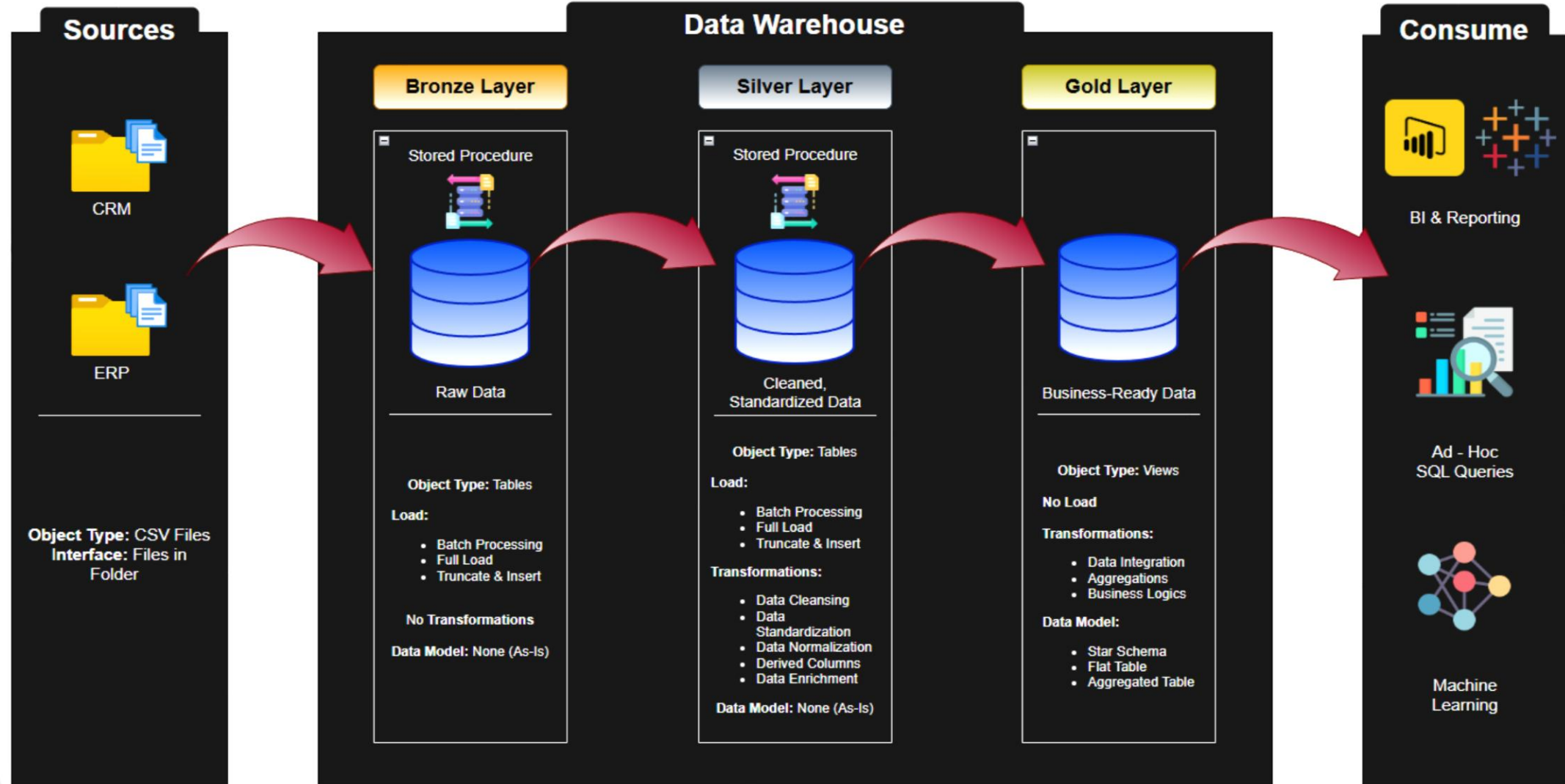


SQL Data Warehouse Project



DW High-Level Architecture



Database Layers Details

	Bronze Layer	Silver Layer	Gold Layer
Definition	Raw, unprocessed data as-is from sources	Cleand and Standardized data	Business-Ready Data
Objective	Traceability and Debugging	Intermediate Layer Prepare data for analysis	Provide data to be consumed for reporting and analytics
Objective Type	Tables	Tables	Views
Load Method	Full Load (Truncate and Insert)	Full Load (Truncate and Insert)	None
Data Transformation	None (As-Is)	<ul style="list-style-type: none"> • Data Cleaning • Data Standardization • Data Normalization • Derive Columns • Data Enrichment 	<ul style="list-style-type: none"> • Data Integration • Data Aggregation • Business Logic and Rules
Data Modeling	None (As-Is)	None (As-Is)	Star Schema Aggregated Objects Flat Tables
Target Audience	Data Engineers	Data Analysts Data Engineers	Data Analysts Business Users

Source System Interview

Business Context & Ownership

- . Who owns the data?
- . What business process it supports?
- . System and data documentation
- . Data model and data catalog

Architecture & Tech Stack

- . How is data stored (SQL Server, Oracle, AWS, Azure...)?
- . What are the integration capabilities?

