



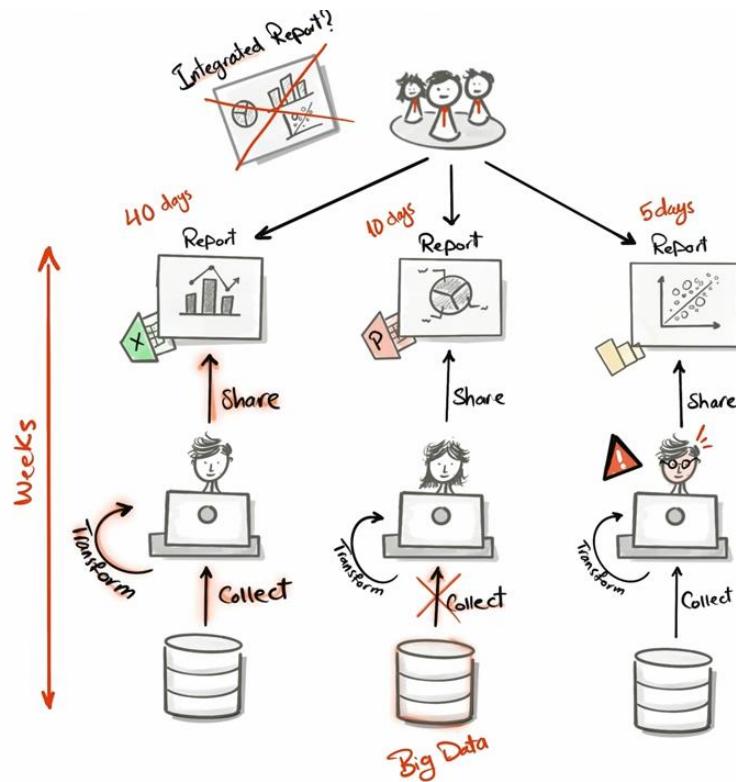
# DATA WAREHOUSE



A subject-oriented, integrated, time-variant, and non-volatile collection of data in support of management's decision-making process.

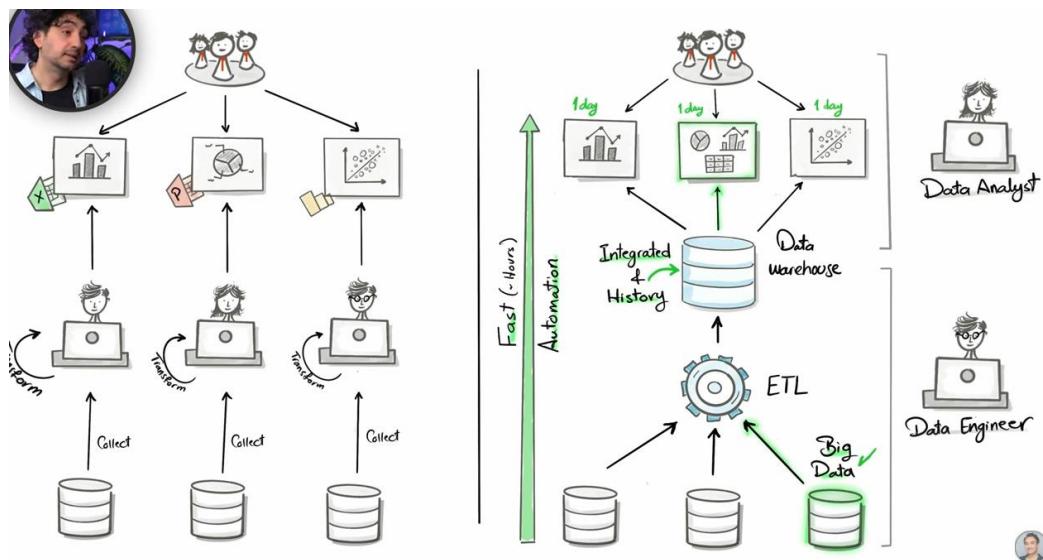
When data is collected from multiple sources, Data Analysts manually compile and generate reports for stakeholders. However, this approach is time-consuming, especially when handling large data sources (Big Data), which can lead to delays in report generation. As a result, this method is inefficient and tedious for large-scale reporting.

Moreover, if we need to create an integrated report that combines insights from multiple reports built on different sources, the process becomes even more complex and inefficient. This setup lacks effective data management, making it challenging to maintain consistency, accuracy, and scalability.



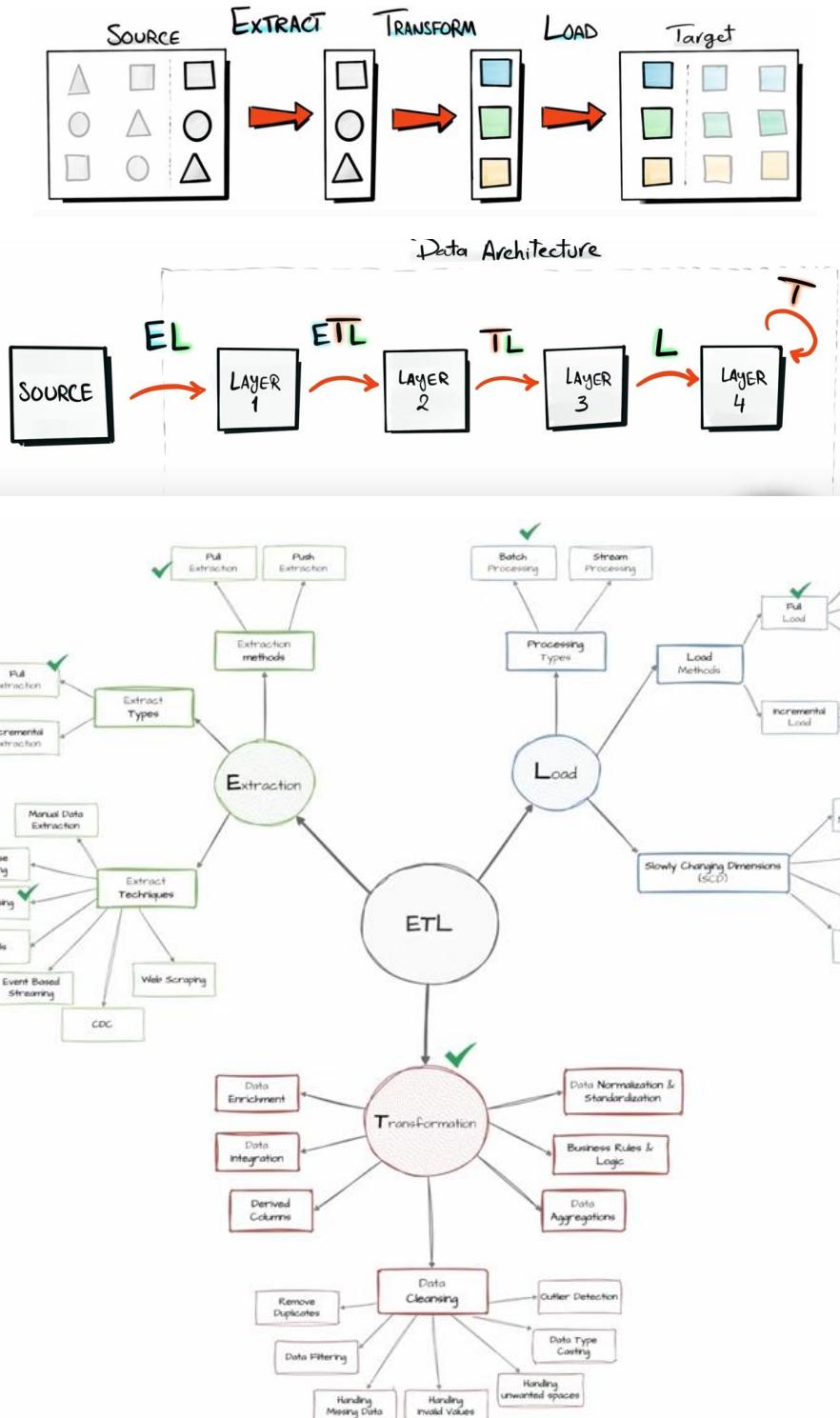
To overcome this challenge, we need a Data Warehouse where a copy of the source data and transformed data is continuously stored in a single database, Delta Lakehouse, or an S3 bucket. Data Engineers perform ETL operations to move data from various sources into the Data Warehouse, ensuring that all relevant information is centralized and up to date.

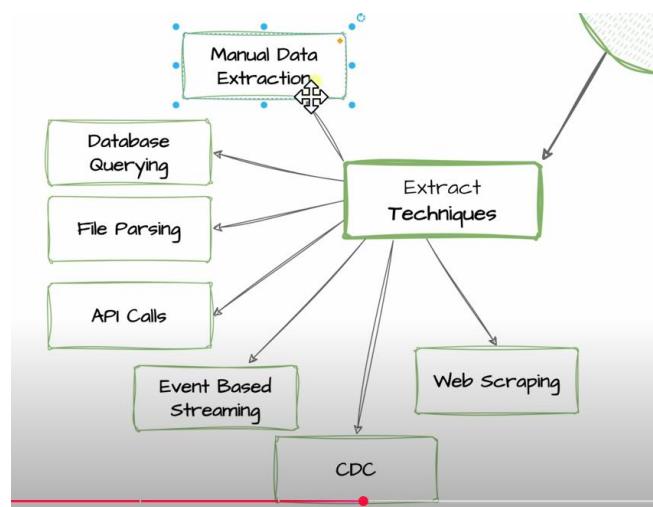
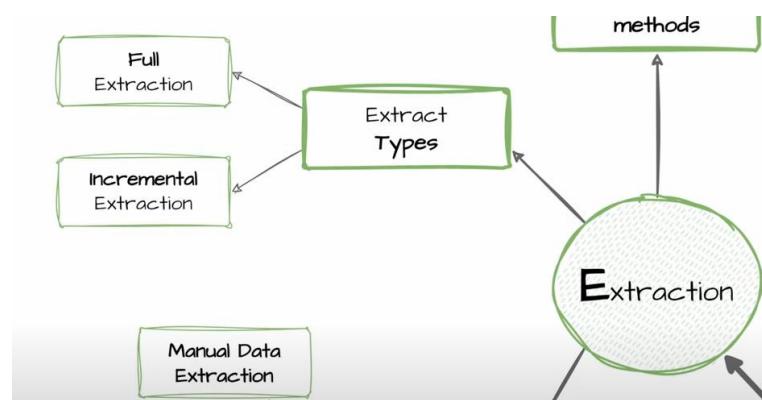
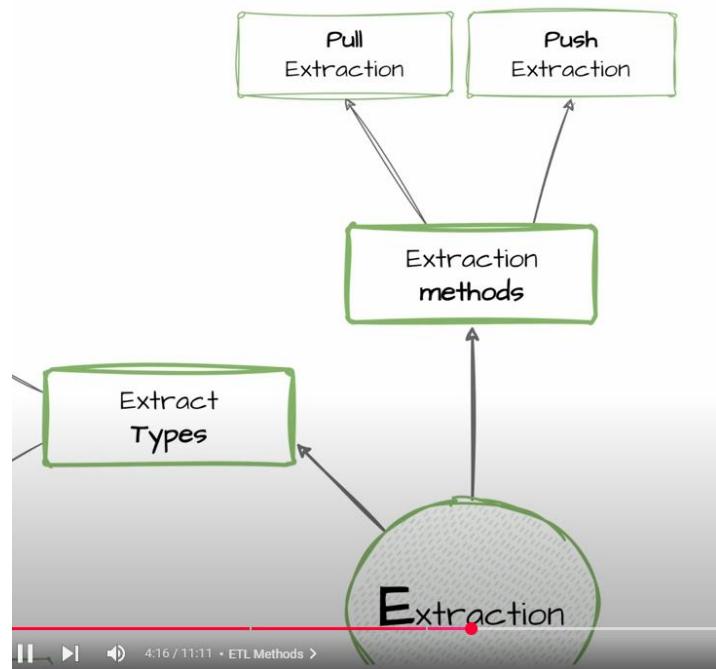
Once the data is available in the Data Warehouse, Data Analysts can efficiently generate reports for stakeholders. Additionally, integrated reports can be built more easily since the Data Warehouse always holds the latest data, thanks to ongoing ETL processes. This approach improves data management, enhances reporting efficiency, and ensures timely insights.

