As this is a fictional situation, the following rules were applied to correct the data quality issues identified for sales, price, and quantity.

- If Sales is negative, zero, or null, derive it using Quantity * Price
- If Price is zero or Null, calculate it using sales / quantity
- If Price is negative, convert to a positive number

The following corrections were added to the transformation query using the above rules.

```

sls_price AS old_sls_price,
CASE WHEN sls_sales IS NULL OR sls_sales <= 0 OR sls_sales != sls_quantity * ABS(sls_price)
      THEN sls_quantity * ABS(sls_price)
      ELSE sls_sales
END AS sls_sales,
sls_quantity,
CASE WHEN sls_price IS NULL OR sls_price <=0
      THEN sls_sales / NULLIF(sls_quantity,0)
      ELSE sls_price
END AS sls_price

```

Finalize transformation query for crm_sales_details, update ddl, and INSERT into the silver table. After adding the above to the transformation query, I also had to update the data types for the date fields in the ddl script for the sales_details table from INT to DATE and re-execute. I then prepared the INSERT INTO script to insert the corrected data into the silver table

```
SQLQuery35.s...lindsay (72)* Silver_Layer_T...lindsay (51)* Silver_Layer_C...lindsay (62) prd_info_trans...lindsay (63)
Cust_Info_Fla...lindsay (68)) cust_info_tran...lindsay (60)
1   INSERT INTO silver.crm_sales_details (
2     sls_ord_num,
3     sls_prd_key,
4     sls_cust_id,
5     sls_order_dt,
6     sls_ship_dt,
7     sls_due_dt,
8     sls_sales,
9     sls_quantity,
10    sls_price
11  )
12
13  SELECT
14    sls_ord_num,
15    sls_prd_key,
16    sls_cust_id,
17    CASE WHEN sls_order_dt = 0 OR LEN(sls_order_dt) != 8 THEN NULL
18      ELSE CAST(CAST(sls_order_dt AS VARCHAR) AS DATE) --Order here is casting the integer values to VARCHAR then
19      --casting the VARCHAR to a date
20    END AS sls_order_dt,
21    CASE WHEN sls_ship_dt = 0 OR LEN(sls_ship_dt) != 8 THEN NULL
22      ELSE CAST(CAST(sls_ship_dt AS VARCHAR) AS DATE) --Order here is casting the integer values to VARCHAR then
23      --casting the VARCHAR to a date
24    END AS sls_ship_dt,
25    CASE WHEN sls_due_dt = 0 OR LEN(sls_due_dt) != 8 THEN NULL
26      ELSE CAST(CAST(sls_due_dt AS VARCHAR) AS DATE) --Order here is casting the integer values to VARCHAR then
27      --casting the VARCHAR to a date
28    END AS sls_due_dt,
29    CASE WHEN sls_sales IS NULL OR sls_sales <= 0 OR sls_sales != sls_quantity * ABS(sls_price)
30      THEN sls_quantity * ABS(sls_price)
31      ELSE sls_sales
32    END AS sls_sales,
33    sls_quantity,
34    CASE WHEN sls_price IS NULL OR sls_price <= 0
35      THEN sls_sales / NULLIF(sls_quantity, 0)
36      ELSE sls_price
37    END AS sls_price
38  FROM bronze.crm_sales_details
39
40  74 % 74 0 ↑ ↓
41  Messages
42
43  (60398 rows affected)
44
45  Completion time: 2025-11-27T21:22:40.4758053-05:00
```

Check quality of data inserted into the the crm_sales_details table post transformation. All of the same quality checks run against the bronze table for crm_sales_details were executed on the silver table for crm_sales_details to ensure that all transformations works, and that there are no additional quality issues requiring correction.

ERP → cust_az12 (The table with customer ids, birthdates and gender)

Checking primary key (cid) to determine what needs to be done to connect to crm_cust_info via the cst_key

Running separate queries against both tables allows me to see that some of the records in the cid field in the cust_az12 have extra characters ("NAS") that do not appear in the cst_key field. I will need to account for this to connect both tables.

The screenshot shows a SQL query window with two queries and their results.

Query 1:

```
SELECT * FROM bronze.erp_cust_az12
```

Query 2:

```
SELECT * FROM silver.crm_cust_info
```

Results:

cid	bdate	gen
1... NASAW00022038	1985-02-14	Female
1... NASAW00022039	1979-09-13	Female
1... NASAW00022040	1980-04-14	Male
1... NASAW00022041	1978-10-31	Female
1... AW00022042	1983-08-18	Female
1... AW00022043	1978-02-09	Female
1... AW00022044	1983-05-13	Male
1... AW00022045	1979-05-04	Male
1... AW00022046	1984-08-16	Female

I wrote the following transformation query to remove the “NAS” from any records that included it in the cid field, then added the CASE WHEN to a WHERE statement to check if there were any records from the correction that did not appear in the silver table of `crm_cust_info`.

```

SQLQuery36...lndslay (51)* -> X
1  SELECT
2    cid,
3    CASE WHEN cid LIKE 'NAS%' THEN SUBSTRING(cid, 4, LEN(cid))
4      ELSE cid
5    END AS cid,
6    bdate,
7    gen
8  FROM bronze.erp_cust_az12
99 % 14 0 ↑ ↓
Results Messages
cid cid bdate gen
1 NASAW00011000 AW00011000 1971-10-06 Male
2 NASAW00011001 AW00011001 1976-05-10 Male
3 NASAW00011002 AW00011002 1971-02-09 Male
4 NASAW00011003 AW00011003 1973-08-14 Female
5 NASAW00011004 AW00011004 1979-08-05 Female
6 NASAW00011005 AW00011005 1976-08-01 Male

```



```

SQLQuery36...lndslay (51)* -> X
1  SELECT
2    cid,
3    CASE WHEN cid LIKE 'NAS%' THEN SUBSTRING(cid, 4, LEN(cid))
4      ELSE cid
5    END AS cid,
6    bdate,
7    gen
8  FROM bronze.erp_cust_az12
9  WHERE CASE WHEN cid LIKE 'NAS%' THEN SUBSTRING(cid, 4, LEN(cid))
10    ELSE cid
11    END NOT IN (SELECT DISTINCT cst_key FROM silver.crm_cust_info)
12
99 % 14 0 ↑ ↓
Results Messages
cid cid bdate gen

```

Check birthdates for values outside of upper and lower date boundaries

The following query was run to identify anyone that is ≥ 100 years old or have birthdates in the future. The results returned records that meet both criteria and would require correction. In a real world situation, I would need to meet with the data owner / business stakeholder to determine how these should be treated. Since this is a fictional situation, I will pretend that I met with the stakeholder and was informed that only birthdates in the future would be adjusted in the data set and they would be replaced with NULL. This was accomplished using a CASE WHEN and GETDATE.

