

# Project: Data Warehouse with PostgreSQL

## Table of Contents

Preparing the environment .....	1
Documentation .....	2
Building the documentation .....	2
Diagram support .....	2
OCR image text extraction support .....	2
Preparing the specifications table .....	3
Preparing the ETL mindmap .....	3
Data architecture .....	3
ETL task breakdown .....	4
Specifications .....	4
Data Flow .....	5
Pipeline .....	5
Bronze layer .....	7
Silver layer .....	7
Data integration analysis and sketch .....	7
Coding the transformations .....	8
CRM customer info table .....	8
CRM product info and sales tables .....	9
ERP tables .....	10
Gold layer .....	10

## Preparing the environment

This project uses Linux on Pop\_OS!, with PostgreSQL as the database management system for the data warehouse.

Before running the data warehouse scripts, ensure that you have PostgreSQL installed and running on your system. You will also need to have access to the PostgreSQL user with sufficient privileges to create databases and tables. Currently, the scripts assume that the PostgreSQL user is `postgres`.

To install PostgreSQL:

```
sudo apt install postgresql
```

By default, the service will start automatically after installation. If it does not, you can start it

manually using the following command:

```
sudo service postgresql start
```

## Documentation

The project documentation is written in AsciiDoc. If you want to modify and rebuild the documentation, you will need Asciidoctor. The method I used to install asciidoctor was to install the latest Ruby gems:

```
sudo apt install ruby  
sudo gem install asciidoctor asciidoctor-diagram asciidoctor-pdf
```

## Building the documentation

To build the documentation, navigate to the `src/docsrc` directory and run:

```
docbuild.sh all
```

## Diagram support

For diagrams, some additional packages are required.

```
sudo apt install graphviz plantuml
```

For mermaid, rather than use `mermaid-cli` from the package manager, I wanted the latest version of mermaid-cli via npm (after installing nvm for the latest version of nodejs).

```
curl -o- https://raw.githubusercontent.com/nvm-sh/nvm/v0.39.5/install.sh | bash  
source ~/.bashrc || source ~/.zshrc  
nvm install --lts  
  
npm install -g @mermaid-js/mermaid-cli
```

## OCR image text extraction support

For OCR image text extraction, Tesseract and the Python Pillow and pytesseract libraries are required.

```
sudo apt install tesseract-ocr
```

```
pip install Pillow pytesseract
```

## Preparing the specifications table

I decided to not make the extraction part of the pipeline and just prep the source beforehand. To recreate the table preparation process:

1. Extract the text:

```
pushd src/docsrc/img-extract
python ocr-extract.py > ocr-output.txt
popd
```

2. ....x.....

3. Okay, nevermind. ChatGPT did it all for me. I uploaded the image and gave it this prompt:

```
I will upload an image that is basically a fancy table.
Extract the text from the image and convert it to a table in AsciiDoc.
Let me know when you're ready.
```

But the OCR was a good exercise anyway.

## Preparing the ETL mindmap

Let's try ChatGPT again.

Well done, ChatGPT! Here's the prompt and you can find the source image in this repository under [src/docsrc/img-extract/etl-map.png](#).

```
New task: Now I'm going to upload an image of a fancy mindmap diagram.
I'd like you to extract the information from the image and recreate it
using mermaid syntax. Let me know when you're ready.
```

I'll clean it up with icons and shapes myself. Can't let AI rob me of all the fun!