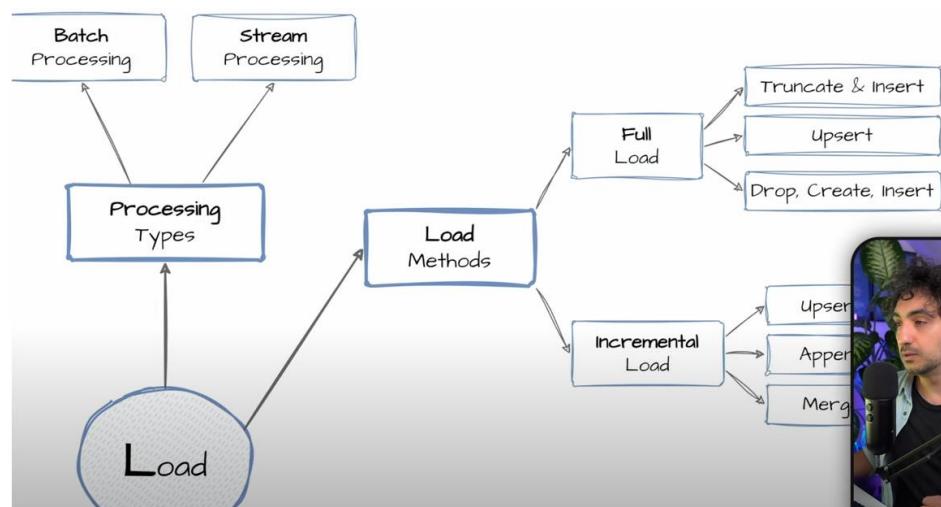
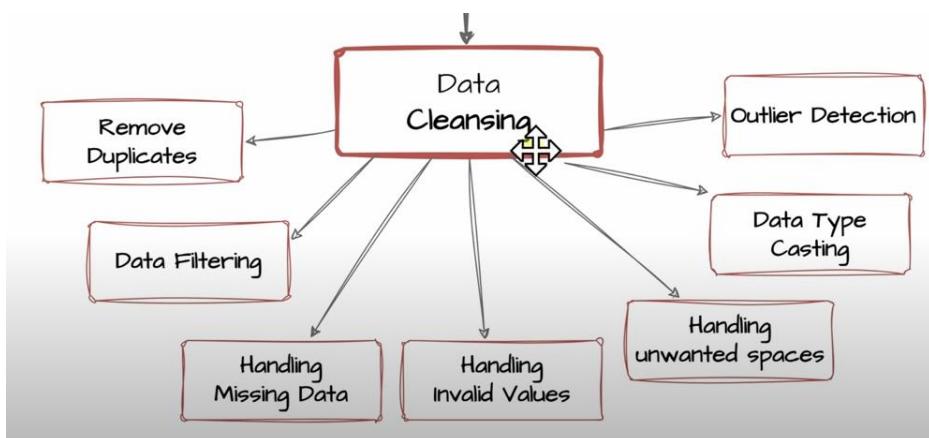
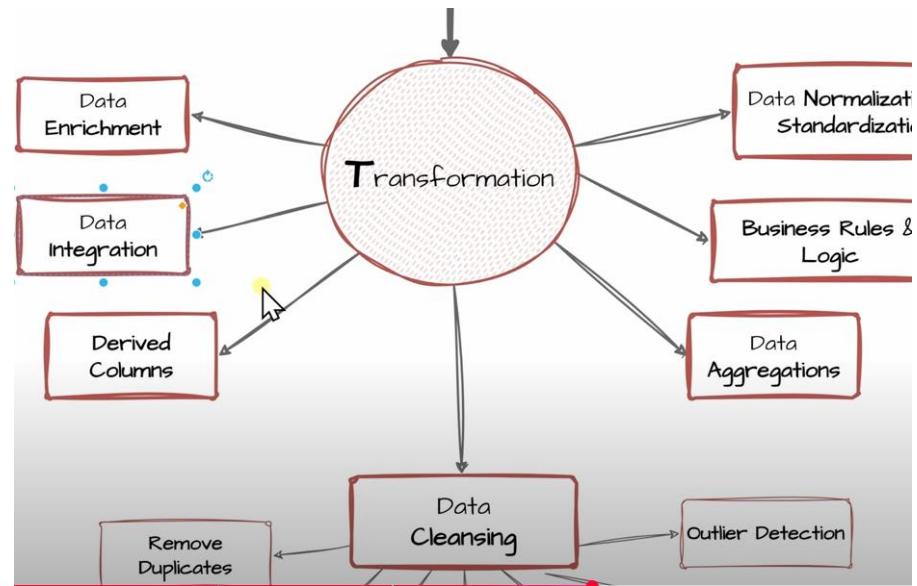


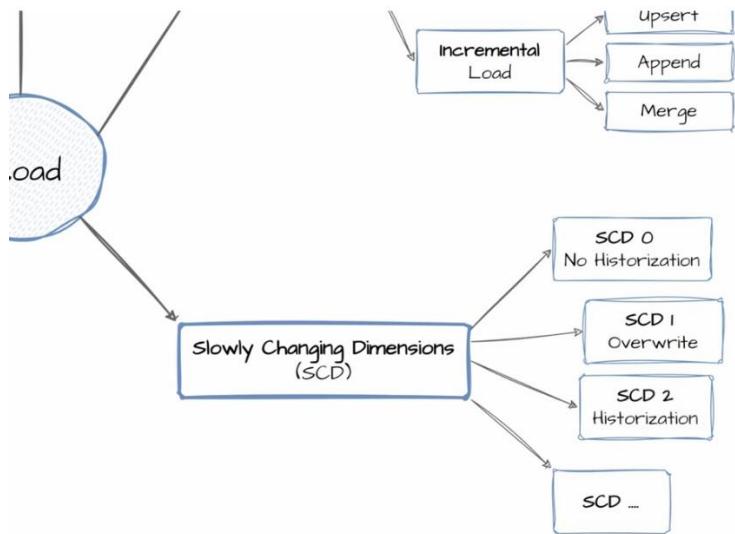
ETL

ETL is a fundamental component of any Data Warehouse project. It plays a crucial role in extracting, transforming, and loading data, but the implementation of ETL varies for each project. The approach to ETL is entirely dependent on the architecture of the system, which dictates how data is sourced, processed, and stored to meet specific business requirements.









## Notion

Notion screenshot showing a project dashboard for a Data Warehouse Project.

**Project Epics:**

Name	Progress
Requirement Analysis	0.00%
Design Data Architecture	0.00%
Project Initialization	0.00%

+ New page

---

**Project Tasks:**

- ↑ Tasks ▾ | ↴ Tasks ▾
- ▶ Requirement Analysis
- ▶ Design Data Architecture
- ▶ Project Initialization

▼ 1 hidden group

+ Add a group

# Requirement Analysis

## 🚀 Project Requirements

### Building the Data Warehouse (Data Engineering)

#### 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



### BI: Analytics & Reporting (Data Analysis)

#### Objective

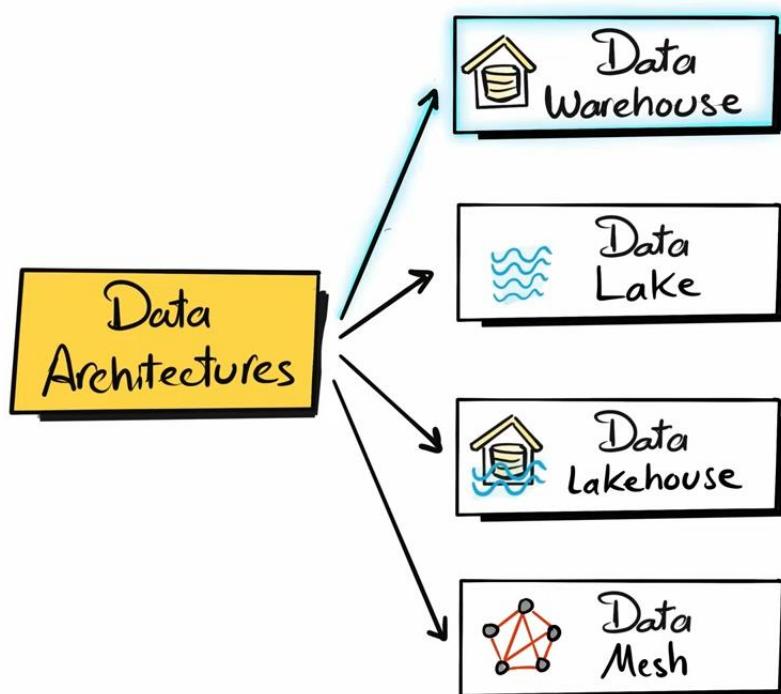


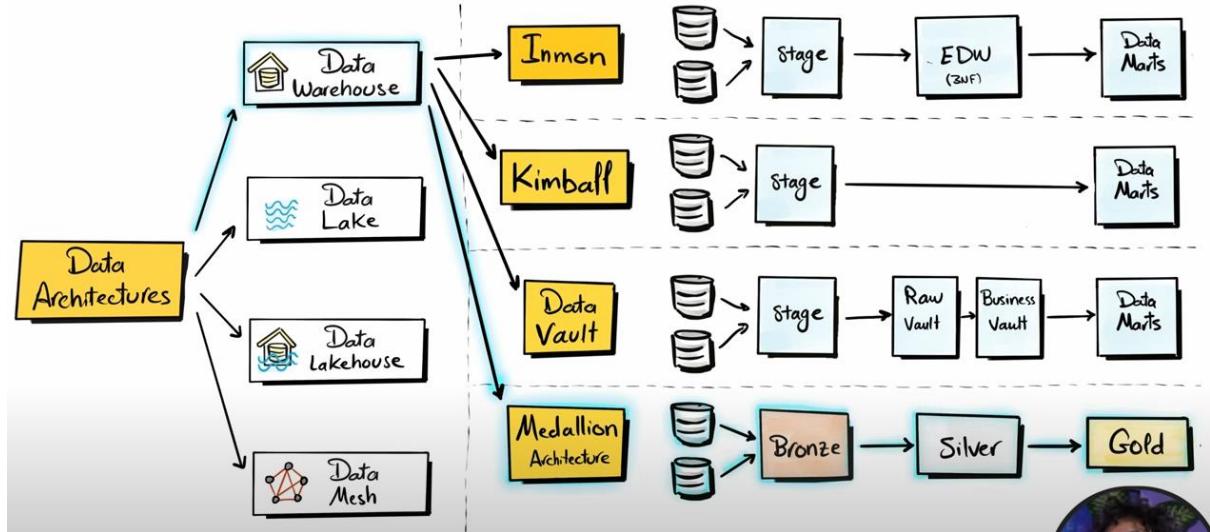
Develop SQL-based analytics to deliver detailed insights into:

- Customer Behavior
- Product Performance
- Sales Trends

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

For more details, refer to [docs/requirements.md](#).