Extract & Load

- . Incremental Vs Full Load?
- . Data Scope and historical needs
- . What are the expected size of the extracts?
- . Are there any data volume limitations?
- . How to avoid impacting the source system's performance?
- . Authentication and authorization (tokens, SSH keys, VPN, IP whitelisting)

Query	Query History
1	-- =====
2	-- Stored Procedure: load_bronze_data
3	-- =====
4	-- Purpose:
5	-- Loads data from CSV files into the bronze layer tables.
6	-- - Truncates tables before load.
7	-- - Uses COPY for bulk insert from CSV.
8	-- - Logs progress and duration.
9	-- =====
10	
11	CREATE OR REPLACE PROCEDURE bronze.load_bronze_data()
12	LANGUAGE plpgsql
13	AS \$\$
14	DECLARE
15	start_time TIMESTAMP;
16	end_time TIMESTAMP;
17	BEGIN
18	RAISE NOTICE '=====';
19	RAISE NOTICE 'Starting Bronze Layer Data Load';
20	RAISE NOTICE '=====';
21	
22	-- CRM TABLES -----
23	RAISE NOTICE '-----';
24	RAISE NOTICE 'Loading CRM Tables';
25	RAISE NOTICE '-----';
26	
27	-- bronze.crm_cust_info
28	start_time := clock_timestamp();
29	RAISE NOTICE '>> Truncating bronze.crm_cust_info...';
30	TRUNCATE TABLE bronze.crm_cust_info;
31	RAISE NOTICE '>> Loading data from /data/csvs/cust_info.csv...';
32	COPY bronze.crm_cust_info
33	FROM '/data/csvs/cust_info.csv'
34	WITH (FORMAT csv, HEADER true);
35	end_time := clock_timestamp();
36	RAISE NOTICE '>> Load duration: % seconds', EXTRACT(epoch FROM (end_time - start_time));
37	
38	-- bronze.crm_prd_info
39	start_time := clock_timestamp();
40	RAISE NOTICE '>> Truncating bronze.crm_prd_info...';
41	TRUNCATE TABLE bronze.crm_prd_info;
42	RAISE NOTICE '>> Loading data from /data/csvs/prd_info.csv...';
43	COPY bronze.crm_prd_info
44	FROM '/data/csvs/prd_info.csv'
45	WITH (FORMAT csv, HEADER true);
46	end_time := clock_timestamp();
47	RAISE NOTICE '>> Load duration: % seconds', EXTRACT(epoch FROM (end_time - start_time));
48	
49	-- bronze.crm sales details

Query	Query History
12	DROP TABLE IF EXISTS silver.crm_cust_info;
13	
14	CREATE TABLE silver.crm_cust_info (
15	cst_id INTEGER,
16	cst_key VARCHAR(50),
17	cst_firstname VARCHAR(50),
18	cst_lastname VARCHAR(50),
19	cst_marital_status VARCHAR(50),
20	cst_gndr VARCHAR(50),
21	cst_create_date DATE,
22	dwh_create_date TIMESTAMP DEFAULT CURRENT_TIMESTAMP
23);
24	
25	-- 2. Drop and create silver.crm_prd_info
26	DROP TABLE IF EXISTS silver.crm_prd_info;
27	
28	CREATE TABLE silver.crm_prd_info (
29	prd_id INTEGER,
30	cat_id VARCHAR(50),
31	prd_key VARCHAR(50),
32	prd_nm VARCHAR(50),
33	prd_cost INTEGER,
34	prd_line VARCHAR(50),
35	prd_start_dt DATE,
36	prd_end_dt DATE,
37	dwh_create_date TIMESTAMP DEFAULT CURRENT_TIMESTAMP
38);
39	
40	-- 3. Drop and create silver.crm_sales_details
41	DROP TABLE IF EXISTS silver.crm_sales_details;
42	
43	CREATE TABLE silver.crm_sales_details (
44	sls_ord_num VARCHAR(50),
45	sls_prd_key VARCHAR(50),
46	sls_cust_id INTEGER,
47	sls_order_dt DATE,
48	sls_ship_dt DATE,
49	sls_due_dt DATE,
50	sls_sales INTEGER,
51	sls_quantity INTEGER,
52	sls_price INTEGER,
53	dwh_create_date TIMESTAMP DEFAULT CURRENT_TIMESTAMP
54);
55	
56	-- 4. Drop and create silver.erp_loc_a101
57	DROP TABLE IF EXISTS silver.erp_loc_a101;
58	
59	CREATE TABLE silver.erp_loc_a101 (
60	cid VARCHAR(50).

```

Query Query History
21 CREATE OR REPLACE PROCEDURE silver.load_silver()
22 LANGUAGE plpgsql
23 AS $$
24 DECLARE
25     start_time TIMESTAMP;
26     end_time TIMESTAMP;
27     batch_start_time TIMESTAMP;
28     batch_end_time TIMESTAMP;
29 BEGIN
30     batch_start_time := NOW();
31     RAISE NOTICE '=====';
32     RAISE NOTICE 'Loading Silver Layer';
33     RAISE NOTICE '=====';
34
35     RAISE NOTICE '-----';
36     RAISE NOTICE 'Loading CRM Tables';
37     RAISE NOTICE '-----';
38
39     -- 1. Loading silver.crm_cust_info
40     start_time := NOW();
41     RAISE NOTICE '>> Truncating Table: silver.crm_cust_info';
42     TRUNCATE TABLE silver.crm_cust_info;
43     RAISE NOTICE '>> Inserting Data Into: silver.crm_cust_info';
44     INSERT INTO silver.crm_cust_info (
45         cst_id,
46         cst_key,
47         cst_firstname,
48         cst_lastname,
49         cst_marital_status,
50         cst_gndr,
51         cst_create_date
52     )
53     SELECT
54         cst_id,
55         cst_key,
56         TRIM(cst_firstname) AS cst_firstname,
57         TRIM(cst_lastname) AS cst_lastname,
58         CASE
59             WHEN UPPER(TRIM(cst_marital_status)) = 'S' THEN 'Single'
60             WHEN UPPER(TRIM(cst_marital_status)) = 'M' THEN 'Married'
61             ELSE 'n/a'
62         END AS cst_marital_status, -- Normalize marital status values to readable format
63         CASE
64             WHEN UPPER(TRIM(cst_gndr)) = 'F' THEN 'Female'
65             WHEN UPPER(TRIM(cst_gndr)) = 'M' THEN 'Male'
66             ELSE 'n/a'
67         END AS cst_gndr, -- Normalize gender values to readable format
68         cst_create_date
69     FROM (