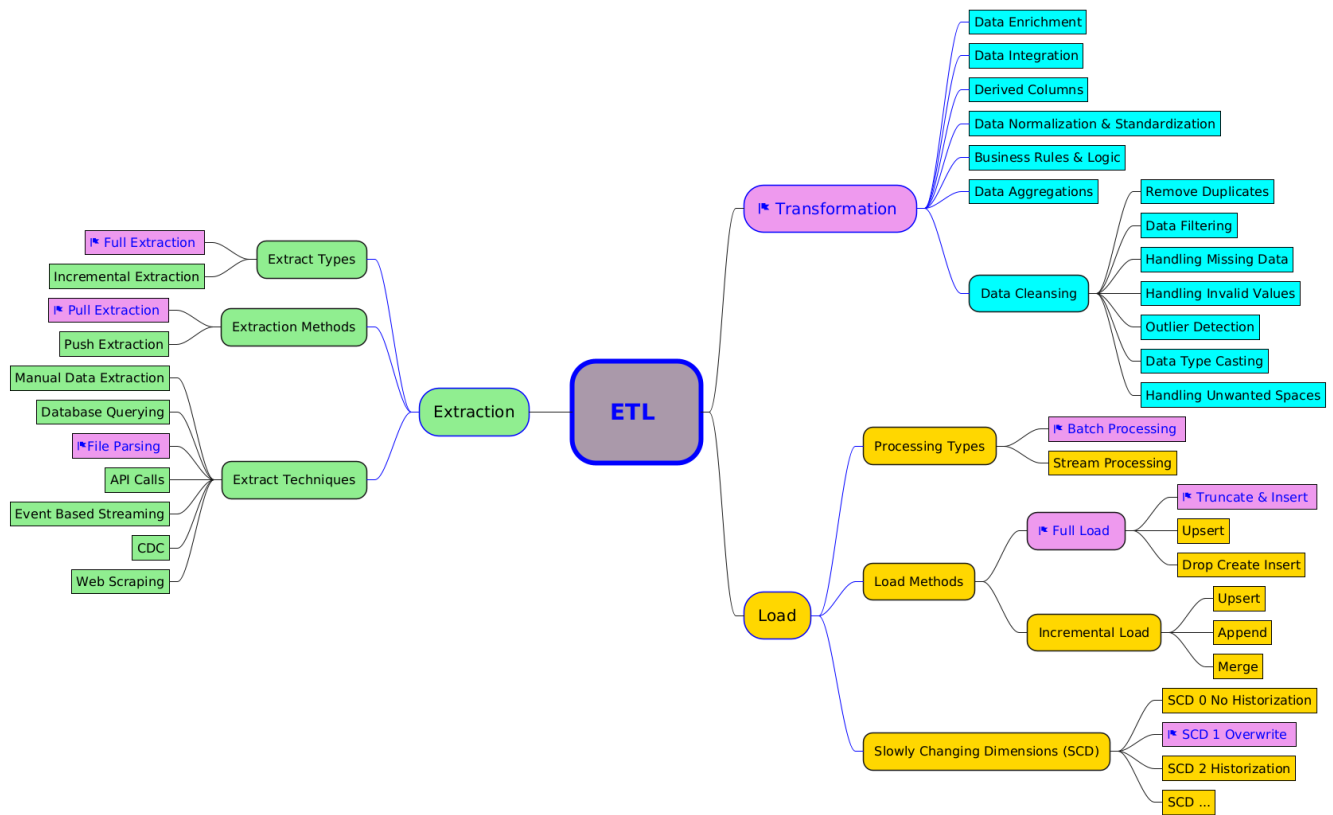
(I later asked ChatGPT to convert to a Drawio file instead, but it had an error about IDs. After debugging with [xmllint](#) and fixing the IDs, I got a valid Drawio file, but it was only two node levels deep and didn't follow the style anyway. So, I'm sticking with mermaid for now.)

# Data architecture

Before defining the data flow, it's essential to establish clear specifications for each layer in the data pipeline. This ensures that each layer serves its intended purpose and meets the requirements of the overall data strategy.

# ETL task breakdown

These are the tasks in the ETL process. I'd go over this with stakeholders to see what's required. For this exercise, Baraa wants us to do the tasks highlighted in pink. (We'll do the entire *Transformation* section and its subtasks.)