# Design The Data Architecture

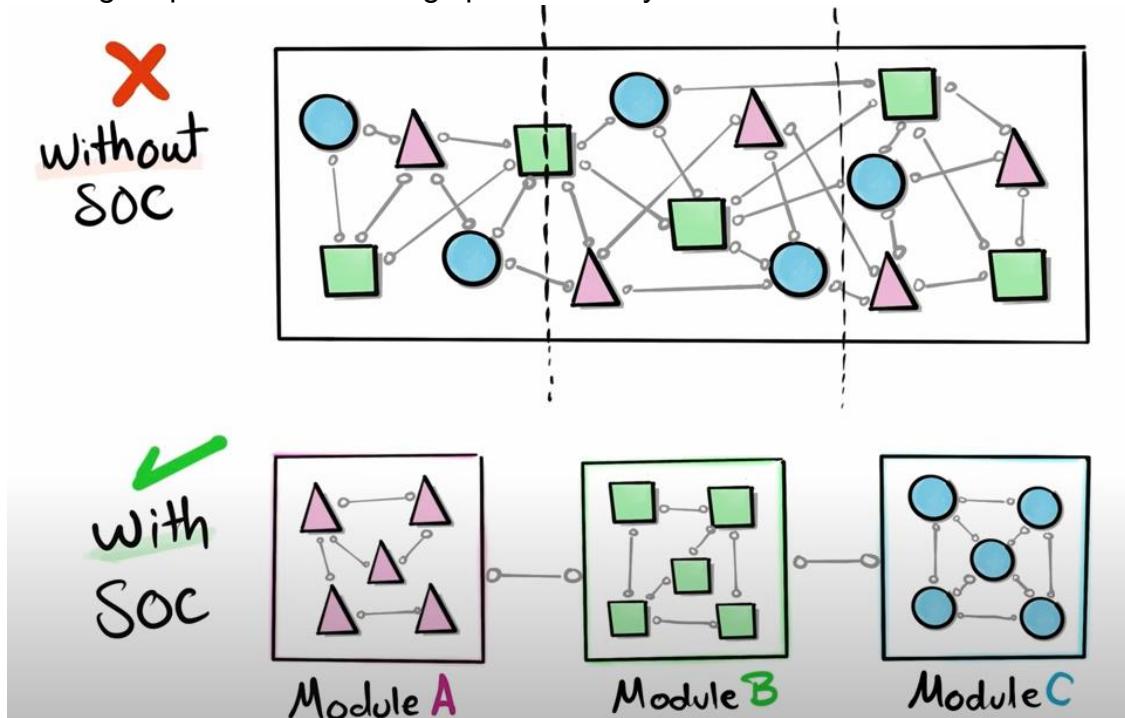
# Design the Layers

	Bronze Layer	Silver Layer	Gold Layer
Definition	Raw, unprocessed data as-is from sources	Clean & 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
Data Transformation	None (as-is) <ul style="list-style-type: none"> <li>- Data Cleaning</li> <li>- Data Standardization</li> <li>- Data Normalization</li> <li>- Derived Columns</li> <li>- Data Enrichment</li> </ul>	<ul style="list-style-type: none"> <li>- Data Integration</li> <li>- Data Aggregation</li> <li>- Business Logic &amp; Rules</li> </ul>	
Data Modeling	None (as-is)	None (as-is)	<ul style="list-style-type: none"> <li>- Start Schema</li> <li>- Aggregated Objects</li> <li>- Flat Tables</li> </ul>
Target Audience	- Data Engineers	- Data Analysts - Data Engineers	- Data Analysts - Business Users

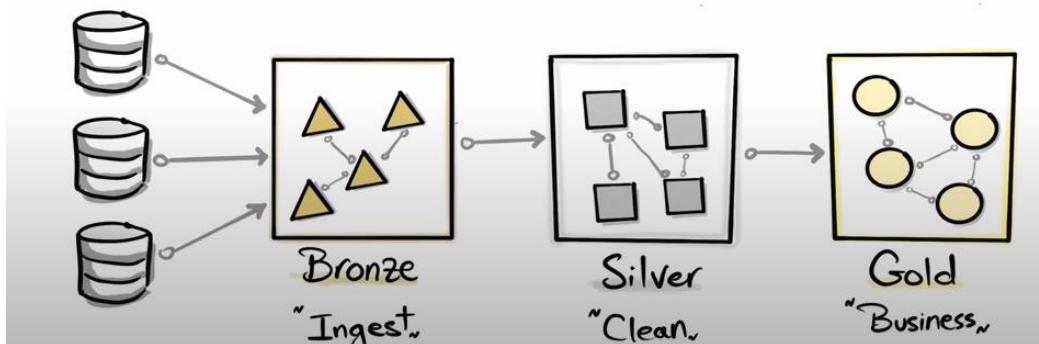
# SOC SEPARATION OF CONCERNS

## Separation of Concerns in Data Architecture

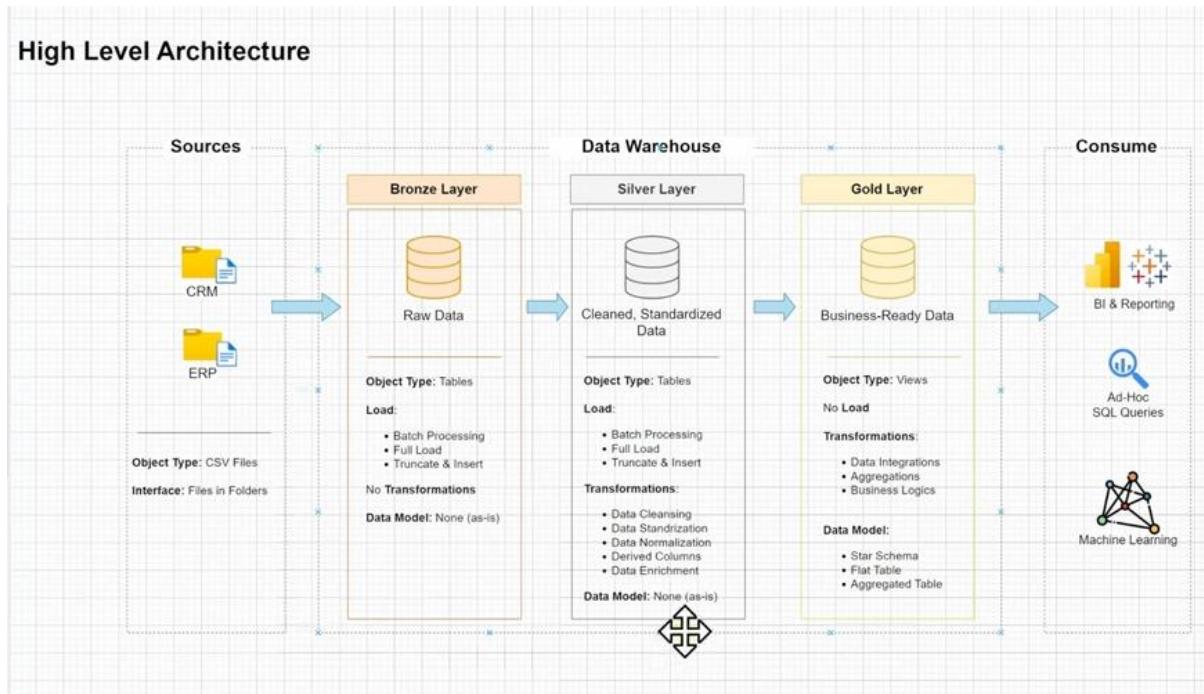
Separation of concerns (SoC) in data architecture ensures that different components of a data system are designed and managed independently, improving scalability, maintainability, and security. It helps in structuring data systems effectively by dividing responsibilities among specialized layers.



A well-designed data project ensures that models and processes operate independently, allowing for better scalability, maintainability, and flexibility.



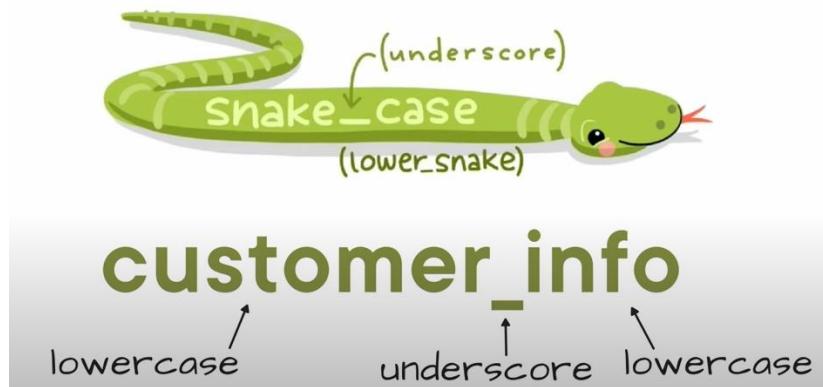
## High Level Architecture



# Naming Conventions

Set of Rules or Guidelines  
for naming anything in the project.

- Database
- Schema
- Tables
- Store Procedures ...



## General Principles

- **Naming Conventions:** Use snake\_case, with lowercase letters and underscores ( `_` ) to separate words.
- **Language:** Use English for all names.
- **Avoid Reserved Words:** Do not use SQL reserved words as object names.

## Table Naming Conventions

### Bronze Rules

- All names must start with the source system name, and table names must match their original names without renaming.
- `<sourcesystem>_<entity>`
  - `<sourcesystem>` : Name of the source system (e.g., `crm`, `erp` ).
  - `<entity>` : Exact table name from the source system.
  - Example: `crm_customer_info` → Customer information from the CRM system.