```

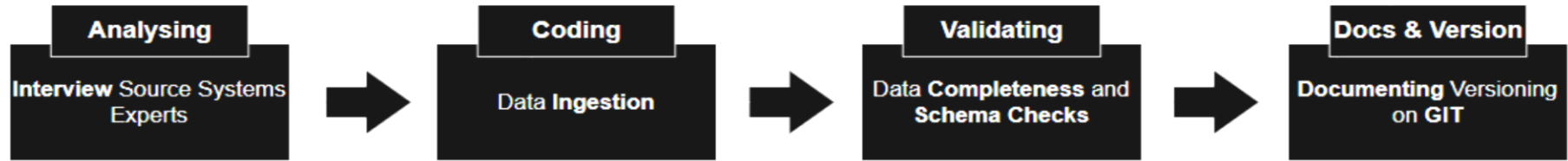
```

Query Query History
17 -- 1. CREATE DIMENSION: gold.dim_customers.
18 -- =====
19
20 DROP VIEW IF EXISTS gold.dim_customers;
21
22 CREATE VIEW gold.dim_customers AS
23 SELECT
24     ROW_NUMBER() OVER (ORDER BY cst_id) AS customer_key, -- Surrogate key
25     ci.cst_id AS customer_id,
26     ci.cst_key AS customer_number,
27     ci.cst_firstname AS first_name,
28     ci.cst_lastname AS last_name,
29     la.centry AS country,
30     ci.cst_marital_status AS marital_status,
31     CASE
32         WHEN ci.cst_gndr != 'n/a' THEN ci.cst_gndr
33         ELSE COALESCE(ca.gen, 'n/a')
34     END AS gender,
35     ca.bdate AS birthdate,
36     ci.cst_create_date AS create_date
37 FROM silver.crm_cust_info AS ci
38 LEFT JOIN silver.erp_cust_az12 AS ca
39 ON ci.cst_key = ca.cid
40 LEFT JOIN silver.erp_loc_a101 AS la
41 ON ci.cst_key = la.cid;
42
43 -- =====
44 -- 2. CREATE DIMENSION: gold.dim_products.
45 -- =====
46
47 DROP VIEW IF EXISTS gold.dim_products;
48
49 CREATE VIEW gold.dim_products AS
50 SELECT
51     ROW_NUMBER() OVER (ORDER BY pn.prđ_start_dt, pn.prđ_key) AS product_key, -- Surrogate key
52     pn.prđ_id AS product_id,
53     pn.prđ_key AS product_number,
54     pn.prđ_nm AS product_name,
55     pn.cat_id AS category_id,
56     pc.cat AS category,
57     pc.subcat AS subcategory,
58     pc.maintenance,
59     pn.prđ_cost AS cost,
60     pn.prđ_line AS product_line,
61     pn.prđ_start_dt AS start_date
62 FROM silver.crm_prđ_info AS pn
63 LEFT JOIN silver.erp_px_cat_glv2 AS pc
64 ON pn.cat_id = pc.id
65 WHERE pn.prđ_end_dt IS NULL; -- Filter out all historical data