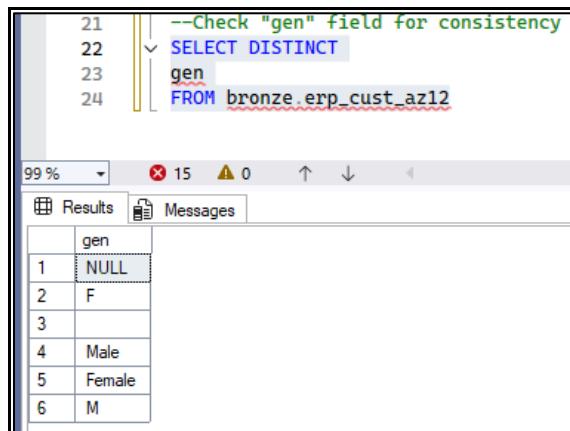
```

--Check for out-of-range-dates
13  SELECT DISTINCT
14    bdate
15  FROM bronze.erp_cust_az12
16  WHERE bdate < '1924-01-01' OR bdate > GETDATE()
17
99 % 12 0 ↑ ↓
Results Messages
bdate
1 1916-02-10
2 1917-02-09
3 1917-06-05
4 1917-09-20
5 1918-02-11
6 1918-11-08
7 1919-02-14
8 1919-03-10
9 1920-11-14
10 1922-01-02
11 1922-04-10
12 1922-06-06
13 1923-04-14
14 1923-08-16
15 1923-11-16
16 2038-10-17
17 2042-02-22
18 2045-03-03
19 2050-05-21

```

Check gen field for data standardization

This is a categorization field and should have standard values. A quick query on distinct values in the gen field resulted in inconsistent values that will require standardization ("F", "Female", "M", "Male", etc.). Additionally, there are empty values and nulls. In a real world situation, I would meet with the business stakeholder to obtain a decision on the appropriate business-friendly correction and how to account for / correct the NULLs and empty fields.

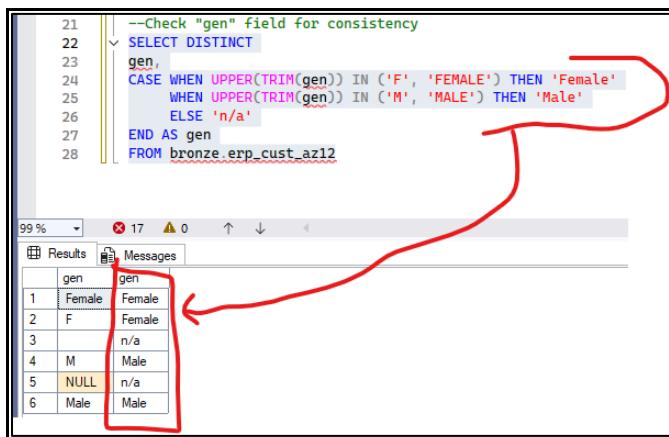


```
--Check "gen" field for consistency
SELECT DISTINCT
    gen
FROM bronze.erp_cust_az12
```

The screenshot shows a SQL query results window. The query is: --Check "gen" field for consistency SELECT DISTINCT gen FROM bronze.erp_cust_az12. The results table has one column 'gen' with 6 rows: 1. NULL, 2. F, 3. (empty), 4. Male, 5. Female, 6. M.

gen
NULL
F
Male
Female
M

The following query was written to standardize values in the gen field.



```
--Check "gen" field for consistency
SELECT DISTINCT
    gen,
    CASE WHEN UPPER(TRIM(gen)) IN ('F', 'FEMALE') THEN 'Female'
        WHEN UPPER(TRIM(gen)) IN ('M', 'MALE') THEN 'Male'
        ELSE 'n/a'
    END AS gen
FROM bronze.erp_cust_az12
```

The screenshot shows a SQL query results window. The query is: --Check "gen" field for consistency SELECT DISTINCT gen, CASE WHEN UPPER(TRIM(gen)) IN ('F', 'FEMALE') THEN 'Female' WHEN UPPER(TRIM(gen)) IN ('M', 'MALE') THEN 'Male' ELSE 'n/a' END AS gen FROM bronze.erp_cust_az12. The results table has two columns: 'gen' and 'gen'. The 'gen' column has 6 rows: 1. Female, 2. F, 3. n/a, 4. M, 5. NULL, 6. Male. The 'gen' column has 6 rows: 1. Female, 2. Female, 3. n/a, 4. Male, 5. n/a, 6. Male. Red arrows highlight the CASE WHEN logic and the resulting standardized values in the second column.

gen	gen
Female	Female
F	Female
n/a	n/a
M	Male
NULL	n/a
Male	Male

Insert erp_cust_az12 into Silver Layer

The corrections identified for erp_cust_az12 were consolidated into a single transformation query along with an INSERT INTO command to insert the corrected data into the silver layer with quality checks completed to ensure data fed correctly

```
--INSERT transformation for erp_cust_az12 into the silver layer
--INSERT INTO silver.erp_cust_az12 (
cid,
bdate,
gen)
SELECT
CASE WHEN cid LIKE 'NAS%' THEN SUBSTRING(cid, 4, LEN(cid))
ELSE cid
END AS cid,
CASE WHEN bdate > GETDATE() THEN NULL
ELSE bdate
END AS bdate,
CASE WHEN UPPER(TRIM(gen)) IN ('F', 'FEMALE') THEN 'Female'
WHEN UPPER(TRIM(gen)) IN ('M', 'MALE') THEN 'Male'
ELSE 'n/a'
END AS gen
FROM bronze.erp_cust_az12

--Check for out-of-range-dates
SELECT DISTINCT
bdate
```

ERP → loca101 (The table with customer id and country)

Checking primary key (cid) to determine what needs to be done to connect to crm_cust_info via the cst_key

Running separate queries against both tables allows me to see that the cid field in the loc_a101 table has a “-” between the first two string characters and the remaining numerical characters in the field. The “-” will need to be removed to enable mapping between the two tables.