### Silver Rules

- All names must start with the source system name, and table names must match their original names without renaming.
- `<sourcesystem>_<entity>`
  - `<sourcesystem>` : Name of the source system (e.g., `crm`, `erp` ).
  - `<entity>` : Exact table name from the source system.
  - Example: `crm_customer_info` → Customer information from the CRM system.

### Gold Rules

- All names must use meaningful, business-aligned names for tables, starting with the category prefix.
- `<category>_<entity>`
  - `<category>` : Describes the role of the table, such as `dim` (dimension) or `fact` (fact table).
  - `<entity>` : Descriptive name of the table, aligned with the business domain (e.g., `customers`, `products`, `sales` ).
  - Examples:
    - `dim_customers` → Dimension table for customer data.
    - `fact_sales` → Fact table containing sales transactions.

### Glossary of Category Patterns

Pattern	Meaning	Example(s)
<code>dim_</code>	Dimension table	<code>dim_customer</code> , <code>dim_product</code>
<code>fact_</code>	Fact table	<code>fact_sales</code>
<code>agg_</code>	Aggregated table	<code>agg_customers</code> , <code>agg_sales_monthly</code>

## Column Naming Conventions

### Surrogate Keys

- All primary keys in dimension tables must use the suffix `_key`.
- `I<table_name>_key`
  - `<table_name>` : Refers to the name of the table or entity the key belongs to.
  - `_key` : A suffix indicating that this column is a surrogate key.
  - Example: `customer_key` → Surrogate key in the `dim_customers` table.

### Technical Columns

- All technical columns must start with the prefix `dwh_`, followed by a descriptive name indicating the column's purpose.
- `dwh_<column_name>`
  - `dwh` : Prefix exclusively for system-generated metadata.
  - `<column_name>` : Descriptive name indicating the column's purpose.
  - Example: `dwh_load_date` → System-generated column used to store the date when the record was loaded.

## Stored Procedure

- All stored procedures used for loading data must follow the naming pattern:
- `load_<layer>`.
  - `<layer>` : Represents the layer being loaded, such as `bronze`, `silver`, or `gold`.
  - Example:
    - `load_bronze` → Stored procedure for loading data into the Bronze layer.
    - `load_silver` → Stored procedure for loading data into the Silver layer.



The screenshot shows a GitHub repository page for 'sql-data-warehouse-project'. The repository has 1 branch and 0 tags. It contains files like README, LICENSE, and README.md. The commit history shows several commits from TharunMedishetty, mostly creating placeholder files. The repository is public and is described as 'Building a modern data warehouse with SQL server including ETL processes, data modeling and analytics.'

# Markdown (.md)

Lightweight markup language that you can use  
to add formatting elements to plaintext text documents.

The diagram illustrates the transformation of raw Markdown code into a formatted document. On the left, the raw Markdown code is shown:

```
## What is Markdown?
see [Markdown](www.markdownlink.com)

> Markdown is a lightweight markup language
----

## List of tips
1. *One asterisk Italicizes*
2. **Two asterisks emphasize**
```

An arrow points from this raw code to the right, where the resulting formatted document is displayed:

**What is Markdown?**

see [Markdown](#)

Markdown is a lightweight markup language

---

**List of tips**

1. *One asterisk Italicizes*
2. **Two asterisks emphasize**

## Project Initialization

# Create Database & Schemas

The screenshot shows two instances of SSMS. The top instance has an Object Explorer window showing a connection to 'LAPTOP-80PBM4NJ\SQLEXPRESS'. Under 'Databases', 'DataWarehouse' is selected. The bottom instance also has an Object Explorer window showing a connection to 'LAPTOP-80PBM4NJ\SQLEXPRESS (SQL Server 16.0.1)'. In this instance, 'DataWarehouse' is expanded, revealing its schema and security settings. Both instances have a 'SQLQuery2.sql' query editor tab open. The top query editor contains the following SQL code:

```

1 USE master;
2
3 CREATE DATABASE DataWarehouse;
4
5 USE DataWarehouse;

```

The bottom query editor contains the following SQL code, which includes comments for context:

```

1 -- Creating Database
2
3 USE master;
4
5 CREATE DATABASE DataWarehouse;
6
7 USE DataWarehouse;
8
9 -- Creating Schemas
10
11 CREATE SCHEMA bronze;

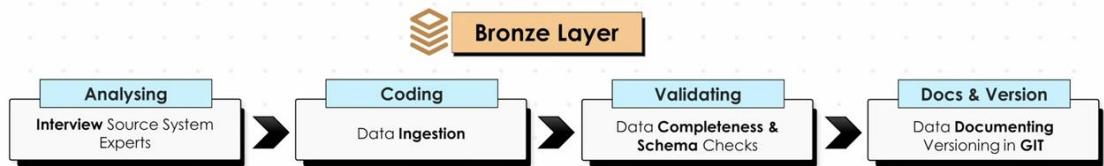
```

The 'Messages' pane at the bottom right of the bottom SSMS instance shows the command completed successfully.

**GO**

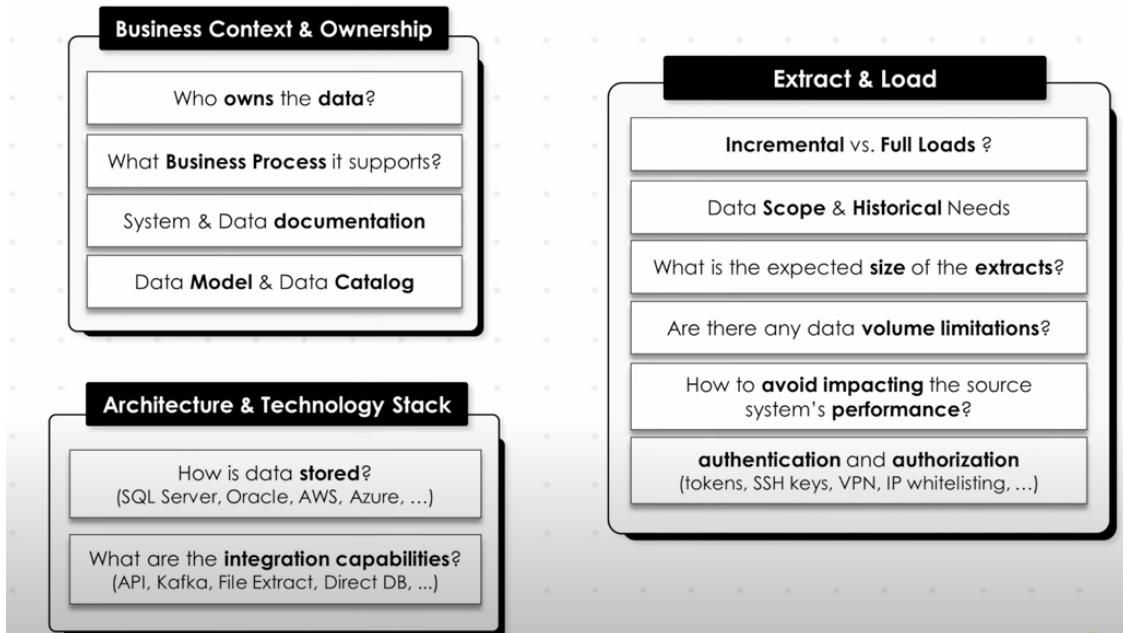
**separate batches when working with multiple SQL statements**

## Bronze Layer



## Build Bronze Layer

### Analyse Source Systems



Bronze Layer	
Definition	Raw, unprocessed data <b>as-is</b> from sources
Objective	Traceability & Debugging
Object Type	Tables
Load Method	Full Load (Truncate & Insert)
Data Transformation	None ( <b>as-is</b> )
Data Modeling	None ( <b>as-is</b> )
Target Audience	- Data Engineers

DDI

Data Definition Language defines the structure of database tables

Consult the technical experts of the source system to understand its metadata.

# DATA PROFILING

Explore the data to identify column names and data types

## Bronze Rules

- All names must start with the source system name, and table names must match their original names without renaming.
  - **<sourcesystem>\_<entity>**
    - **<sourcesystem>** : Name of the source system (e.g., `crm` , `erp` ).
    - **<entity>** : Exact table name from the source system.
    - Example: `crm_customer_info` → Customer information from the CRM system.

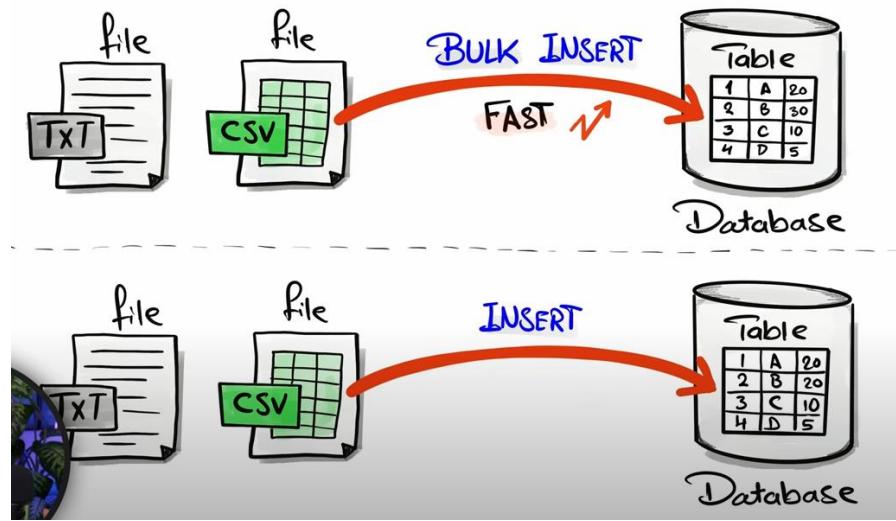


The screenshot shows the Microsoft SQL Server Management Studio (SSMS) interface. The Object Explorer on the left displays the database structure under 'LAPTOP-80PBM4NJ\SQLEXPRESS (SQL Server 16.0.1)'. It includes sections for Databases, System Databases, Database Snapshots, DataWarehouse, Database Diagrams, and Tables. Under Tables, several tables are listed: bronze.crm\_cust\_info, bronze.crm\_prd\_info, bronze.crm\_sales\_details, bronze.erp\_cust\_az12, bronze.erp\_loc\_a101, and bronze.erp\_px\_cat\_g1v2. The right pane shows a SQL query window with the following code:

```
1 USE DataWarehouse;
2
3
4 IF OBJECT_ID('bronze.crm_cust_info', 'U') IS NOT NULL
5     DROP TABLE bronze.crm_cust_info
6
7 CREATE TABLE bronze.crm_cust_info (
8     cst_id INT,
9     cst_key NVARCHAR(50),
10    cst_firstname NVARCHAR(50),
11    cst_lastname NVARCHAR(50),
12    cst_marital_status NVARCHAR(50),
13    cst_gndr NVARCHAR(50),
14    cst_create_date DATE
15 );
16
17 IF OBJECT_ID('bronze.crm_prd_info', 'U') IS NOT NULL
18     DROP TABLE bronze.crm_prd_info
19
20 CREATE TABLE bronze.crm_prd_info (
21     prd_id INT,
22     prd_key NVARCHAR(50),
23     prd_name NVARCHAR(50),
24     prd_desc NVARCHAR(250),
25     prd_qty INT,
26     prd_price DECIMAL(10, 2),
27     prd_status NVARCHAR(50),
28     prd_create_date DATE
29 );
```

### **Build Bronze Layer**

## Develop SQL Load Scripts



In a **bulk insert**, all the data from the files is loaded into the database at once, making the process faster and more efficient. In contrast, a standard **INSERT** statement inserts data **row by row**, which can be significantly slower, especially for large datasets.