```

Layers Creation

Bronze Layer



Silver Layer

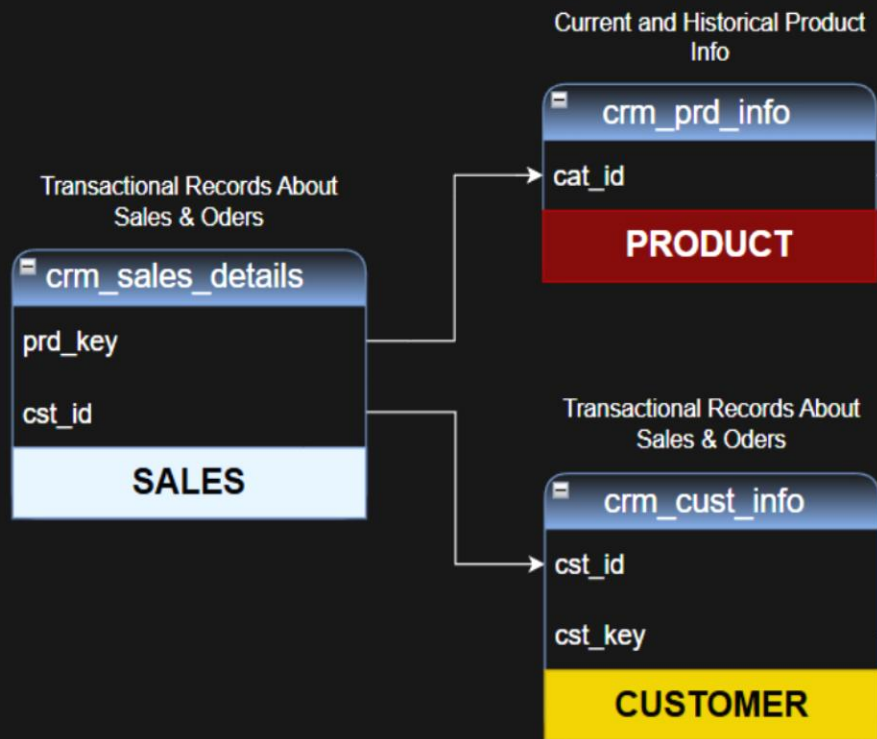


Gold Layer

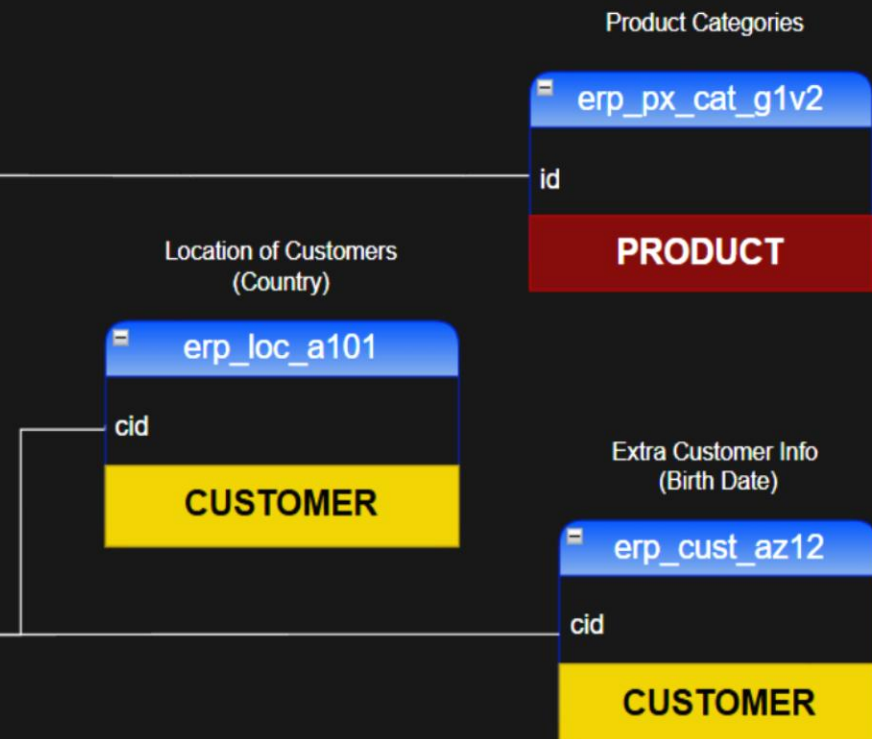


Integration Model

CRM

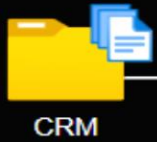


ERP

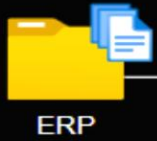


Data Flow Diagram

Sources



CRM



ERP

Bronze Layer

crm_sales_details

crm_cust_info

crm_prd_info

erp_cust_az12

erp_loc_a101

erp_px_cat_g1v2

Silver Layer

crm_sales_details

crm_cust_info

crm_prd_info

erp_cust_az12

erp_loc_a101

erp_px_cat_g1v2

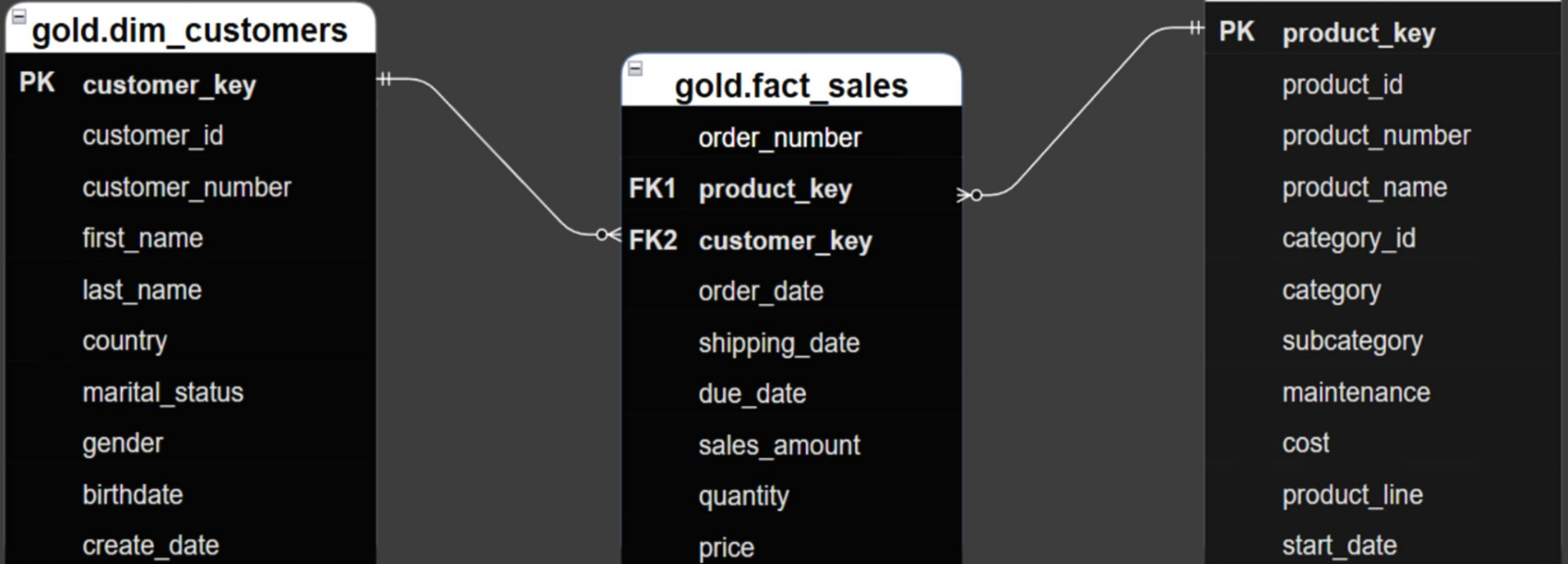
Gold Layer

fact_sales

dim_customers

dim_products

Sales Data Mart (Star Schema)











SQL Data Warehouse Project

Table

DWH Epics

Epics

Progress

 Requirements Analysis	100.00% <div><div></div></div>	
 Design Data Architecture	100.00% <div><div></div></div>	
 Project Execution	100.00% <div><div></div></div>	
 Build Bronze Layer	100.00% <div><div></div></div>	
 Build Silver Layer	100.00% <div><div></div></div>	
 Build Gold Layer	100.00% <div><div></div></div>	

+ New page

SQL Data Warehouse Project

1. Introduction

This project aims to build a modern Data Warehouse using SQL, incorporating ETL processes, data modeling, and analytics. It serves as a portfolio project, showcasing industry best practices in data engineering and analytics. The main goal is to centralize, standardize, and integrate data from multiple systems into a single repository, enabling reliable analysis and management reporting. Clear objectives and structured methodologies were defined to guide the work, ensuring organization and traceability.

- **Data Warehouse**
 - What is a Data Warehouse?
A centralized and structured repository that integrates data from multiple sources. Its main purpose is to provide a single, consistent view to support business analysis and decision-making.
- **ETL**
 - What is ETL (Extract, Transform, Load)?
A process that extracts data from different systems, transforms it into a consistent format, and loads it into the Data Warehouse.