The following script was written to remove the “-” from the cid field in loc_a101. I also ran the same query with a NOT IN subquery to determine whether any of the transformed values did not have matches in the cust_info table

```
SQLQuery37.s...indsley (65)* - X erp_cust_az12...lind
1  SELECT
2    REPLACE(cid, '-', '') cid,
3    ctry
4    FROM bronze.erp_loc_a101
5

99 %   x 6  ▲ 0  ↑  ↓  ◀
Results Messages



|   | cid        | ctry      |
|---|------------|-----------|
| 1 | AW00011000 | Australia |
| 2 | AW00011001 | Australia |
| 3 | AW00011002 | Australia |
| 4 | AW00011003 | Australia |


```

Checking the country field to identify consistency / needs for standardization

I ran a SELECT DISTINCT query to identify all possible values and identify needs for standardization. There is a mix of abbreviations, fully spelled out country names, NULLs, and empty values.

```

SQLQuery37...lindsay (65)* -> X erp_cust_az12...lindsay (51)*
1   SELECT
2     REPLACE(cid, '-', '') cid,
3     ctry
4   FROM bronze.erp_loc_a101;
5
6   SELECT DISTINCT
7     ctry
8   FROM bronze.erp_loc_a101
9

```

	cntry
1	DE
2	USA
3	Germany
4	United States
5	NULL
6	Australia
7	United Kingdom
8	
9	Canada
10	France
11	US

The following query was written and added to the transformation query to correct for the inconsistent country identifiers. I also added the same CASE WHEN to the previous query to identify how the matches mapped to their original values to confirm correction.

```

SQLQuery37.s...lindsay (65)* -> X erp_cust_az12...lindsay (51)*
1   SELECT
2     REPLACE(cid, '-', '') cid,
3     CASE WHEN TRIM(ctry) = 'DE' THEN 'Germany'
4       WHEN TRIM(ctry) IN ('US', 'USA') THEN 'United States'
5       WHEN TRIM(ctry) = '' or ctry IS NULL THEN 'n/a'
6       ELSE TRIM(ctry)
7     END AS ctry
8   FROM bronze.erp_loc_a101;

```

```

SQLQuery37.s...lindsay (65)* -> X erp_cust_az12...lindsay (51)*
9
10
11
12   --CHECK COUNTRY FOR CONSTITENCY & STANDARDIZATION--
13
14   SELECT DISTINCT
15     ctry AS old_ctry,
16     CASE WHEN TRIM(ctry) = 'DE' THEN 'Germany'
17       WHEN TRIM(ctry) IN ('US', 'USA') THEN 'United States'
18       WHEN TRIM(ctry) = '' or ctry IS NULL THEN 'n/a'
19       ELSE TRIM(ctry)
20     END AS ctry
21   FROM bronze.erp_loc_a101

```

	old_ctry	ctry
1	US	United States
2	United States	United States
3	DE	Germany
4	Canada	Canada
5	Australia	Australia
6	France	France
7	USA	United States
8	NULL	n/a
9	United Kingdom	United Kingdom
10		n/a
11	Germany	Germany

Insert loc_a101 into silver layer

I wrote an INSERT INTO query to insert the transformed data from the bronze layer into the silver layer and ran a query to review and confirm that corrections to cid and cntry fields appeared correctly in the silver layer

```
1  ✓ INSERT INTO silver.erp_loc_a101 (
2    cid,
3    cntry
4  )
5  SELECT
6    REPLACE(cid, '-', '') cid,
7    CASE WHEN TRIM(cntry) = 'DE' THEN 'Germany'
8      WHEN TRIM(cntry) IN ('US', 'USA') THEN 'United States'
9      WHEN TRIM(cntry) = '' or cntry IS NULL THEN 'n/a'
10     ELSE TRIM(cntry)
11   END AS cntry
12  FROM bronze.erp_loc_a101;
```

	cid	cntry	dwh_create_date
1...	AW00029465	Australia	2025-11-27 22:48:44.3600000
1...	AW00029466	Germany	2025-11-27 22:48:44.3600000
1...	AW00029467	Germany	2025-11-27 22:48:44.3600000
1...	AW00029468	United King...	2025-11-27 22:48:44.3600000
1...	AW00029469	United King...	2025-11-27 22:48:44.3600000
1...	AW00029470	France	2025-11-27 22:48:44.3600000
1...	AW00029471	France	2025-11-27 22:48:44.3600000
1...	AW00029472	France	2025-11-27 22:48:44.3600000
1...	AW00029473	United King...	2025-11-27 22:48:44.3600000
1...	AW00029474	Germany	2025-11-27 22:48:44.3600000
1...	AW00029475	Germany	2025-11-27 22:48:44.3600000

ERP → px_cat_g1v2 (The table with product category ids, product categories, subcategories, and maintenance identifier)

Checking primary key (id) to determine what needs to be done to connect to crm_prod_info via the cat_id that was created earlier in for the silver layer of this table

a. Running separate queries against both tables allows me to see the values cat_id in crm_prod_info appears to be in the same format and structure as the values listed in the id field of px_cat_g1v2 and is ready to be used for mapping

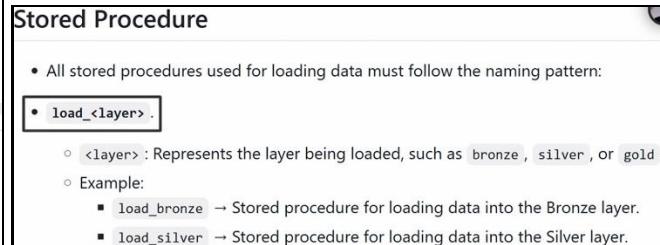
Checking the “cat” field to ensure consistency, removal of unwanted spaces, etc.

- Running a query to identify unwanted spaces across all fields returned zero results
- Running a query to identify inconsistencies in the cat field resulted in no consistent value usage
- Running a query to identify inconsistencies in the subcat field resulted in no consistent value usage
- Running a query to identify inconsistencies in the maintenance field resulted in no consistent value usage

With all checks complete, I wrote the INSERT INTO query to place the data into the silver layer. Then built a query that consolidated the INSERT queries for all tables in the silver layer with TRUNCATE TABLE commands for each.

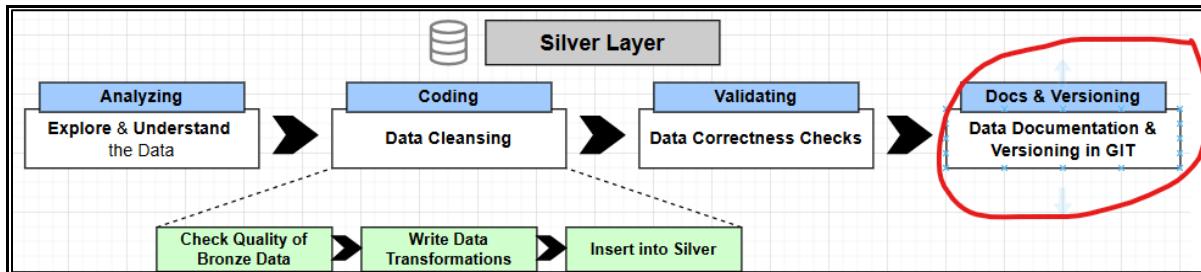
Creating a stored procedure:

a. I wrote a stored procedure for the loading of the data into the database as it is intended to be a recurring action to bring in the data.