DWH\_Project\_Extra...PBM4NJ\rehan (60) DWH\_Project\_Creat...PBM4NJ\rehan (80) SQLQuery2.sql - L...PBM4NJ\rehan (61)\*

```

1 USE DataWarehouse;
2
3 BULK INSERT [bronze].[crm_cust_info]
4 FROM 'D:\Tharun__Personal\SQL by Baara\Project\sql-data-warehouse-project\datasets\source_crm\cust_info.csv'
5 WITH (
6     FIRSTROW = 2,
7     FIELDTERMINATOR = ',',
8     TABLOCK
9 );

```

**HINT**

Save frequently used SQL code in stored procedures in database

BULK INSERT bronze.erp cust az12

**Add PRINTS**

Add Prints to track execution, debug issues, and understand its flow

**Add TRY...CATCH**

Ensures error handling, data integrity, and issue logging for easier debugging

**TRY...CATCH**

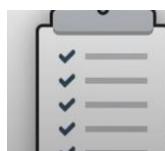
SQL runs the TRY block, and if it fails, it runs the CATCH block to handle the error

**Track ETL Duration**

Helps to identify bottlenecks, optimize performance, monitor trends, detect issues

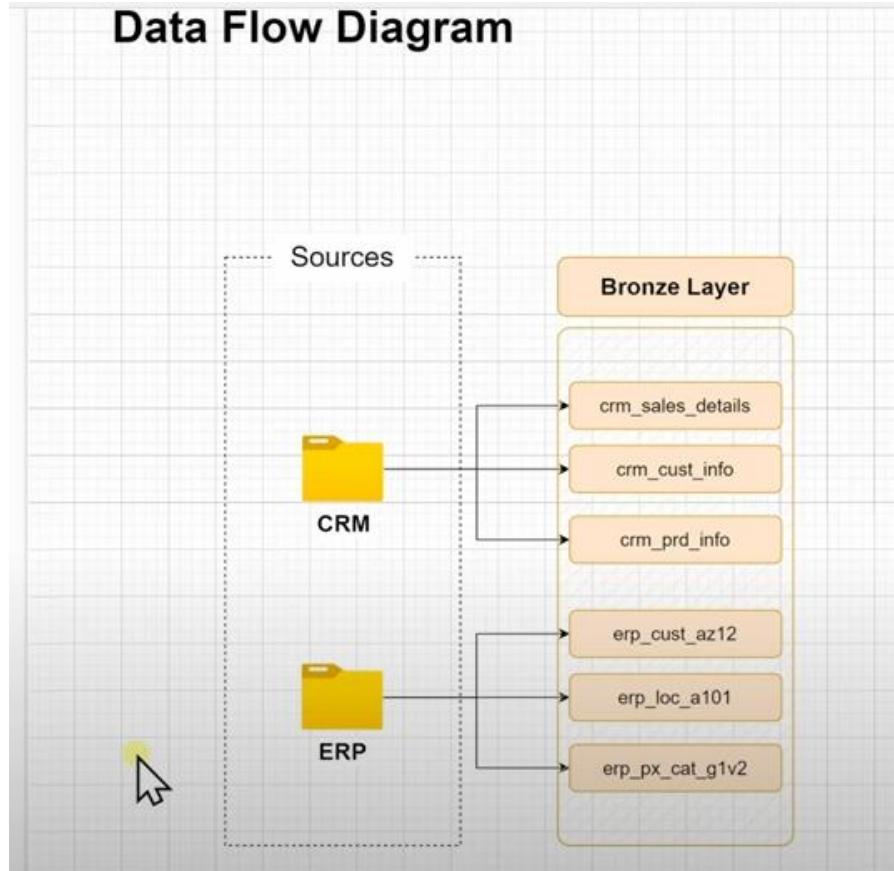
## DATEDIFF ()

calculates the difference between two dates, returns days, months, or years

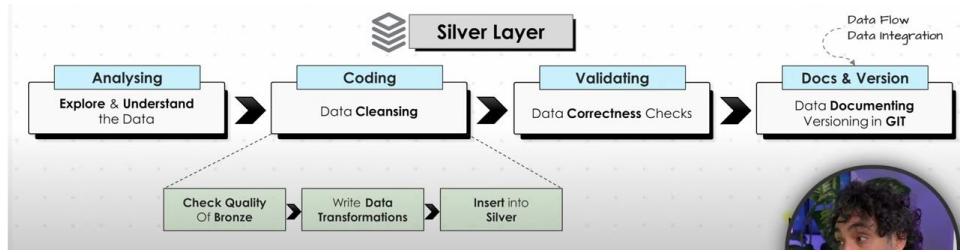


Calculate the Duration of Loading Bronze Layer "Whole Batch"

## Data Flow Diagram



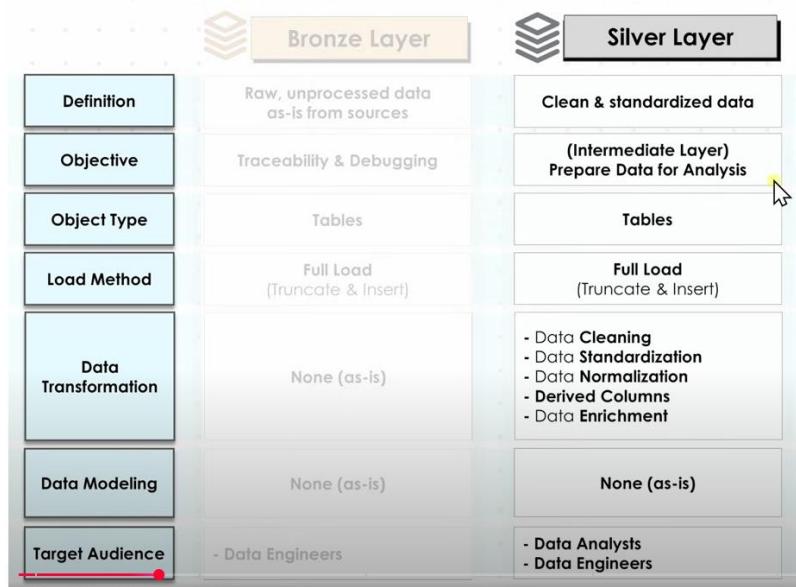
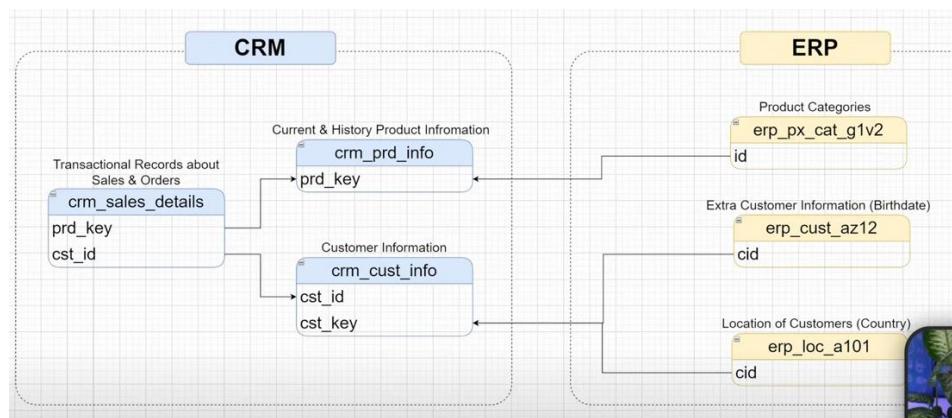
# Silver Layer



# Build Silver Layer

# Explore & Understand The Data

## Document & Visualize What You Understand from Data



## METADATA COLUMNS

Extra columns added by data engineers  
that do not originate from the source data.

**create\_date** : The record's load timestamp.

**update\_date**: The record's last update timestamp.

**source\_system**: The origin system of the record.

**file\_location**: The file source of the record.

Initially, we create silver tables with the same structure as bronze tables. After reviewing the data across different tables, we determine whether additional metadata columns are needed and incorporate them accordingly. Similarly, in our project, we follow this approach by first replicating the bronze tables in Silver and later enhancing them as required.

Table : silver.load\_silver

	cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date
1	11000	AW00011000	Jon	Yang	M	M	2025-10-06
2	11001	AW00011001	Eugene	Huang	S	M	2025-10-06
3	11002	AW00011002	Ruben	Torres	M	M	2025-10-06
4	11003	AW00011003	Christy	Zhu	S	F	2025-10-06
5	11004	AW00011004	Elizabeth	Johnson	S	F	2025-10-06
6	11005	AW00011005	Julio	Ruiz	S	M	2025-10-06
7	11006	AW00011006	Isaac	Akesson	S	F	2025-10-06

## Quality Check

A Primary Key must be unique and not null

We have duplicates in the **primary key** and **NULL values** in the cst\_id column. Upon closer inspection, the duplicates have different **cst\_create\_date** values. To clean the data:

1. Remove records where cst\_id is NULL
2. For duplicates, retain only the latest record based on cst\_create\_date

```

quality_checks_si...OPBM4NJ\rehan (52)* # X SQLQuery5.sql - L...PBM4NJ\rehan (71))^ DWH_Proj
1 -----silver.crm_cust_info-----
2 -- Check for Nulls or Duplicates in Primary Key
3 -- Expectation: No Result
4
5 SELECT
6 cst_id,
7 COUNT(*)
8 FROM bronze.crm_cust_info
9 GROUP BY cst_id
10 HAVING COUNT(*)>1 or cst_id IS NULL;
11
12 SELECT * FROM bronze.crm_cust_info WHERE cst_id = 29466;
13

```

100 % ▾

Results Messages

	cst_id	(No column name)
1	29449	2
2	29473	2
3	29433	2
4	NULL	3
5	29483	2
6	29466	3

	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

```

15
16 --query to remove the duplicates and Nulls in primary key
17 SELECT *
18
19 FROM
20   (SELECT
21     *
22     ROW_NUMBER() OVER(PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_latest
23   FROM
24     bronze.crm_cust_info
25     WHERE cst_id IS NOT NULL) t
26 WHERE flag_latest = 1
--
```