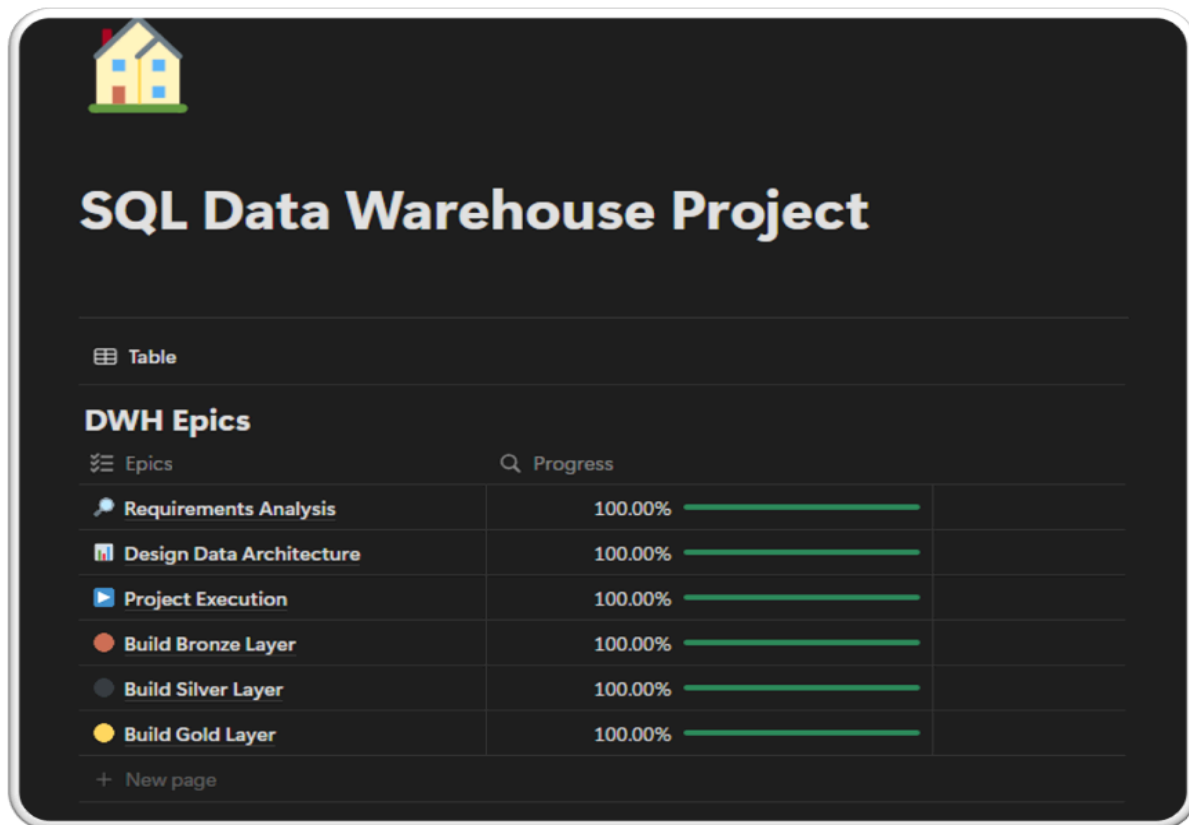
2. Initiation and Planning

We started by gathering all project requirements and defining the best approach to handle the complexities involved. We used Notion as a planning tool, documenting each step and creating a clear roadmap for the next phases. We also set up a GitHub repository for version control and documentation, and we aligned the architecture and essential conventions for building the Data Warehouse.

- **Requirements Gathering**
Defined the project scope and business needs.

- **Project Plan (Notion)**

Centralized all tasks, timelines, and responsibilities in a single platform.



- **Layered Approach**

Defined the layers of the DW: Bronze (raw data), Silver (cleaned and standardized data), and Gold (business model).

- **GIT Repository**

Created for version control and documentation:

- <https://github.com/JohnCustodio/sql-data-warehouse-project>

- **Architecture and Data Modeling**

We designed the high-level architecture of the DW in Draw.io and defined the data model for each layer.

- **Naming Conventions**

Established consistent naming rules for schemas, tables, views, and columns.