## Specifications

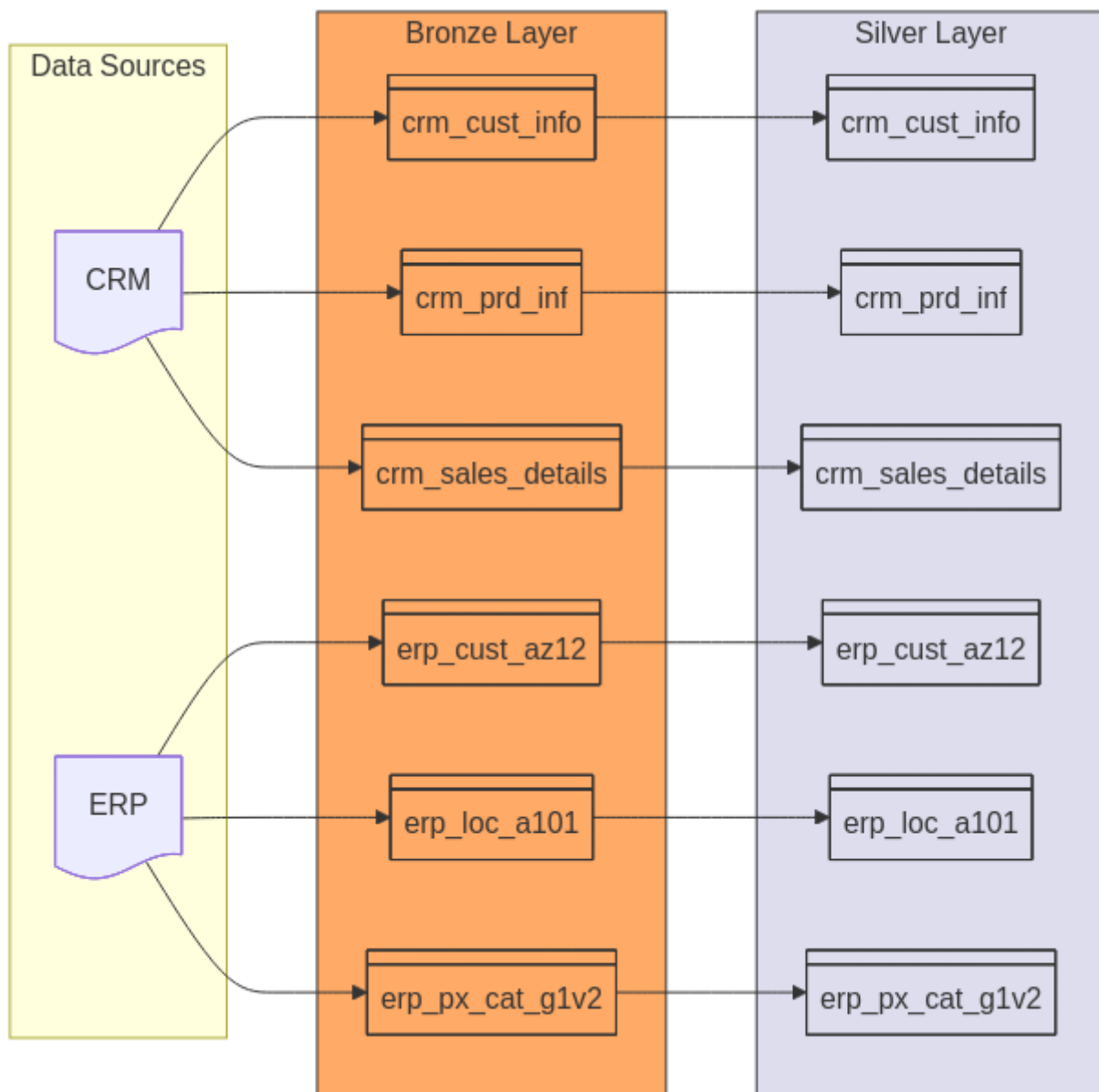
Given the chosen tasks in the ETL process, we have the following specifications.

Category	Bronze Layer	Silver Layer	Gold Layer
<b>Definition</b>	Raw, unprocessed data as-is from sources	Clean & standardized data	Business-Ready data
<b>Objective</b>	Traceability & Debugging	(Intermediate Layer) Prepare Data for Analysis	Provide data to be consumed for reporting & Analytics
<b>Object Type</b>	Tables	Tables	Views
<b>Load Method</b>	Full Load (Truncate & Insert)	Full Load (Truncate & Insert)	None
<b>Data Transformation</b>	None (as-is)	Data Cleaning Data Standardization Data Normalization Derived Columns Data Enrichment	Data Integration Data Aggregation Business Logic & Rules

Category	Bronze Layer	Silver Layer	Gold Layer
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

## Data Flow

Looking at the source data we have, this will be the flow to address our specifications.



## Pipeline

The ETL pipeline (`pipeline.sh`) is run with Bash, calling SQL scripts run with psql, with output to the command line and a log in `/tmp/de-dwh-sql`.

To run the full pipeline:

```
sudo -u postgres ./pipeline.sh full
```

There are other options besides 'full', if you want to run isolated stages of the pipeline. Here are the stages each option runs:

```
case $1 in
  "init")
    init_db
    ;;
  "bronze-create")
    bronze_create
    ;;
  "bronze-load")
    bronze_load
    ;;
  "bronze-all")
    bronze_create
    bronze_load
    ;;
  "silver-create")
    silver_create
    ;;
  "silver-load")
    silver_load
    ;;
  "silver-validate")
    silver_validate
    ;;
  "silver-all")
    silver_create
    silver_load
    silver_validate
    ;;
  "gold-create")
    gold_create
    ;;
  "gold-validate")
    gold_validate
    ;;
  "gold-all")
    gold_create
    gold_validate
    ;;
  "full")
    init_db
    bronze_create
    bronze_load
```

```

silver_create
silver_load
silver_validate
gold_create
gold_validate
;;
*) echo -ne "\nUsage: $0 {init\n |bronze-create|bronze-load|bronze-all\n
|silver-create|silver-load|silver-validate|silver-all\n |gold-create|gold-
validate|gold-all|full}\n"
esac

```

## Bronze layer

The Bronze layer is the initial storage area for raw data ingested from various sources. This layer is designed to store data in its original format, preserving its integrity and providing a foundation for further processing and transformation.