Now, we need to check for the string columns (cst\_firstname,cst\_lastname)

```
-- Check for unwanted Spaces  
SELECT cst_firstname  
FROM bronze.crm_cust_info  
WHERE cst_firstname != TRIM(cst_firstname)
```

### TRIM()

Removes leading and trailing spaces from a string

If the original value is not equal to  
the same value after trimming,  
it means there are spaces!

```
--query to remove the unwanted spaces from the string columns  
  
=SELECT  
cst_id,  
cst_key,  
TRIM(cst_firstname) AS cst_firstname,  
TRIM(cst_lastname) AS cst_lastname,  
cst_marital_status,  
cst_gndr,  
cst_create_date  
FROM  
bronze.crm_cust_info
```

cst\_create\_date

### Quality Check

Check the consistency of values in low cardinality columns

	cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date
97	11096	AW00011096	Andr®s	Anand	M	M	2025-10-07
98	11097	AW00011097	Edwin	Nara	M	M	2025-10-07
99	11098	AW00011098	Mallory	Rubio	S	F	2025-10-07
100	11099	AW00011099	Adam	Ross	M	M	2025-10-07
101	11100	AW00011100	Latasha	Navarro	S	F	2025-10-07
102	11101	AW00011101	Abby	Sai	S	F	2025-10-07
103	11102	AW00011102	Julia	Nelson	S	F	2025-10-07
104	11103	AW00011103	Cassie	Chande	S	F	2025-10-07
105	11104	AW00011104	Edgar	Sara	M	M	2025-10-07
106	11105	AW00011105	Candace	Fernandez	S	F	2025-10-07
107	11106	AW00011106	Jessie	Liu	S	M	2025-10-07

```

52  -- Data standardization and Consistency
53  SELECT DISTINCT
54    cst_marital_status,cst_gndr
55  FROM
56    bronze.crm_cust_info

```

100 %

Results Messages

	cst_marital_status	cst_gndr
1	M	M
2	S	M
3	S	F
4	S	NULL
5	M	F
6	M	NULL
7	NULL	NULL

In our project, we need to use **abbreviations** for values instead of single letters. Additionally, we must determine the appropriate way to handle **NULL values**, which will be finalized after discussions with the **source system manager**.

```

2  --query to standardize the low cardinality columns
3
4  SELECT
5    cst_id,
6    cst_key,
7    TRIM(cst_firstname) AS cst_firstname,
8    TRIM(cst_lastname) AS cst_lastname,
9    CASE
10      WHEN UPPER(cst_marital_status) = 'M' THEN 'Married'
11      WHEN UPPER(cst_marital_status) = 'S' THEN 'Single'
12      ELSE 'n/a'
13    END AS cst_marital_status,
14    CASE
15      WHEN UPPER(cst_gndr) = 'F' THEN 'Female'
16      WHEN UPPER(cst_gndr) = 'M' THEN 'Male'
17      ELSE 'n/a'
18    END AS cst_gndr
19  FROM
20  (SELECT
21    *,
22    ROW_NUMBER() OVER(PARTITION BY cst_id ORDER BY cst_create_date DESC) AS flag_latest
23    FROM
24      bronze.crm_cust_info
25    WHERE cst_id IS NOT NULL) t
26  WHERE flag_latest = 1
27

```

From the above code, we've completed cleaning the data for the table **bronze.crm\_cust\_info** and we can store it into **silver.load\_silver**

**Table : silver.crm\_prd\_info**

6 | SELECT TOP (20) \* FROM [bronze].[crm\_prd\_info]

100 %

Results Messages

	prd_id	prd_key	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt
1	210	CO-RF-FR-R92B-58	HL Road Frame - Black- 58	NULL	R	2003-07-01	NULL
2	211	CO-RF-FR-R92R-58	HL Road Frame - Red- 58	NULL	R	2003-07-01	NULL
3	212	AC-HE-HL-U509-R	Sport-100 Helmet- Red	12	S	2011-07-01	2007-12-28
4	213	AC-HE-HL-U509-R	Sport-100 Helmet- Red	14	S	2012-07-01	2008-12-27
5	214	AC-HE-HL-U509-R	Sport-100 Helmet- Red	13	S	2013-07-01	NULL
6	215	AC-HE-HL-U509	Sport-100 Helmet- Black	12	S	2011-07-01	2007-12-28
7	216	AC-HE-HL-U509	Sport-100 Helmet- Black	14	S	2012-07-01	2008-12-27
8	217	AC-HE-HL-U509	Sport-100 Helmet- Black	13	S	2013-07-01	NULL
9	218	CL-SO-SO-B909-M	Mountain Bike Socks- M	3	M	2011-07-01	2007-12-28
10	219	CL-SO-SO-B909-L	Mountain Bike Socks- L	3	M	2011-07-01	2007-12-28
11	220	AC-HE-HL-U509-B	Sport-100 Helmet- Blue	12	S	2011-07-01	2007-12-28
12	221	AC-HE-HL-U509-B	Sport-100 Helmet- Blue	14	S	2012-07-01	2008-12-27
13	222	AC-HE-HL-U509-B	Sport-100 Helmet- Blue	13	S	2013-07-01	NULL
14	223	CL-CA-CA-1098	AWC Logo Cap	6	S	2011-07-01	2007-12-28
15	224	CL-CA-CA-1098	AWC Logo Cap	5	S	2012-07-01	2008-12-27
16	225	CL-CA-CA-1098	AWC Logo Cap	7	S	2013-07-01	NULL
17	226	CL-JE-LJ-0192-S	Long-Sleeve Logo Jersey- S	32	S	2011-07-01	2007-12-28
18	227	CL-JE-LJ-0192-S	Long-Sleeve Logo Jersey- S	29	S	2012-07-01	2008-12-27
19	228	CL-JE-LJ-0192-S	Long-Sleeve Logo Jersey- S	38	S	2013-07-01	NULL
20	229	CL-JE-LJ-0192-M	Long-Sleeve Logo Jersey- M	32	S	2011-07-01	2007-12-28

prd\_id: It's primary key and we don't have any nulls or duplicates values in it.

```

95 -- Check for Nulls or Duplicates in Primary Key
96 -- Expectation: No Result
97
98
99 SELECT
100 prd_id,
101 COUNT(*)
102 FROM
103 bronze.crm_prd_info
104 GROUP BY prd_id
105 HAVING COUNT(*)>1 OR prd_id IS NULL
106

```

prd\_key: The prd\_key column contains metadata from another table, specifically **sales\_details (sls\_prd\_key)** and **px\_cat\_g1v2 (id)**. We need to split this column into two separate columns to correctly store and manage the respective data.

```

6 | SELECT TOP(20) * FROM [bronze].[crm_prd_info];
7 | SELECT TOP(20) * FROM [bronze].[crm_sales_details];
8 | SELECT TOP(20) * FROM [bronze].[erp_px_cat_g1v2];
9 | SELECT TOP(20) * FROM [silver].[erp_cust_az12];
10| SELECT TOP(20) * FROM [silver].[erp_loc_a101];
11|

```

100 %

Results Messages

	prd_id	prd_key	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt		
1	210	CO-RF-FR-R92B-58	HL Road Frame - Black- 58	NULL	R	2003-07-01	NULL		
2	211	CO-RF-FR-R92R-58	HL Road Frame - Red- 58	NULL	R	2003-07-01	NULL		
3	212	AC-HE-HL-U509-R	Sport-100 Helmet- Red	12	S	2011-07-01	2007-12-28		
4	213	AC-HE-HL-U509-R	Sport-100 Helmet- Red	14	S	2012-07-01	2008-12-27		
5	214	AC-HE-HL-U509-R	Sport-100 Helmet- Red	13	S	2013-07-01	NULL		
6	215	AC-HE-HL-U509	Sport-100 Helmet- Black	12	S	2011-07-01	2007-12-28		
7	216	AC-HE-HL-U509	Sport-100 Helmet- Black	14	S	2012-07-01	2008-12-27		
8	217	AC-HE-HL-U509	Sport-100 Helmet- Black	13	S	2013-07-01	NULL		
9	218	CL-SO-SO-B909-M	Mountain Bike Socks- M	3	M	2011-07-01	2007-12-28		
10	219	CL-SO-SO-B909-L	Mountain Bike Socks- L	3	M	2011-07-01	2007-12-28		
11	220	AC-HE-HL-U509-B	Sport-100 Helmet- Blue	12	S	2011-07-01	2007-12-28		
12	221	AC-HE-HL-U509-B	Sport-100 Helmet- Blue	14	S	2012-07-01	2008-12-27		
	sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_qtyunatly	sls_price
1	SO43697	BK-R93R-62	21768	20101229	20110105	20110110	3578	1	3578
2	SO43698	BK-M82S-44	28389	20101229	20110105	20110110	3400	1	3400
3	SO43699	BK-M82S-44	25863	20101229	20110105	20110110	3400	1	3400
4	SO43700	BK-R50B-62	14501	20101229	20110105	20110110	699	1	699
5	SO43701	BK-M82S-44	11003	20101229	20110105	20110110	3400	1	3400
6	SO43702	BK-R93R-44	27645	20101230	20110106	20110111	3578	1	3578
7	SO43703	BK-R93R-62	16624	20101230	20110106	20110111	3578	1	3578
8	SO43704	BK-M82B-48	11005	20101230	20110106	20110111	3375	1	3375
9	SO43705	BK-M82S-38	11011	20101230	20110106	20110111	3400	1	3400
10	SO43706	BK-R93R-48	27621	20101231	20110107	20110112	3578	1	3578
11	SO43707	BK-R93R-48	27616	20101231	20110107	20110112	3578	1	3578
12	SO43708	BK-R50R-52	20042	20101231	20110107	20110112	699	1	699
	id	cat	subcat	maintenance					
1	AC_BR	Accessories	Bike Racks	Yes					
2	AC_BS	Accessories	Bike Stands	No					
3	AC_BC	Accessories	Bottles and Cages	No					
4	AC_CL	Accessories	Cleaners	Yes					
5	AC_FE	Accessories	Fenders	No					
6	AC_HE	Accessories	Helmets	Yes					
7	AC_HP	Accessories	Hydration Packs	No					
8	AC_LI	Accessories	Lights	Yes					
9	AC_LO	Accessories	Locks	Yes					
10	AC_PA	Accessories	Panniers	No					
11	AC_PU	Accessories	Pumps	Yes					
12	AC_TT	Accessories	Tires and Tubes	Yes					

prd\_nm: It's a string column, so we've checked for unwanted spaces and we've found nothing.

```

108  -- Check for Nulls
109  -- Expectation: No Result
110
111  SELECT
112    prd_key
113  FROM bronze.crm_prd_info
114  WHERE prd_key IS NULL
115
116  -- splitting prd_key into two columns
117
118  SELECT
119    prd_id,
120    SUBSTRING(prd_key,7,LEN(prd_key)) AS prd_key,
121    REPLACE(SUBSTRING(prd_key,1,5),'-','_') AS cat_id,
122    prd_nm,
123    prd_cost,
124    prd_line,
125    prd_start_dt,
126    prd_end_dt
127  FROM bronze.crm_prd_info
128
129  -- Check for Unwanted Spaces for String Columns
130  -- Expectation: No Result
131
132  SELECT
133    prd_key
134  FROM bronze.crm_prd_info
135  WHERE prd_key != TRIM(prd_key)
136
137  SELECT
138    prd_key
139  FROM bronze.crm_prd_info
140  WHERE prd_nm != TRIM(prd_nm)
141