- **Database and Schemas**

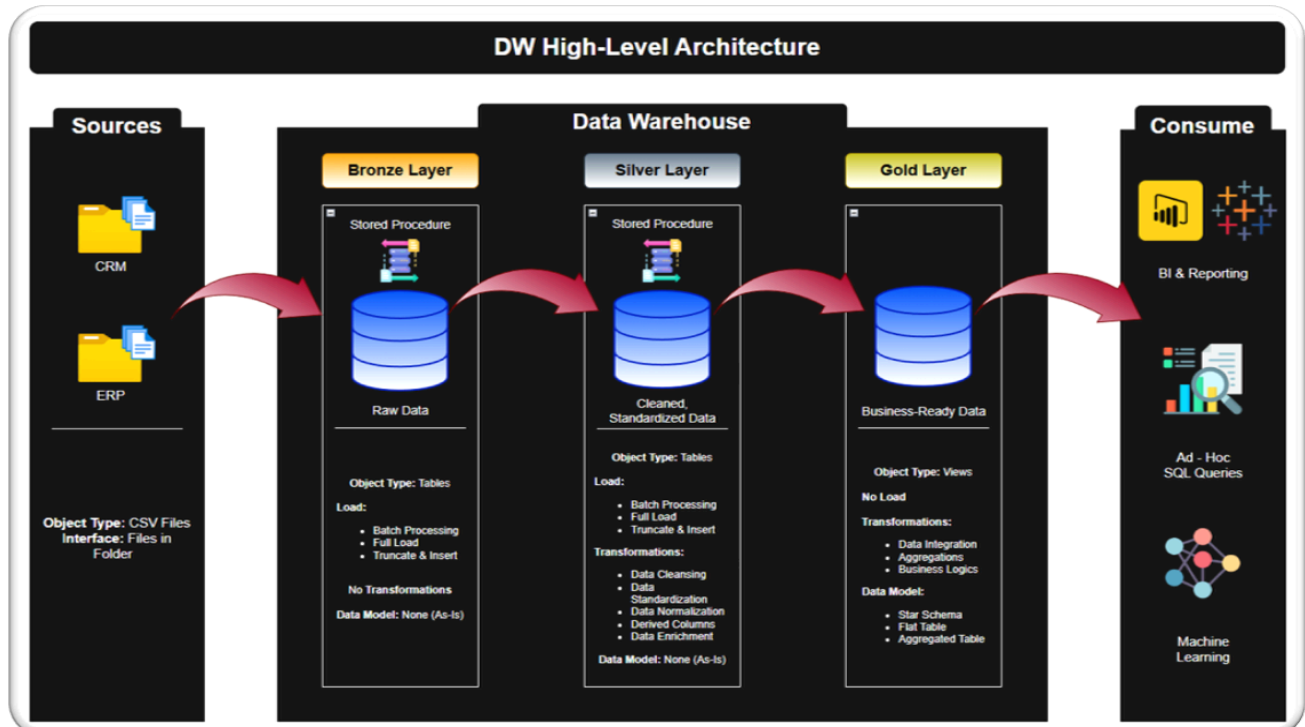
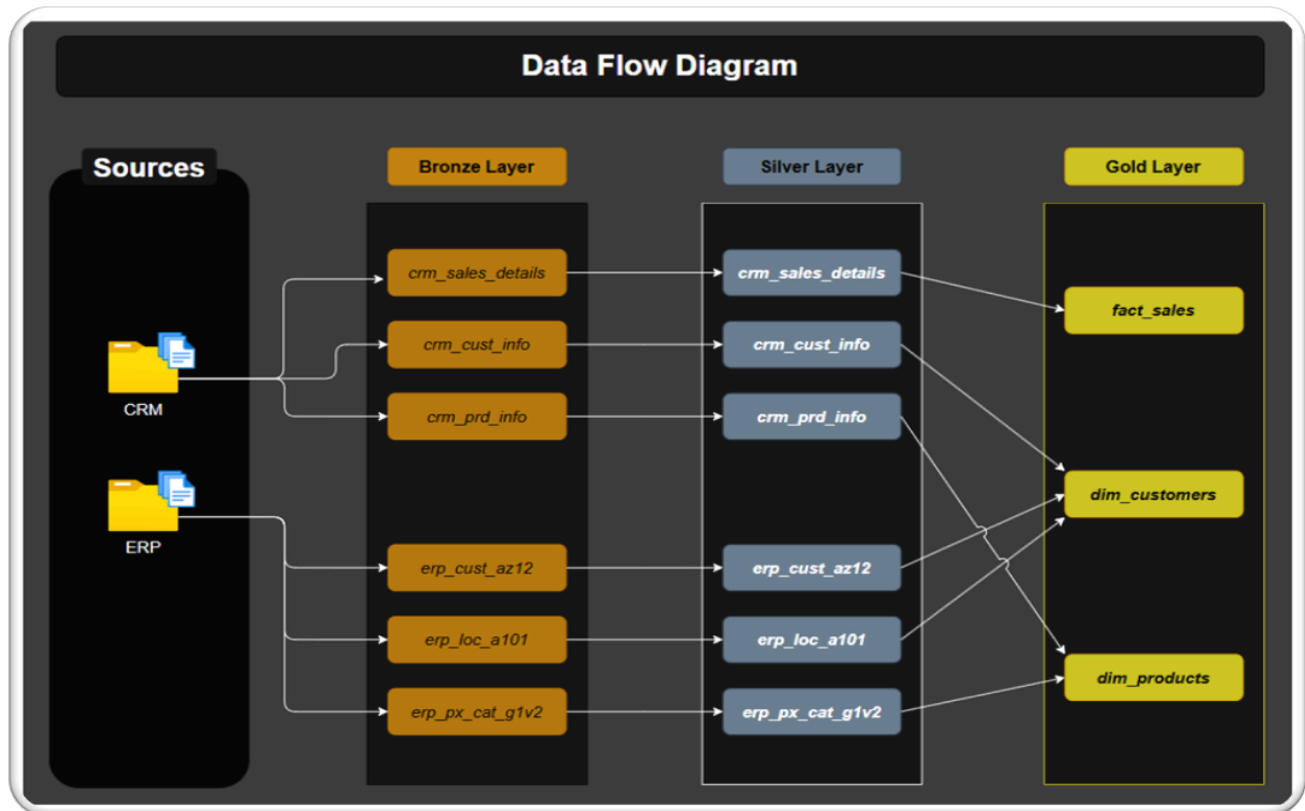
It was structured the PostgreSQL environment with separated schemas for each layer,

running the server in a Docker container, and accessing/manipulating data via PgAdmin.

3. Execution

The execution phase was organized into three main steps, each representing a DW layer. We followed an iterative and incremental approach, allowing for quick adjustments and partial deliveries that were continuously validated.