## Silver layer

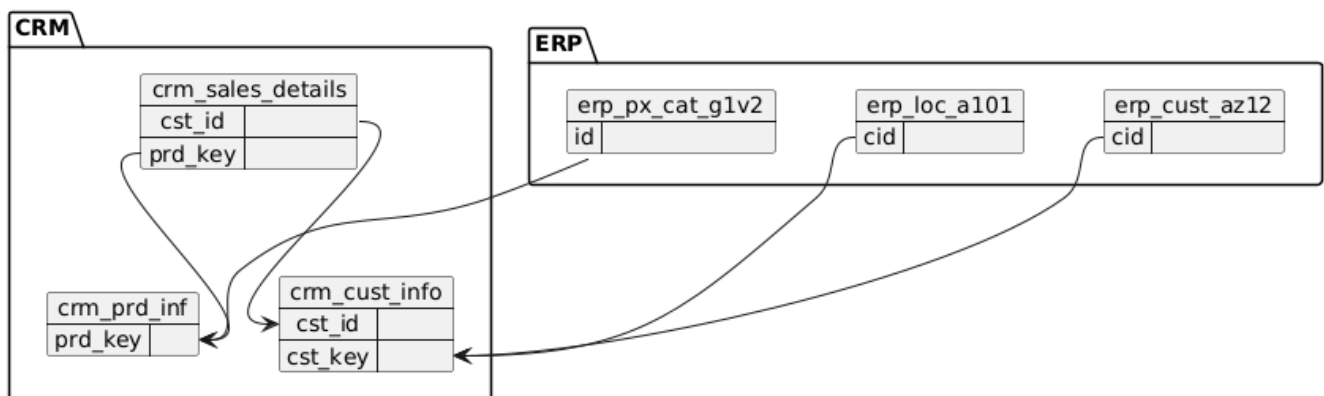
The Silver layer serves as the intermediate processing stage in the data pipeline. In this layer, raw data from the Bronze layer is cleaned, transformed, and enriched to enhance its quality and usability for analysis.

The key functions of the Silver layer include:

- **Data Cleaning:** Removing duplicates, handling missing values, and correcting inconsistencies in the data.
- **Data Transformation:** Converting data into a more structured format, applying business rules, and aggregating information as needed.
- **Data Enrichment:** Integrating additional data sources to provide more context and insights.

By implementing these processes in the Silver layer, we ensure that the data is reliable and ready for advanced analytics and reporting in the subsequent Gold layer.

## Data integration analysis and sketch



# Coding the transformations

Add metadata columns to track data lineage and help with debugging. We'll add a **created** column to each silver table. Following the data engineer naming conventions, and Postgres syntax, we add this column to the create script:

```
dwh_created    TIMESTAMP DEFAULT CURRENT_TIMESTAMP;
```

## CRM customer info table

Next, we gotta check the data quality of the bronze tables.

First, we'll check for nulls in the primary key columns.

```
select cst_id, count(*) from silver.crm_cust_info
group by cst_id having count(*) > 1 or cst_id is null
```

Reveals some duplicates.

	cst_id	↕	count	↕
1	NULL		4	
2	29473		2	
3	29449		2	
4	29433		2	
5	29466		3	
6	29483		2	

Focus on one to see what the problem is.

```
select * from bronze.crm_cust_info
where cst_id = '29466';
```

Results:

	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

This shows that there are multiple entries with different dates. Since this data warehouse specification doesn't need historical data, we'll take the latest.

But remember, we leave bronze the same and transform (clean) as we go. To do this, let's grab the latest record per **cst\_id** and rank them.



```
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 = 29466;
```

We'll apply this to the whole table to see if this same procedure flags all duplicates.

```
select * from (
    select *, row_number() over (
        partition by cst_id order by cst_create_date desc
    ) as flag_last
    from bronze.crm_cust_info
) sub
where flag_last != 1;
```

Checking against the duplicates found earlier, when we change `flag_last = 1`, we see that all duplicates are removed.

That's just one example of data cleaning in the silver layer. See the `src/silver` folder for the full SQL scripts.