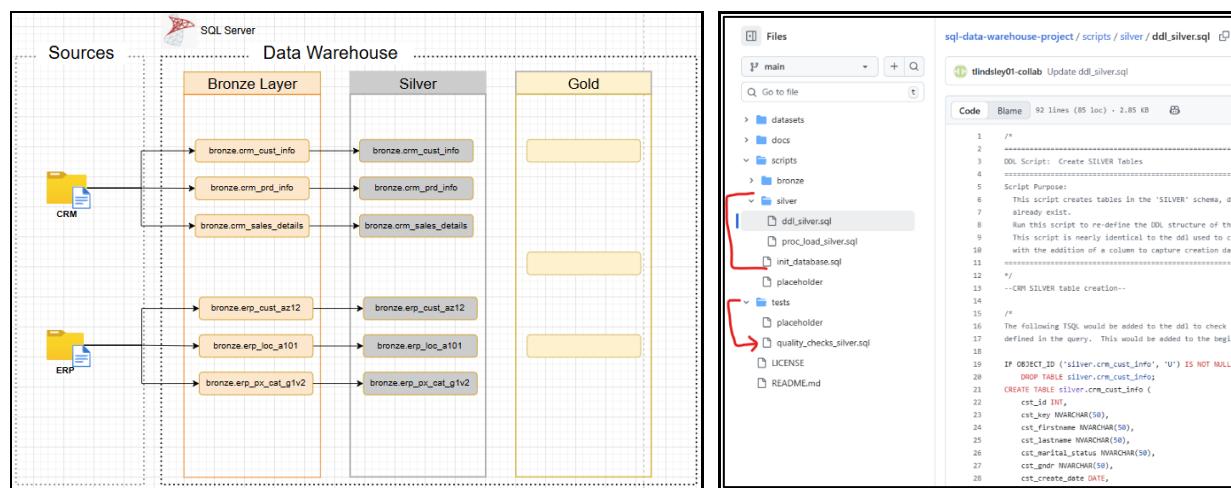


Once executing the stored procedure and refreshing the DB to visually confirm its existence, I then created a new query to test the execution of the procedure. The following SQL statement was used: EXEC silver..load silver

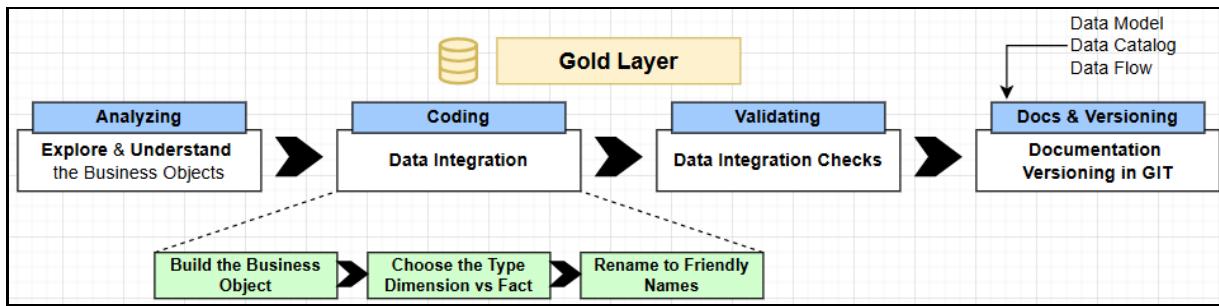
d. Documenting & Commit Code in GIT



I updated the data flow to reflect the lineage of data from source systems to bronze to silver. I also loaded and committed my ddl script, stored procedure, and quality checks for the silver layer to my Git repo.



8. Build Gold Layer



a. Analyzing: Explore Business Objects

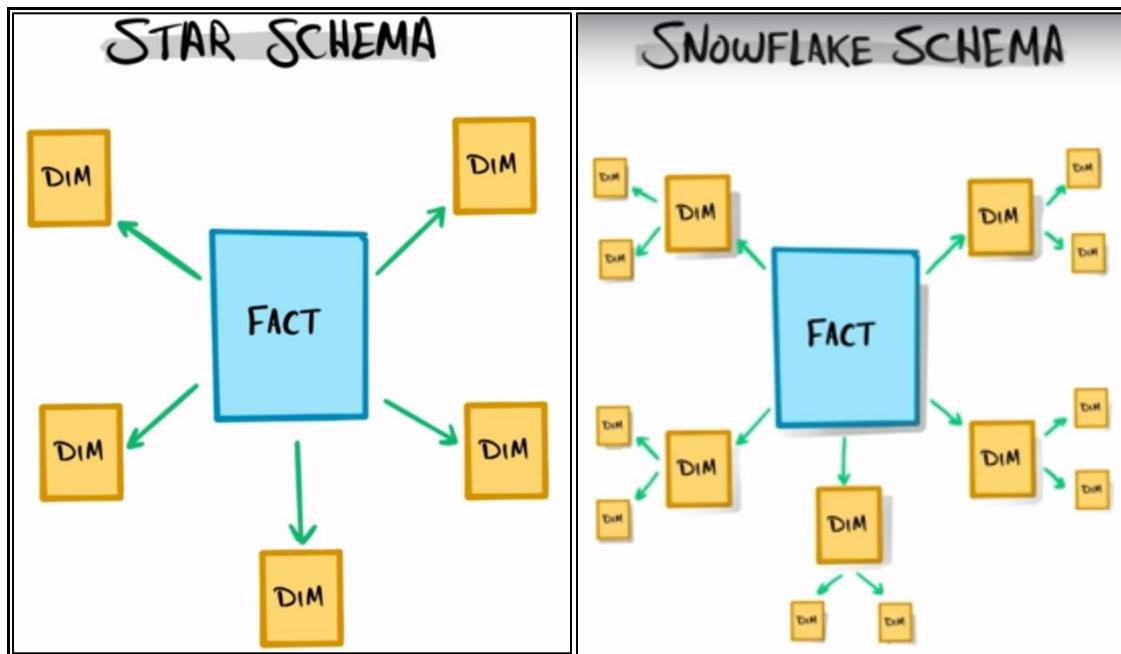
Data Model Selection

The first step in this task is to select the logical data model to enable development of relationships between data sets. Additionally, it serves as a reference for future users of the data warehouse.

The data model needs to be optimized for data analytics. It also needs to be flexible, scalable and easy to understand. The types of data models typically used are Star Schema and Snowflake Schema.

Star Schema: A central fact table in the center and surrounded by dimensions. The fact table contains transactional / event-oriented data, and the dimensions contain descriptive information and the relationship between the fact table and the dimensions around it form a star shape when depicted visually....hence “Star Schema”. It is simple in design and easy to query; however dimensions might have duplicates and they get larger over time. Commonly used and ideal for BI tools like Tableau. Star Schema will be used for this project

Snowflake Schema: Similar to the Star Schema in that the fact table is the central component of the model and is surrounded by dimensions; however, each dimension is broken out further into sub-dimensions. This extension into sub-dimensions results in a snowflake pattern when depicted visually. Hence...the name “Snowflake Schema”. The Snowflake Schema is more complex, requires deeper knowledge of the underlying data sets and requires more effort to query. The main advantage comes through normalization from breaking redundancies into smaller tables which optimizes storage.