```

prd\_cost: It's a numeric column, so we need to check for negative values and nulls. We've found few null values here. We need to handle them after discussing with the team. Here, in our case we are making them zero.

```

142  -- Check for Nulls or Negative Numbers
143  -- Expectation: No Result
144
145  SELECT
146    prd_cost
147  FROM bronze.crm_prd_info
148  WHERE prd_cost IS NULL OR prd_cost < 0

```

**ISNULL()**

You can use COALESCE as well

Replaces NULL values with a specified replacement value

```

.50 --Replacing Nulls with 0
.51
.52 SELECT
.53 prd_id,
.54 SUBSTRING(prd_key,7,LEN(prd_key)) AS prd_key,
.55 REPLACE(SUBSTRING(prd_key,1,5) ,'-','_') AS cat_id,
.56 prd_nm,
.57 ISNULL(prd_cost,0) AS prd_cost,
.58 prd_line,
.59 prd_start_dt,
.60 prd_end_dt
.61 FROM bronze.crm_prd_info

```

prd\_line: It is low cardinality column. So, we need to fill with the abbreviation instead of single letter.

```

.63 -- Data standardization and Consistency
.64 SELECT DISTINCT
.65 prd_line
.66 FROM
.67 bronze.crm_prd_info
.68
.69 --Replacing with Abbrevations
.70 SELECT
.71 prd_id,
.72 SUBSTRING(prd_key,7,LEN(prd_key)) AS prd_key,
.73 REPLACE(SUBSTRING(prd_key,1,5) ,'-','_') AS cat_id,
.74 prd_nm,
.75 ISNULL(prd_cost,0) AS prd_cost,
.76 --CASE
.77     --WHEN UPPER(TRIM(prd_line)) = 'M' THEN 'Mountain'
.78     --WHEN UPPER(TRIM(prd_line)) = 'R' THEN 'Road'
.79     --WHEN UPPER(TRIM(prd_line)) = 'S' THEN 'Other Sales'
.80     --WHEN UPPER(TRIM(prd_line)) = 'T' THEN 'Touring'
.81     --ELSE 'n/a'
.82 --END AS prd_line,
.83 CASE UPPER(TRIM(prd_line))
.84     WHEN 'M' THEN 'Mountain'
.85     WHEN 'R' THEN 'Road'
.86     WHEN 'S' THEN 'Other Sales'
.87     WHEN 'T' THEN 'Touring'
.88     ELSE 'n/a'
.89 END AS prd_line,
.90 prd_start_dt,
.91 prd_end_dt
.92 FROM bronze.crm_prd_info
.93

```

prd\_start\_dt and prd\_end\_dt: We need to validate the prd\_end\_dt column to ensure that the end dates are not earlier than the corresponding prd\_start\_dt. Upon review, we found some records with this issue. To correct them, we will update the prd\_end\_dt by setting it to one day before the prd\_start\_dt of the same product. Before using it, we need to discuss with the team.

```

193
194 -- Check for Invalid Date Orders
195 -- End date must not be earlier than start date
196 SELECT *
197
198 FROM bronze.crm_prd_info
199 WHERE prd_end_dt < prd_start_dt
200
201
202 -- check for nulls in prd_start_dt
203
204 SELECT *
205
206 FROM bronze.crm_prd_info
207 WHERE prd_start_dt IS NULL

```

## #2 Solution

End Date = Start Date of the 'NEXT' Record -1

**LEAD ()**



Access values from the next row within a window

```

209 ---Replacing End date with day before of start date for each prd_key
210
211 SELECT
212     prd_id,
213     REPLACE(SUBSTRING(prd_key,1,5), ',', '-') AS cat_id,
214     SUBSTRING(prd_key,7,LEN(prd_key)) AS prd_key,
215     prd_nm,
216     ISNULL(prd_cost,0) AS prd_cost,
217     CASE UPPER(TRIM(prd_line))
218         WHEN 'M' THEN 'Mountain'
219         WHEN 'R' THEN 'Road'
220         WHEN 'S' THEN 'Other Sales'
221         WHEN 'T' THEN 'Touring'
222         ELSE 'n/a'
223     END AS prd_line,
224     prd_start_dt,
225     DATEADD(DAY, -1, LEAD(prd_start_dt) OVER (PARTITION BY prd_key ORDER BY prd_start_dt)) AS prd_end_dt
226
227 FROM bronze.crm_prd_info;

```

Now, we need to add cat\_id column to silver.crm\_prd\_info as we've derived this column from bronze.crm.prд\_info.

**Table : silver.crm\_sales\_details**

	sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_qty	sls_price
10	SO43706	BK-R93R-48	27621	20101231	20110107	20110112	3578	1	3578
11	SO43707	BK-R93R-48	27616	20101231	20110107	20110112	3578	1	3578
12	SO43708	BK-R50R-52	20042	20101231	20110107	20110112	699	1	699
13	SO43709	BK-R93R-52	16351	20101231	20110107	20110112	3578	1	3578
14	SO43710	BK-R93R-56	16517	20101231	20110107	20110112	3578	1	3578
15	SO43711	BK-R93R-56	27606	20110101	20110108	20110113	3578	1	3578
16	SO43712	BK-R93R-44	13513	20110101	20110108	20110113	3578	1	3578
17	SO43713	BK-R93R-62	27601	20110102	20110109	20110114	3578	1	3578

This table does not have a primary key, and some of the first three columns contain no null values. The next three columns are of type **INT** but represent **dates**, so we need to convert them into the **DATE** format.

```

247 -- checking for NULL
248 SELECT *
249
250 FROM bronze.crm_sales_details
251 WHERE sls_ord_num IS NULL
252
253 SELECT *
254
255 FROM bronze.crm_sales_details
256 WHERE sls_prd_key IS NULL
257
258 SELECT *
259
260 FROM bronze.crm_sales_details
261 WHERE sls_cust_id IS NULL
262

```

Negative numbers or zeros can't be cast to a date

```

264 -- check for Invalid Dates
265
266 SELECT sls_order_dt
267 FROM
268 bronze.crm_sales_details
269 WHERE sls_order_dt <= 0
270
271

```

	sls_order_dt
1	0
2	0
3	0
4	0
5	0
6	0
7	0
8	0
9	0
10	0

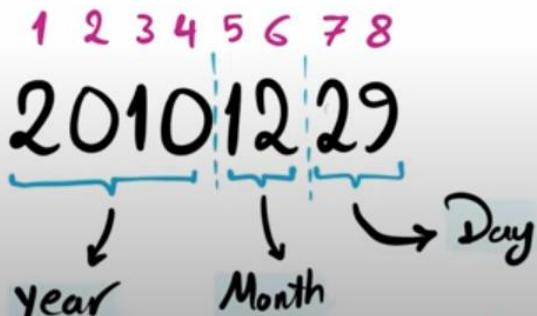
We've found some zero values and we need to replace them with NULL

### NULLIF()

Returns NULL if two given values are equal; otherwise, it returns the first expression.

In this scenario, the length of the date must be 8

	sls_order_dt
1	20101229
2	20101229
3	20101229
4	20101229
5	20101229
6	20101230
7	20101230
8	20101230
9	20101230
10	20101231



Check for outliers by validating the boundaries of the date range

We can't convert INT to DATE directly. So first we need to cast it to VARCHAR then DATE.

```
272 --Replacing with Nulls when dt <=0 or len(dt) != 8
273
274 SELECT
275     sls_ord_num,
276     sls_prd_key,
277     sls_cust_id,
278     CASE
279         WHEN sls_order_dt <= 0 OR LEN(sls_order_dt) != 8 THEN NULL
280         ELSE CAST(CAST(sls_order_dt AS VARCHAR) AS DATE)
281     END AS sls_order_dt,
282     CASE
283         WHEN sls_ship_dt <= 0 OR LEN(sls_ship_dt) != 8 THEN NULL
284         ELSE CAST(CAST(sls_ship_dt AS VARCHAR) AS DATE)
285     END AS sls_ship_dt,
286     CASE
287         WHEN sls_due_dt <= 0 OR LEN(sls_due_dt) != 8 THEN NULL
288         ELSE CAST(CAST(sls_due_dt AS VARCHAR) AS DATE)
289     END AS sls_due_dt,
290     sls_sales,
291     sls_qunatity,
292     sls_price
293     FROM bronze.crm_sales_details
294
295
```

Order Date must always be earlier than the Shipping Date or Due Date

```
294 ---Invalid order dates
295
296
297 SELECT *
298
299 FROM bronze.crm_sales_details
300 WHERE sls_order_dt > sls_ship_dt
301 OR sls_order_dt > sls_due_dt
```

100 %

Results Messages

sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_qtyunatly	sls_price
-------------	-------------	-------------	--------------	-------------	------------	-----------	---------------	-----------

We don't have any invalid order dates.

## Business Rules

$$\sum \text{Sales} = \text{Quantity} * \text{Price}$$

🚫 Negative, Zeros, Nulls are Not Allowed!

```

302
303 -- Data consistency with sales, quantity and price columns
304 -- sales = quantity * price
305 -- Values must not be zero, negative or null
306
307
308 SELECT
309 sls_sales,
310 sls_quantity,
311 sls_price
312 FROM bronze.crm_sales_details
313 WHERE sls_sales != sls_quantity * sls_price
314 OR sls_sales IS NULL OR sls_quantity IS NULL OR sls_price IS NULL
315 OR sls_sales <= 0 OR sls_quantity <= 0 OR sls_price <= 0
316

```

100 %

Results Messages

	sls_sales	sls_quantity	sls_price
1	10	2	NULL
2	25	5	NULL
3	70	2	NULL
4	9	1	NULL
5	35	1	NULL
6	100	10	NULL
7	16	2	NULL
8	769	1	-769
9	30	1	-30
10	22	1	-22
11	35	2	35
12	1701	1	-1701
13	21	1	-21
14	50	2	50
15	NULL	1	9
16	NULL	1	35
17	NULL	1	8
18	NULL	1	22
19	-18	1	9
20	0	1	10
21	50	2	50
22	NULL	1	10
23	5	2	5
24	64	4	64
25	25	3	25
26	40	1	2
27	2	1	50
28	33	1	24

## #1 Solution

Data Issues will be fixed direct in source system

## #2 Solution

~~Data Issues has to be fixed in data warehouse~~

### Rules

If Sales is negative, zero, or null, derive it using Quantity and Price.