In addition to **removing duplicates in primary keys**, we will do things like:

- **Remove unwanted spaces in text fields**
- **Data normalization or standardization**, e.g. consistency of values in low cardinality columns
- **Handling missing values appropriately**, e.g. replacing NULLs with 'Unknown' or default values

## CRM product info and sales tables

For the `crm_prd_info` table, we will, in addition:

- **Derive new columns**, e.g. extracting category from product codes
- **Data type conversions**, e.g. casting timestamps with `00:00:00` times to date
- **Data enrichment**, e.g. making sales end dates make sense and not overlap with start dates

For the `crm_sales_details` table, the sales, quantity, and price columns are tricky. After doing some quality checks, we discover quite a few places where things don't add up. Before doing **data enrichment**, we would verify with the data owner what to do with these bad values.

The rules Baraa came up with are:

- If sales is negative, zero, null, or not equal to derived sales from quantity \* price, set it to derived value
- If price is zero or null, derive it from sales / quantity (or null if 0 quantity to prevent div by 0)

- If prices is negative, set it to absolute value (positive)

## ERP tables

In addition to deriving a valid foreign key to `crm_cust_info`, there are birthdays that are in the future and as old as 1917! After consulting the data owner, we decide to set future birthdays to null and leaving old birthdays as is.

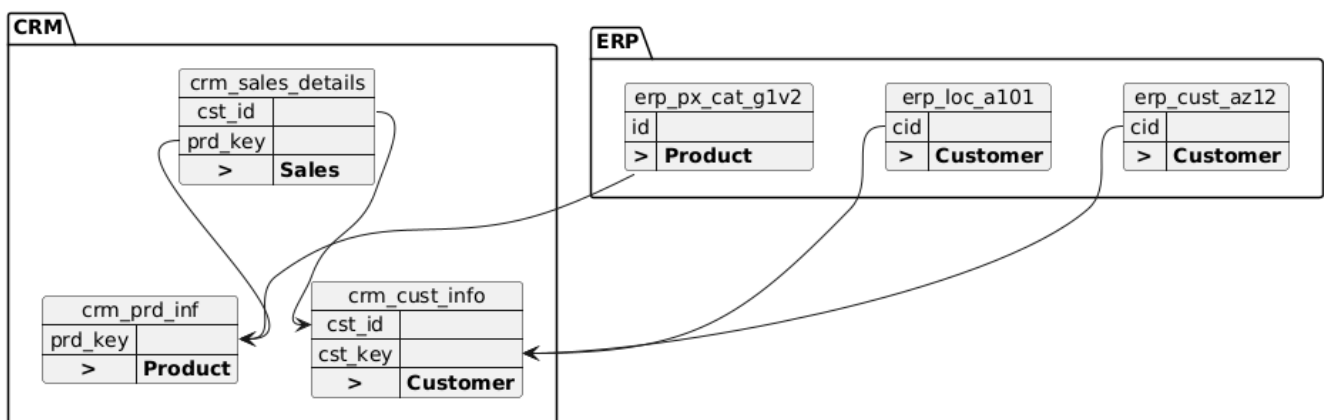
For the `erp_loc_a101` table, besides deriving the foreign key, we also normalize the country fields.

For the `erp_px_cat_g1v2` table, we already derived a category id column in the `cust_prod_info` silver table that matches the id, and when checking spaces, standardization, and consistency, we find no issues.

## Gold layer

The Gold layer represents the final stage in the data pipeline, where the data is refined and optimized for business intelligence and analytics, separated into facts and dimensions. For this project, the star schema model is sufficient as there are only a couple of product and customer *dimensions* surrounding the sales orders *facts*.

First step on how to determine how to build our gold data model is to examine the silver data model and add some labels to group things logically.



Referring back to the [Specifications](#), we see that rather than loading data, we'll be creating views into the cleaned silver tables, views that might aggregate and integrate information.

Another thing to add to the views are surrogate keys to use as the primary keys. TIL that these are exposed in the fact view SELECTs, while it's fact view JOINS that create the links between the relevant columns gold views (rather than using those columns as foreign keys). Only after JOINing the fact table are the foreign keys resolved. That's why it seemed that the surrogates could be sorted however you want. (I think.)

To verify the joins (foreign key integrity) execute joins on the keys and see if any keys are NULL:

```
SELECT COUNT(*) FROM gold.fact_sales f
LEFT JOIN gold.customers c ON f.customer_key = c.customer_key
```

```
LEFT JOIN gold.products c ON f.product_key = c.product_key  
WHERE p.product_key IS NULL OR c.customer_key IS NULL;
```