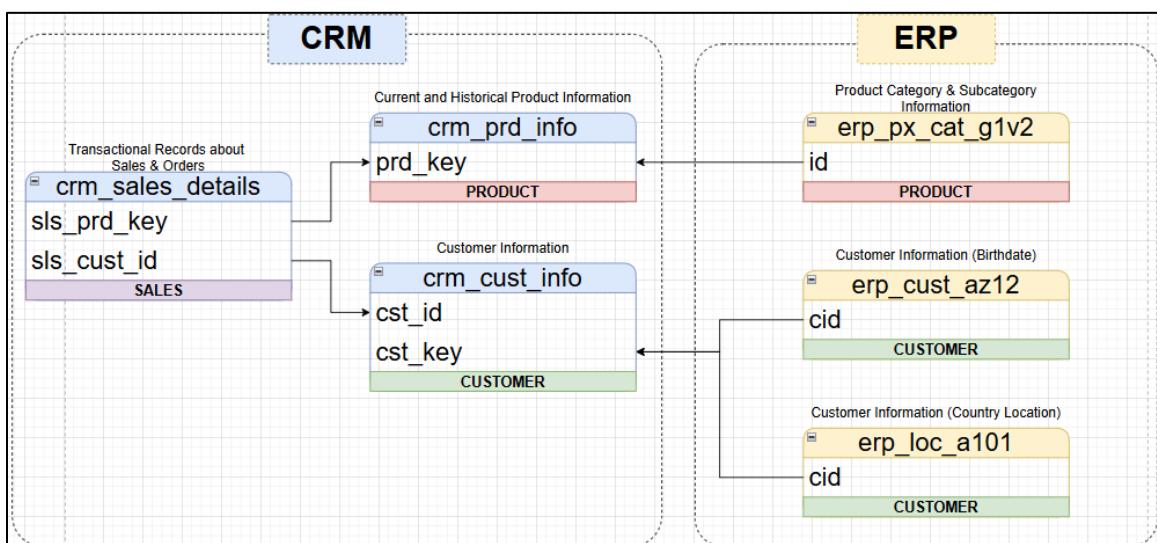


Data Model Design & Build

The dimensions and facts need to be defined now that Star Schema has been selected. Dimension tables will contain descriptive information to give context to the data (example: product info). If the data is oriented around who, what, when and where...then it belongs in the dimension table. Fact tables contain quantitative information that represent events / transactions (Example: sales orders, shipments, returns with dates, and other numerical data). If the data is oriented around how much, when, etc....then it belongs in the Fact table.

The first step here is to understand the business context of the data to ensure data is correctly identified / assigned to dimension tables and fact table. In a real-world situation, I would have met with the business stakeholders / process SMEs / process owners to walk-through the data and its alignment to relevant business processes.

Before doing so, it is important to explore the data and develop an initial understanding and documenting it to enable review and validation with the business SMEs. As this is a fictional situation without stakeholders to review, I merely prepared the below visual representation of the data objects and their relationships that will serve as the foundation for building out the Star Schema for the Gold Layer.



b. Coding & Validating: Data Integration & Quality Checks

Build of the Gold Layer is focused on creating the Star Schema in the data warehouse with dimensions and facts defined to build Views. This will be done by integrating Silver Layer data into business-friendly objects to enable data consumption for reporting & analytics.

Integrating Customer Data

There are three tables that need to be integrated to bring customer data together into a single view in the Gold layer from the Silver layer. Those are:

- silver.crm_cust_info
- silver.erp_cust_az12
- silver.erp_loc_a101

I created a join of these three tables with 'silver.crm_cust_info' serving as the master table.

Once created, I ran the query and ran a check to ensure there are no duplicates in the primary keys. If duplicates appear as a result of the check, then that is an indication that there is an issue with the Silver layer data that would need to be resolved before moving forward with the Gold Layer.

Initial Query to Join Customer Data

```

SQLQuery41...lndley (51)*  □ X
1   --Integrate Customer Data From Silver Layer per the data model
2   ✓ SELECT
3       ci.cst_id,
4       ci.cst_key,
5       ci.cst_firstname,
6       ci.cst_lastname,
7       ci.cst_marital_status,
8       ci.cst_gndr,
9       ci.cst_create_date,
10      ca.bdate,
11      ca.gen,
12      la.cntry
13  ✓ FROM silver.crm_cust_info ci -- silver.crm_cust_info serves as the master table to which the other customer tables will be joined
14      -- "ci" is the alias being assigned to this table. To enable use in the query
15  ✓ LEFT JOIN silver.erp_cust_az12 ca -- "ca" is the alias being assigned to this table. To enable use in the query
16      ON ci.cst_key = ca.cid -- joining on these two fields in alignment with the data model and data clean-up work that was done in the Silver layer
17  ✓ LEFT JOIN silver.erp_loc_az01 la -- "la" is the alias being assigned to this table. To enable use in the query
18      ON ci.cst_key = la.cid -- joining on these two fields in alignment with the data model and data clean-up work that was done in the Silver layer

```

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date	bdate	gen	cntry
29469	AW00029469	Dominique	Saunders	Married	Female	2026-01-25	1977-10-20	n/a	United King...
29470	AW00029470	Nathan	Roberts	Single	n/a	2026-01-25	1971-08-22	n/a	France
29471	AW00029471	Dana	Oteaga	Single	n/a	2026-01-25	1965-09-23	n/a	France
29472	AW00029472	Lacey	Shama	Married	n/a	2026-01-25	1965-09-11	n/a	France
29473	AW00029473	Carmen	Subram	Single	n/a	2026-01-26	1965-11-19	n/a	United King...
29474	AW00029474	Jane	Raje	Married	n/a	2026-01-25	1965-04-02	n/a	Germany
29475	AW00029475	Jaed	Ward	Single	n/a	2026-01-25	1976-03-18	n/a	Germany
29476	AW00029476	Elizabeth	Bradley	Married	n/a	2026-01-25	1964-12-30	n/a	Germany
29477	AW00029477	Neil	Ruiz	Married	n/a	2026-01-25	1965-01-02	n/a	United King...
29478	AW00029478	Daren	Carlson	Single	n/a	2026-01-25	1964-11-21	n/a	United King...
29479	AW00029479	Tomas	Tapo	Married	n/a	2026-01-25	1965-06-30	n/a	France

