

# Project: Data Warehouse with PostgreSQL

## Table of Contents

Preparing the environment .....	1
Documentation .....	1
Data Flow .....	2
Bronze layer .....	3
Silver layer .....	3
Data integration analysis and sketch .....	3
Gold layer .....	3

## 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
```

For diagrams, some additional packages are required.

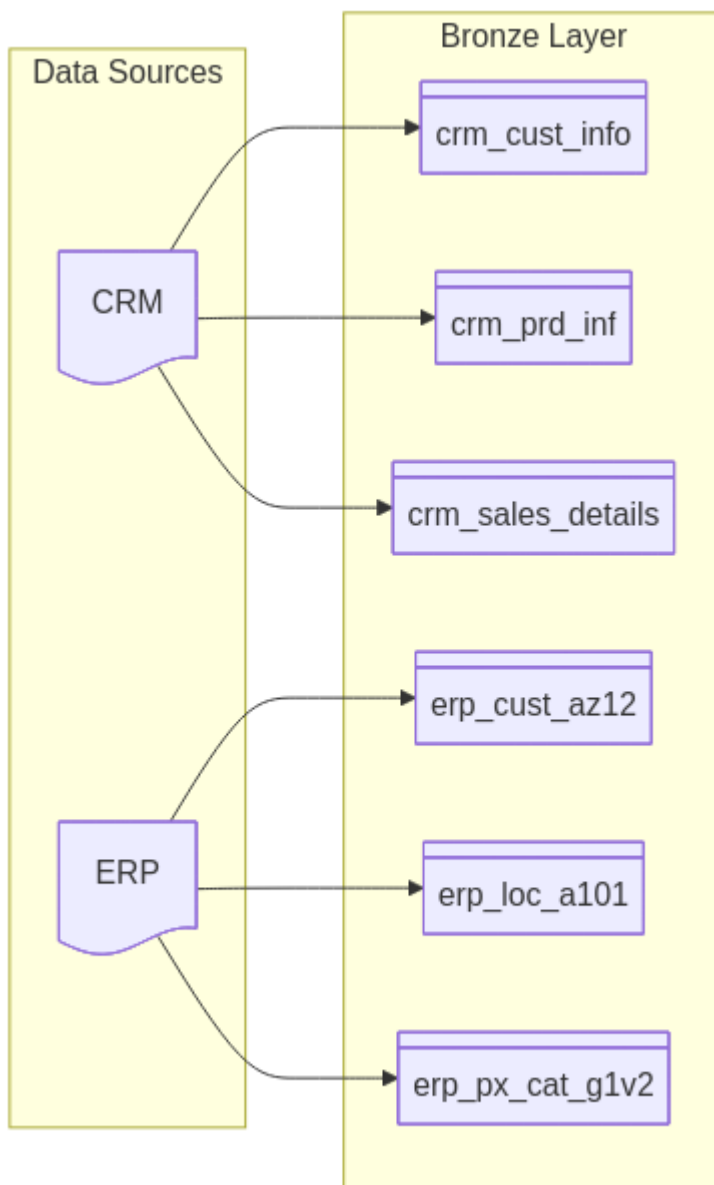
```
sudo apt install graphviz plantuml
```

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
```

## Data Flow



# 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

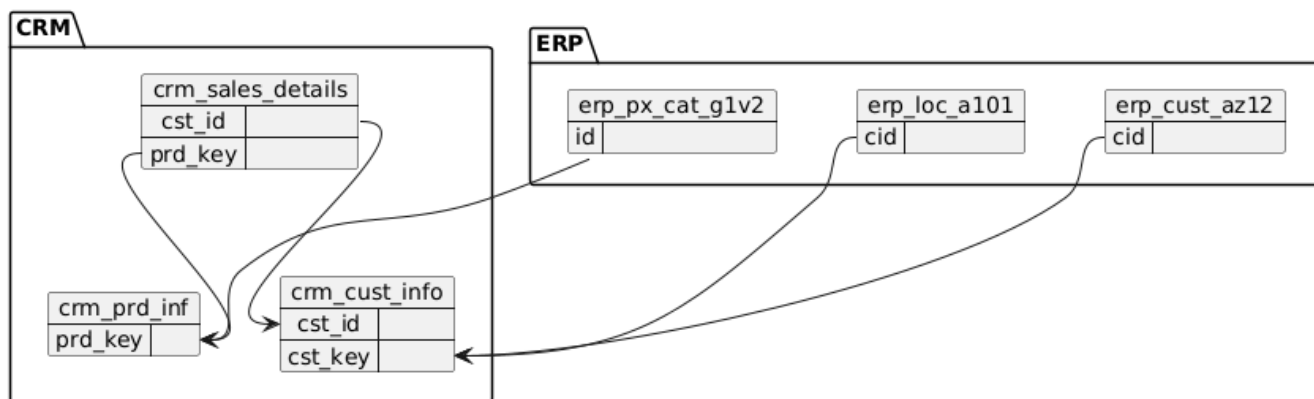
The Silver layer serves as the intermediate processing stage in the data pipeline. In this layer, raw data from the Bronze layer is cleaned, transformed, and enriched to enhance its quality and usability for analysis.

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



# Gold layer

The Gold layer represents the final stage in the data pipeline, where data is fully refined and optimized for business intelligence and analytics. In this layer, data from the Silver layer is aggregated, summarized, and structured to meet specific reporting and analytical needs.

The key functions of the Gold layer include:

- **Data Aggregation:** Summarizing data to provide high-level insights, such as totals, averages, and other key performance indicators (KPIs).

- Data Structuring: Organizing data into star or snowflake schemas to facilitate efficient querying and reporting.
- Data Optimization: Enhancing performance through indexing, partitioning, and other techniques to ensure fast access to critical information.

By implementing these processes in the Gold layer, we ensure that the data is not only accurate and reliable but also tailored to support strategic decision-making and advanced analytics across the organization.