- **Bronze Layer**
 - Analysis of the data sources and creation of corresponding tables in PostgreSQL.
 - SQL scripts for initial data load, without transformations.
 - Creation of stored procedures to automate data loads.
 - Documented the data flow for this layer.
- **Silver Layer**
 - Understanding the data and creating the Silver tables.
 - Cleaning and standardizing data from CRM and ERP sources, including normalization, derivations, and enrichment.
 - Refined load scripts to remove inconsistencies.
 - Stored procedures for repetitive processes.
 - Documented the data flow for the Silver layer.
- **Gold Layer**
 - Applied data modeling concepts to meet business needs.
 - Created dimension tables (dim_customers, dim_products) and fact table (fact_sales).
 - Built the Star Schema model to facilitate analysis and reporting.
 - Created a Data Catalog and documented the final data flow.



4. Project Closure

We completed the project with thorough validations and reviews to ensure data quality and usability of the DW. All materials were documented in Notion and versioned in GitHub, forming a complete guide for future maintenance and enhancements.

- **Final Validation**
Checked data integrity and consistency by running SQL scripts for validation at each layer.
- **Complete Documentation**
All steps and deliverables were documented in Notion, and all scripts and documents were versioned in Git.
- **Next Steps**
This Data Warehouse will serve as a foundation for future exploratory and advanced data analysis projects.

Portuguese Version

Projeto de Data Warehouse SQL

1. Introdução

Este projeto tem como objetivo construir um Data Warehouse moderno utilizando SQL, incorporando processos de ETL, modelagem de dados e análises. Ele serve como um projeto de portfólio, demonstrando as melhores práticas de engenharia e análise de dados. O foco está em centralizar, padronizar e integrar dados de múltiplos sistemas em um único repositório, facilitando análises confiáveis e relatórios gerenciais. Foram definidos objetivos claros e metodologias que guiaram o trabalho, garantindo organização e rastreabilidade.

- **Data Warehouse**

- O que é um Data Warehouse?

Um repositório centralizado e estruturado de dados integrados de várias fontes. Seu objetivo é oferecer uma visão única e consistente para apoiar análises e decisões de negócio.

- **ETL**

- O que é ETL (Extract, Transform, Load)?

Um processo que extrai dados de diferentes sistemas, transforma-os em formatos consistentes e carrega-os no Data Warehouse.

2. Iniciação e Planejamento

Iniciamos com um levantamento completo dos requisitos do projeto e definimos a abordagem ideal para lidar com as complexidades envolvidas. Utilizamos o Notion como ferramenta de planejamento, registrando cada etapa e garantindo um caminho claro para as próximas fases. Também criamos um repositório no GitHub para controle de versões e documentação, e alinhamos a arquitetura e convenções essenciais para a criação do Data Warehouse.

- **Levantamento de Requisitos**

Definição do escopo e das necessidades do negócio.

- **Plano de Projeto (Notion)**

Organização de tarefas, cronogramas e responsáveis em um só lugar.

- **Abordagem de Camadas**

Definição das camadas do DW: Bronze (dados brutos), Silver (dados limpos e ajustados) e Gold (modelo de negócio).

- **Repositório GIT**

Criado para versionamento e documentação:

- <https://github.com/JohnCustodio/sql-data-warehouse-project>

- **Arquitetura e Modelagem**

Desenhamos a arquitetura de alto nível do DW no Draw.io e definimos a modelagem de dados para cada camada.

- **Naming Conventions**

Regras de nomenclatura para schemas, tabelas, views e colunas.

- **Banco de Dados e Schemas**

Estruturamos o PostgreSQL com schemas separados para cada camada, rodando o servidor em um container Docker, e acessando/manipulando os dados via PgAdmin.

3. Execução

A fase de execução foi dividida em três grandes etapas, cada uma representando uma camada do DW. Adotamos um desenvolvimento iterativo e incremental, permitindo ajustes rápidos e entregas parciais que foram validadas continuamente.

- **Camada Bronze**

- Análise das fontes de dados e criação de tabelas correspondentes no PostgreSQL.
- Scripts SQL para carga inicial, sem transformações.
- Criação de Stored Procedures para automatizar as cargas.
- Documentação do fluxo de dados dessa camada.

- **Camada Silver**

- Entendimento dos dados e criação de tabelas Silver.
- Limpeza e padronização dos dados das fontes CRM e ERP, incluindo normalização, derivações e enriquecimento.
- Scripts de carga refinados para remover inconsistências.
- Stored Procedures para processos repetitivos.
- Documentação do fluxo de dados Silver.

- **Camada Gold**

- Aplicação de conceitos de Data Modeling para atender às necessidades de negócio.
 - Criação de tabelas de dimensão (dim_customers, dim_products) e de fato (fact_sales).
 - Construção do modelo Star Schema para facilitar análises e relatórios.
 - Criação de um Catálogo de Dados e documentação do fluxo final.
-

4. Encerramento

Finalizamos com validações e revisões para garantir a qualidade dos dados e a usabilidade do DW. Todo o material gerado foi documentado via Notion e versionado no GitHub, formando um guia completo para futuras manutenções e evoluções.

- **Validação Final**

Conferimos a integridade e consistência dos dados, executando scripts SQL para validação em cada camada.

- **Documentação Completa**

Tudo registrado no Notion e scripts versionados no GIT.

- **Próximos Passos**

Este Data Warehouse servirá como base para futuros projetos de Análise Exploratória e Avançada de Dados.