Query & Results Checking for Duplicate Primary Keys

```

SQLQuery41...lndley (51)*  □ X
1   --Integrate Customer Data From Silver Layer per the data model
2   ✓ SELECT cst_id, COUNT(*) FROM
3       (SELECT
4           ci.cst_id,
5           ci.cst_key,
6           ci.cst_firstname,
7           ci.cst_lastname,
8           ci.cst_marital_status,
9           ci.cst_gndr,
10          ci.cst_create_date,
11          ca.bdate,
12          ca.gen,
13          la.cntry
14      ✓ FROM silver.crm_cust_info ci -- silver.crm_cust_info serves as the master table to which the other customer tables will be joined
15      -- "ci" is the alias being assigned to this table. To enable use in the query
16  ✓ LEFT JOIN silver.erp_cust_az12 ca -- "ca" is the alias being assigned to this table. To enable use in the query
17      ON ci.cst_key = ca.cid -- joining on these two fields in alignment with the data model and data clean-up work that was done in the Silver layer
18  ✓ LEFT JOIN silver.erp_loc_az01 la -- "la" is the alias being assigned to this table. To enable use in the query
19      ON ci.cst_key = la.cid -- joining on these two fields in alignment with the data model and data clean-up work that was done in the Silver layer
20
21  ✓ GROUP BY cst_id
22  ✓ HAVING COUNT(*) > 1

```

cst_id	(No column name)
	1

The fields that contain customer gender in tables 'crm_cust_info' and 'silver.erp_cust_az12' were not cleansed for consistency during the build of the Silver layer. I would think that the correct approach would be to go back to the silver layer and make the corrections there, but....as I am merely following along with the instructor, I will cleanse them for consistency in the gold layer as he does.

Below is a view of the inconsistency that needs to be corrected. In a real world situation, I would have the business stakeholder / data owner confirm the system that serves as the source of truth for customer gender. In this fictitious situation, the crm is the source of truth. So, I will need to create a query that applies consistency and build it into the primary query for customer data.

The query written to apply consistency to customer gender data will leverage a CASE WHEN to use the value in silver.crm_cust_info if the value is not 'n/a'. If the value is 'n/a', the it will pull the gender from silver.erp_cust_az12. If it is n/a in both, then the value remains as 'n/a'.

The query built to apply consistency was added to the customer info data integration query. I then applied business-friendly names to each field, reordered the columns to improve readability.

The integrated customer info is a dimension as it contains descriptive data rather than events, transactions, or any measurements.

As dimension object, it will need a primary key defined for each record. In this project, a surrogate key will be created using a Window function (Row_Number) so that there is no dependency on use of the primary keys of integrated tables.

Surrogate Key = System-generated unique identifier assigned to each record in a table. They can be defined in the ddl or they can be query based using a Window function (Row_Number)

Query Identifying Inconsistency in Gender Identifiers

The screenshot shows two tabs in SSMS: 'SQLQuery42.s...indsley (61)*' and 'SQLQuery41.sq...indsley (51)*'. The code in the first tab is:

```

1  SELECT DISTINCT
2      ci.cst_gndr,
3      ca.gen
4  FROM silver.crm_cust_info ci -- silver.crm_cu
-- "ci" is t
5      LEFT JOIN silver.erp_cust_az12 ca -- "ca" is
6          ON ci.cst_key = ca.cid -- joining on t
7      LEFT JOIN silver.erp_loc_a101 la -- "la" is t
8          ON ci.cst_key = la.cid
9
10     ORDER BY 1, 2

```

The results grid shows the following data:

	cst_gndr	gen
1	Female	Female
2	Female	Male
3	Female	n/a
4	Male	Female
5	Male	Male
6	Male	n/a
7	n/a	NULL
8	n/a	Female
9	n/a	Male
10	n/a	n/a

Query Applying Consistency to Gender Data

The screenshot shows two tabs in SSMS: 'SQLQuery42.s...indsley (61)*' and 'SQLQuery41.sq...indsley (51)*'. The code in the first tab is:

```

1  SELECT DISTINCT
2      ci.cst_gndr,
3      ca.gen
4      CASE WHEN ci.cst_gndr != 'n/a' THEN ci.cst_gndr -- CRM is master
5          ELSE COALESCE(ca.gen, 'n/a') END AS new_gen
6
7  FROM silver.crm_cust_info ci -- silver.crm_cust_info serves as the
-- "ci" is the alias being assign
8      LEFT JOIN silver.erp_cust_az12 ca -- "ca" is the alias being assign
9          ON ci.cst_key = ca.cid -- joining on these two fields in al
10     LEFT JOIN silver.erp_loc_a101 la -- "la" is the alias being assign
11         ON ci.cst_key = la.cid
12
13     ORDER BY 1, 2

```

The results grid shows the following data:

	cst_gndr	gen	new_gen
1	Female	Female	Female
2	Female	Male	Female
3	Female	n/a	Female
4	Male	Female	Male
5	Male	Male	Male
6	Male	n/a	Male
7	n/a	NULL	n/a
8	n/a	Female	Female
9	n/a	Male	Male
10	n/a	n/a	n/a

CREATE CUSTOMER VIEW

The screenshot shows the Object Explorer on the left with a tree view of database objects. A red box highlights the 'gold.dim_customer' view under the 'DataWarehouse' database. The main pane displays the T-SQL code for creating the view:

```
CREATE VIEW gold.dim_customer AS
SELECT
    ROW_NUMBER() OVER (ORDER BY cst_id) AS customer_key, --Assigns surrogate key (primary key to all records)
    ci.cst_id AS customer_id,
    ci.cst_key AS customer_number,
    ci.cst_firstname AS first_name,
    ci.cst_lastname AS last_name,
    la.ctry AS country,
    ci.cst_marital_status AS marital_status,
    CASE WHEN ci.cst_gndr != 'n/a' THEN ci.cst_gndr -- CRM is master for gender info
        ELSE COALESCE(ca.gen, 'n/a')
    END AS gender,
    ca.bdate AS birthdate,
    ci.cst_create_date AS create_date

FROM silver.crm_cust_info ci -- silver.crm_cust_info serves as the master table to which the other customer tab
-- "ci" is the alias being assigned to this table. To enable use in the query
LEFT JOIN silver.erp_cust_ar12 ca -- "ca" is the alias being assigned to this table. To enable use in the quer
ON ci.cst_key = ca.cid -- joining on these two fields in alignment with the data model and data clean-up
LEFT JOIN silver.erp_loc_ar01 la -- "la" is the alias being assigned to this table. To enable use in the quer
ON ci.cst_key = la.cid -- joining on these two fields in alignment with the data model and data clean-up
```

The status bar at the bottom indicates 'Commands completed successfully.' and 'Completion time: 2020-12-01T09:07:41.0146617-05:00'.

CREATE CUSTOMER VIEW

The final step in the creation of the Customer View is to conduct a quality check of the data. To do so, I ran a 'SELECT *' query and conducted a quick scan to ensure all data appeared in the correct order and that there are no unexpected NULL values. I also ran a 'SELECT DISTINCT' query on the 'gender' field to ensure that there are no values in the field other than 'Male', 'Female', or 'n/a'. Results of both queries indicated success.

Integrating Product Data

1. Product data exists in both the CRM and the ERP data sources. Integrating these into a single, usable view is necessary to enable analytics on the data.

The product info table contains historical and current product information. In a real-world scenario, a decision would be needed as to whether the view must contain both historical and current product information for analysis. In this fictitious case, only the current product data will be used in the view. The product key is being used in the mapping of product data between tables. This also allows for the historical data to be filtered out so only current product data is included.

To filter out historical data, a WHERE condition was applied to the query on the prd_end_dt field to filter out records where the value != NULL. All records with NULL in the prd_end_dt field are active products because the product end date has not occurred.

Query Joining Product Data from CRM and ERP tables

The screenshot shows a SQL Server Management Studio window with three tabs: 'SQLQuery43.s...lindsay (62)*', 'SQLQuery42.sq...lindsay (61)*', and 'gold_layer_cre...lindsay (51)'. The main query window contains the following code:

```

1  SELECT
2      pn_prd_id,
3      pn_cat_id,
4      pn_prd_key,
5      pn_prd_nm,
6      pn_prd_cost,
7      pn_prd_line,
8      pn_prd_start_dt,
9      pc_cat,
10     pc_subcat,
11     pc_maintenance
12  FROM silver.crm_prd_info pn -- silver.crm_prd_info serves as the master table to which the ERP product table will be joined
13      -- "pn" is the alias being assigned to this table. To enable use in the query
14  LEFT JOIN silver.erp_px_cat_glv2 pc -- "pc" is the alias being assigned to this table. To enable use in the query
15  ON pn.cat_id = pc.id
16  WHERE prd_end_dt IS NULL -- Filter out all historical data

```

The results grid displays 15 rows of product data, including columns like prd_id, cat_id, prd_key, prd_nm, prd_cost, prd_line, prd_start_dt, cat, subcat, and maintenance.

Conduct Data Quality Checks

The check for data quality, I wrote a `SELECT COUNT(*)` query with grouping by `prd_key` for records with `prd_key` appearing more than once. There were no records displayed in results the quality check query. That tells me that:

- There are no duplicates resulting from the join
- Each product has only one record
- There are no historical data

The screenshot shows a SQL Server Management Studio window with three tabs: 'SQLQuery43.s...lindsay (62)*', 'SQLQuery42.sq...lindsay (61)*', and 'gold_layer_cre...lindsay (51)'. The main query window contains the following code:

```

2  SELECT prd_key, COUNT(*) FROM (
3      SELECT
4          pn_prd_id,
5          pn_cat_id,
6          pn_prd_key,
7          pn_prd_nm,
8          pn_prd_cost,
9          pn_prd_line,
10         pn_prd_start_dt,
11         pc_cat,
12         pc_subcat,
13         pc_maintenance
14    FROM silver.crm_prd_info pn -- silver.crm_prd_info serves as the master table to which the ERP product table will be joined
15        -- "pn" is the alias being assigned to this table. To enable use in the query
16    LEFT JOIN silver.erp_px_cat_glv2 pc -- "pc" is the alias being assigned to this table. To enable use in the query
17    ON pn.cat_id = pc.id
18    WHERE prd_end_dt IS NULL -- Filter out all historical data
19 )t
20 GROUP BY prd_key HAVING COUNT(*) > 1

```

The results grid shows a single row with the column `prd_key` and value `(No column name)`.

The field names were then assigned business-friendly names and columns were reordered for readability. Additionally a the query was enhanced to include the creation of a surrogate key (primary key) for each product using the `ROW_NUMBER` window function (similar to what was done for the customer info dimension).

CREATE VIEW was added to the query, then the query was executed...resulting in the creation of the view in the data warehouse.

The screenshot shows the SQL Server Management Studio interface. The Object Explorer on the left lists databases, including 'PF32ZQ8SQLEXPRESS' (SQL Server 16.0) which contains 'gold.dim_products'. The query editor window on the right displays a T-SQL CREATE VIEW statement for 'gold.dim_products'. A red arrow points from the 'gold.dim_products' entry in Object Explorer to the corresponding table reference in the query editor.

```
CREATE VIEW gold.dim_products AS
BEGIN
    SELECT
        ROW_NUMBER() OVER (ORDER BY pn.prd_start_dt, pn.prd_key) AS product_key, --Create and assign primary key
        pn.prd_id AS product_id,
        pn.prd_key AS product_number,
        pn.prd_nm AS product_name,
        pn.cat_id AS category_id,
        pc.cat AS category,
        pc.subcat AS subcategory,
        pc.maintenance,
        pn.prd_cost AS cost,
        pn.prd_line AS product_line,
        pn.prd_start_dt AS start_date
    FROM silver.crm_prd_info pn -- silver.crm_prd_info serves as the master table to which the ERP product table will be joined
        -- "pn" is the alias being assigned to this table. To enable use in the query
    LEFT JOIN silver.erp_px_cat_glv2 pc -- "pc" is the alias being assigned to this table. To enable use in the query
        ON pn.cat_id = pc.id
    WHERE prd_end_dt IS NULL -- Filter out all historical data
END
```

Building the View of Sales Data as a Facts Table

What I did - connected Customer and Product dimensions (product ID, customer ID) to the sales data

DIM keys **DATES** **MEASURES**

```
--Create Sales Facts Table Joining Customer and Product Dimension to Sales Data
CREATE VIEW gold_fact_sales AS
SELECT
    sd.sls_ord_num AS order_number,
    pr.product_key,
    cu.customer_key,
    sd.sls_order_dt AS order_date,
    sd.sls_ship_dt AS shipping_date,
    sd.sls_due_dt AS due_date,
    sd.sls_sales AS sales_amount,
    sd.sls_quantity AS quantity,
    sd.sls_price AS price
FROM silver_crm_sales_details sd
LEFT JOIN gold.dim_product pr
ON sd.sls_prd_key = pr.product_number
LEFT JOIN gold.dim_customer cu
ON sd.sls_cust_id = cu.customer_id

--Create Sales Facts Table Joining Customer and Product Dimension to Sales Data
CREATE VIEW gold_fact_sales AS
SELECT
    sd.sls_ord_num AS order_number,
    pr.product_key,
    cu.customer_key,
    sd.sls_order_dt AS order_date,
    sd.sls_ship_dt AS shipping_date,
    sd.sls_due_dt AS due_date,
    sd.sls_sales AS sales_amount,
    sd.sls_quantity AS quantity,
    sd.sls_price AS price
FROM silver_crm_sales_details sd
LEFT JOIN gold.dim_product pr
ON sd.sls_prd_key = pr.product_number
LEFT JOIN gold.dim_customer cu
ON sd.sls_cust_id = cu.customer_id
```

Checking Foreign Key Integrity (dimensions) to see if all dimension tables can successfully join to the fact table: No bad results....good to go

```
SQLQuery46.s...indsley (63)* ✎ x SQLQuery45.sql...indsley (61)* gold_layer_cre...  
1 ---Foreign Key Integrity (Dimensions)  
2 SELECT * FROM gold_fact_sales f  
3 LEFT JOIN gold_dim_customer c  
4 ON c.customer_key = f.customer_key  
5 LEFT JOIN gold_dim_products p  
6 ON p.product_key = f.product_key  
7 WHERE p.product_key IS NULL  
  
90% ▾ 8 ✘ 0 ⌈ ⌋ ⌄ ⌅   
Results Messages  

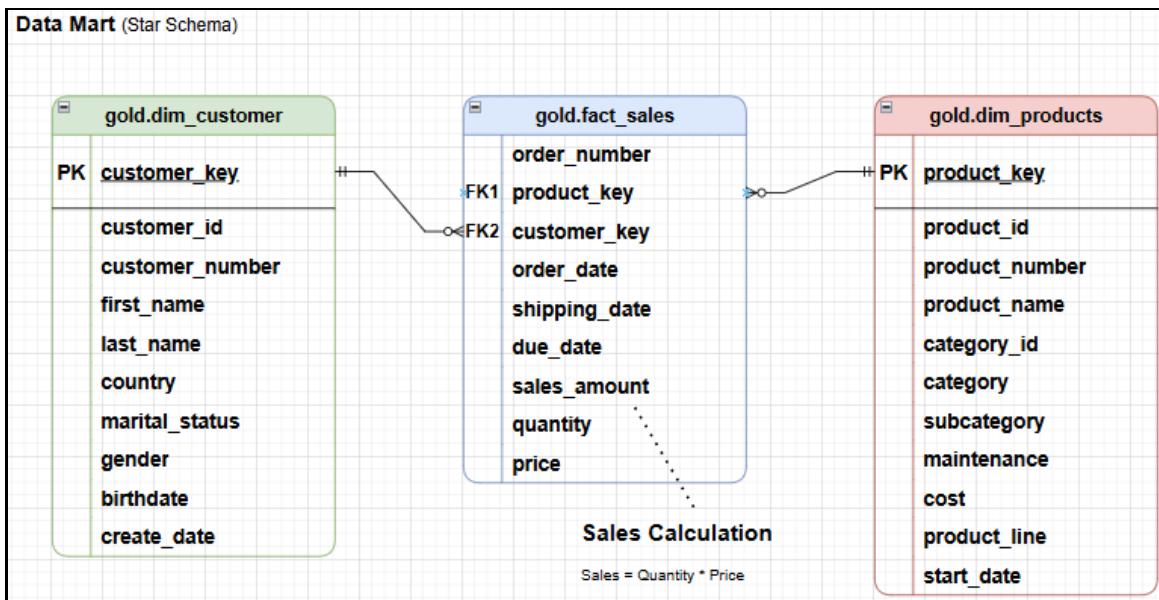

| order_number | product_key | customer_key | order_date | shipping_date | due_date | sales_ |
|--------------|-------------|--------------|------------|---------------|----------|--------|
|--------------|-------------|--------------|------------|---------------|----------|--------|


```

c. Document & Commit Code to Git

Document Star Schema

The Star Schema was modeled to visually depict the relationships defined between the customer dimension, the products dimension and the sales Facts views in the Gold Layer. It includes identification of primary keys, foreign keys, and the relationships (many to one / one to many)



Create Data Catalog & Document Versioning in Git

A Data Catalog is needed to support the data product we have created (the data warehouse) to be developed based on the data warehouse. The Data Catalog describes the data model (it's columns, tables, relationships, etc) to enable users to derive meaningful insights and to avoid repeat questions from users. The data catalog is documented in an md file in my Git repository.

The Data Catalog has a section for each table in the Gold layer (`gold.dim_customer`, `gold.dim_products`, and `gold.facts_sales`). Each section includes the purpose of the table, the column names, data types, and column descriptions (with examples).

I defined and documented the data catalog in a markdown file and committed it to my Git repo for this project under the 'docs' folder.

Data Dictionary for Gold Layer

Overview

The Gold layer is the business-level data representation, structure to support analytical and reporting use cases. It consists of dimension tables and fact tables for specific business metrics.

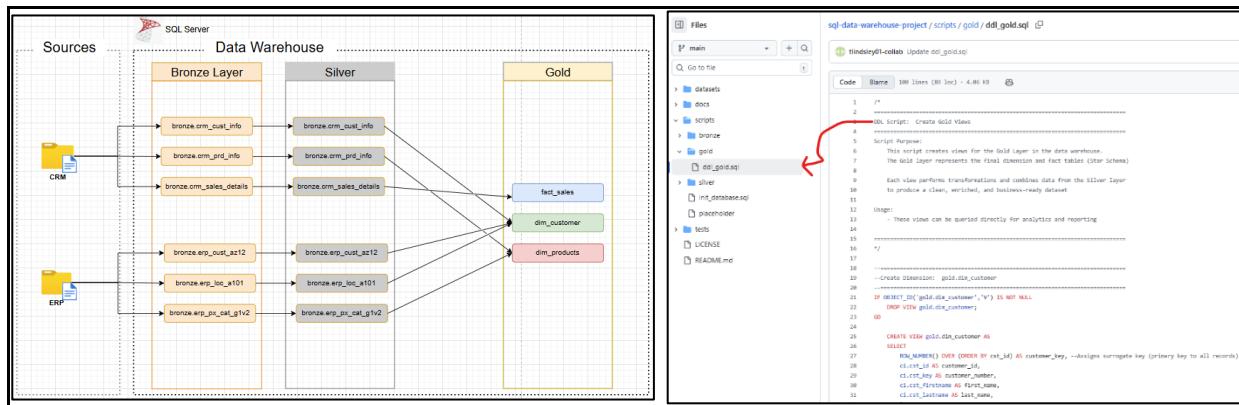
1. gold.dim_customer

- Purpose: Stores customer details enriched with demographic and geographic data
- Columns:

Column Name	Data Type	Description
customer_key	INT	Surrogate key uniquely identifying each customer record in the dimension table
customer_id	INT	Unique numerical identifier assigned to each customer
customer_number	NVARCHAR(50)	Alphanumeric identifier representing the customer, used for tracking and referencing
first_name	NVARCHAR(50)	The customer's first name as recorded in the CRM
last_name	NVARCHAR(50)	The customer's last name as recorded in the CRM
country	NVARCHAR(50)	The customer's country of residence (e.g., 'France')
marital_status	NVARCHAR(50)	The marital status of the customer (e.g. 'Married', 'Single')
gender	NVARCHAR(50)	The customer's gender (e.g., 'Male', 'Female', 'n/a')

Update the Data Flow Diagram to show Gold Layer and commit Gold Layer view ddl to Git: I updated the data flow diagram to include the dimension tables for customer and product data along with the fact table for sales data. This view provides a visual representation of the data lineage for the data warehouse.

I also combined all gold layer creation scripts into a single ddl sql and uploaded to git with the scripts written for quality checks in the gold layer and all final diagrams and data flows created throughout this project.



9. Exploratory SQL Data Analysis - Soon to Come

10. SQL Advanced Analytics - Soon to Come