If Price is zero or null, calculate it using Sales and Quantity.

If Price is negative, convert it to a positive value

### ABS ()

Returns absolute value of a number

```

318 ----fix
319 SELECT
320     sls_sales AS old_sales,
321     sls_qunatity,
322     sls_price AS old_prices,
323     CASE
324         WHEN sls_sales IS NULL
325             OR sls_sales <= 0
326             OR sls_sales != sls_qunatity * ABS(sls_price)
327             THEN sls_qunatity * ABS(sls_price)
328         ELSE sls_sales
329     END AS sls_sales,
330     CASE
331         WHEN sls_price IS NULL OR sls_price = 0
332             THEN sls_sales/NULLIF(sls_qunatity,0)
333         WHEN sls_price < 0 THEN ABS(sls_price)
334         ELSE sls_price
335     END AS sls_price
336     FROM bronze.crm_sales_details
337     WHERE sls_sales != sls_qunatity * sls_price
338     OR sls_sales IS NULL OR sls_qunatity IS NULL OR sls_price IS NULL
339     OR sls_sales <= 0 OR sls_qunatity <= 0 OR sls_price <= 0
340
341

```

100 %

Results Messages

	old_sales	sls_qunatity	old_prices	sls_sales	sls_price
8	769	1	-769	769	769
9	30	1	-30	30	30
10	22	1	-22	22	22
11	35	2	35	70	35
12	1701	1	-1701	1701	1701
13	21	1	-21	21	21
14	50	2	50	100	50
15	NULL	1	9	9	9
16	NULL	1	35	35	35
17	NULL	1	8	8	8
18	NULL	1	22	22	22
19	-18	1	9	9	9
20	0	1	10	10	10
21	50	2	50	100	50
22	NULL	1	10	10	10
23	5	2	5	10	5
24	64	4	64	256	64
25	25	3	25	75	25

```

342 -----Corrected Query for the whole table
343
344 SELECT
345     sls_ord_num,
346     sls_prd_key,
347     sls_cust_id,
348     CASE
349         WHEN sls_order_dt <= 0 OR LEN(sls_order_dt) != 8 THEN NULL
350         ELSE CAST(CAST(sls_order_dt AS VARCHAR) AS DATE)
351     END AS sls_order_dt,
352     CASE
353         WHEN sls_ship_dt <= 0 OR LEN(sls_ship_dt) != 8 THEN NULL
354         ELSE CAST(CAST(sls_ship_dt AS VARCHAR) AS DATE)
355     END AS sls_ship_dt,
356     CASE
357         WHEN sls_due_dt <= 0 OR LEN(sls_due_dt) != 8 THEN NULL
358         ELSE CAST(CAST(sls_due_dt AS VARCHAR) AS DATE)
359     END AS sls_due_dt,
360     CASE
361         WHEN sls_sales IS NULL
362             OR sls_sales <= 0
363             OR sls_sales != sls_qunatity * ABS(sls_price)
364             THEN sls_qunatity * ABS(sls_price)
365             ELSE sls_sales
366     END AS sls_sales,
367     sls_qunatity,
368     CASE
369         WHEN sls_price IS NULL OR sls_price = 0
370             THEN sls_sales/NULLIF(sls_qunatity,0)
371         WHEN sls_price < 0 THEN ABS(sls_price)
372         ELSE sls_price
373     END AS sls_price
374 FROM bronze.crm_sales_details

```

Table : silver.erp\_cust\_az12

	cid	bdate	gen
1	NASA00011000	1971-10-06	Male
2	NASA00011001	1976-05-10	Male
3	NASA00011002	1971-02-09	Male
4	NASA00011003	1973-08-14	Female
5	NASA00011004	1979-08-05	Female
6	NASA00011005	1976-08-01	Male
7	NASA00011006	1976-12-02	Female
8	NASA00011007	1969-11-06	Male
9	NASA00011008	1975-07-04	Female
10	NASA00011009	1969-09-29	Male
11	NASA00011010	1969-08-05	Female
12	NASA00011011	1969-05-03	Male
13	NASA00011012	1979-01-14	Female
14	NASA00011013	1979-08-03	Male
15	NASA00011014	1973-11-06	Female
16	NASA00011015	1984-08-26	Female
17	NASA00011016	1984-10-25	Male
18	NASA00011017	1949-12-24	Female
19	NASA00011018	1955-10-06	Male
20	NASA00011019	1983-09-04	Male

We have a **cid** column like the **cst\_key** in **bronze.crm\_cust\_info**, but in the current table, it includes a "NAS" prefix. We need to remove this prefix using the **SUBSTRING** function.

	cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date
1	11000	AW00011000	Jon	Yang	M	M	2025-10-06
2	11001	AW00011001	Eugene	Huang	S	M	2025-10-06
3	11002	AW00011002	Ruben	Torres	M	M	2025-10-06
4	11003	AW00011003	Christy	Zhu	S	F	2025-10-06

We've found few birth days those are ahead of current time, and we can replace this values with Nulls.

```

452
453 -- Out of range birsth day dates
454
455 SELECT
456 bdate
457 FROM bronze.erp_cust_az12
458 WHERE bdate > GETDATE()

```

	bdate
1	2050-07-06
2	2042-02-22
3	2050-05-21
4	2038-10-17
5	2045-03-03
6	2050-11-22
7	2066-06-16
8	9999-09-13
9	2065-12-12
10	9999-11-20
11	9999-09-11
12	9999-05-10
13	2050-09-07
14	2080-03-15
15	2055-01-23
16	2980-03-09

```

460 --fixing bday dates
461
462 SELECT
463 CASE
464     WHEN cid LIKE 'NAS%' THEN SUBSTRING(cid,4,LEN(cid))
465     ELSE cid
466 END AS cid,
467 CASE
468     WHEN bdate > GETDATE() THEN NULL
469     ELSE bdate
470 END AS bdate,
471 gen
472 FROM bronze.erp_cust_az12

```

```

475
476     ----Data Standardization and consistency
477
478     SELECT
479         DISTINCT gen
480     FROM    bronze.erp_cust_az12

```

Results

gen	
1	NULL
2	F
3	
4	Male
5	Female
6	M

In gender column, we've found inconsistencies in the data, so we need to standardize it.

```

482     --fixing standardization and whole query
483
484     SELECT
485         CASE
486             WHEN cid LIKE 'NAS%' THEN SUBSTRING(cid,4,LEN(cid))
487             ELSE cid
488         END AS cid,
489         CASE
490             WHEN TRY_CAST(bdate AS DATE) > GETDATE() THEN NULL -- WE don't get bdate always as DATE
491             ELSE bdate
492         END AS bdate,
493         CASE
494             WHEN UPPER(TRIM(gen)) IN ('M','Male') THEN 'Male'
495             WHEN UPPER(TRIM(gen)) IN ('F','Female') THEN 'Female'
496             ELSE 'n/a'
497         END gen
498     FROM bronze.erp_cust_az12
499

```

**Table : silver.erp\_loc\_a101**

```

9      SELECT TOP(20) * FROM [bronze].[erp_loc_a101];
10     SELECT TOP(20) * FROM [bronze].[erp_px_cat_g1v2];
11

```

Results

cid	cntry
1	Australia
2	Australia
3	Australia
4	Australia
5	Australia
6	Australia
7	Australia
8	Australia
9	Australia
10	Australia
11	Australia
12	Australia
13	US
14	US
15	US
16	US
17	US
18	Australia
19	Australia
20	Canada

```

5  SELECT TOP(5) * FROM [bronze].[crm_cust_info];
6  SELECT TOP(5) * FROM [bronze].[erp_loc_a101];
100 %

```

Results

	cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date
1	11000	AW00011000	Jon	Yang	M	M	2025-10-06
2	11001	AW00011001	Eugene	Huang	S	M	2025-10-06
3	11002	AW00011002	Ruben	Torres	M	M	2025-10-06
4	11003	AW00011003	Christy	Zhu	S	F	2025-10-06
5	11004	AW00011004	Elizabeth	Johnson	S	F	2025-10-06

	cid	cntry
1	AW-00011000	Australia
2	AW-00011001	Australia
3	AW-00011002	Australia
4	AW-00011003	Australia
5	AW-00011004	Australia

In the **cid** column of the **bronze.erp\_loc\_a101** table, we need to remove the hyphen (-) to ensure it matches the **cst\_key** column in the **bronze.crm\_cust\_info** table.

We don't have any NULL values for cid column

```

504
505  -- checking for NULL or Duplicates
506
507  SELECT
508  cid,
509  --cntry
510  COUNT(*)
511  FROM
512  bronze.erp_loc_a101
513  GROUP BY cid
514  HAVING COUNT(*) > 1 OR cid IS NULL
515

```

We've found inconsistencies in the cntry column

```

516
517  --Data Standardization
518  SELECT
519  DISTINCT cntry
520  FROM
521  bronze.erp_loc_a101
522

```

Results

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

```

2   -- fixing standardization and query for whole table
3
4   SELECT
5     REPLACE(cid, '-', '') AS cid,
6     CASE
7       WHEN UPPER(TRIM(cntry)) = 'DE' THEN 'Germany'
8       WHEN UPPER(TRIM(cntry)) IN ('US', 'USA') THEN 'United States'
9       WHEN cntry IS NULL OR cntry = '' THEN 'n/a'
10      ELSE TRIM(cntry)
11    END AS cntry
12  FROM
13  bronze.erp_loc_a101

```

Table : silver.erp\_px\_cat\_g1v2

	id	cat	subcat	maintenance
1	AC_BR	Accessories	Bike Racks	Yes
2	AC_BS	Accessories	Bike Stands	No
3	AC_BC	Accessories	Bottles and Cages	No
4	AC_CL	Accessories	Cleaners	Yes
5	AC_FE	Accessories	Fenders	No
6	AC_HE	Accessories	Helmets	Yes
7	AC_HP	Accessories	Hydration Packs	No
8	AC_LI	Accessories	Lights	Yes
9	AC_LO	Accessories	Locks	Yes
10	AC_PA	Accessories	Panniers	No
11	AC_PU	Accessories	Pumps	Yes
12	AC_TT	Accessories	Tires and Tubes	Yes
13	BI_MB	Bikes	Mountain Bikes	Yes
14	BI_RB	Bikes	Road Bikes	Yes
15	BI_TB	Bikes	Touring Bikes	Yes
16	CL_BS	Clothing	Bib-Shorts	No
17	CL_CA	Clothing	Caps	No
18	CL_GL	Clothing	Gloves	No
19	CL_JE	Clothing	Jerseys	No
20	CL_SH	Clothing	Shorts	No

```

558
559  --- check for unwanted spaces
560
561  SELECT
562
563  FROM bronze.erp_px_cat_g1v2
564  WHERE cat != TRIM(cat) OR subcat != TRIM(subcat) OR maintenance != TRIM(maintenance)
565

```

	id	cat	subcat	maintenance

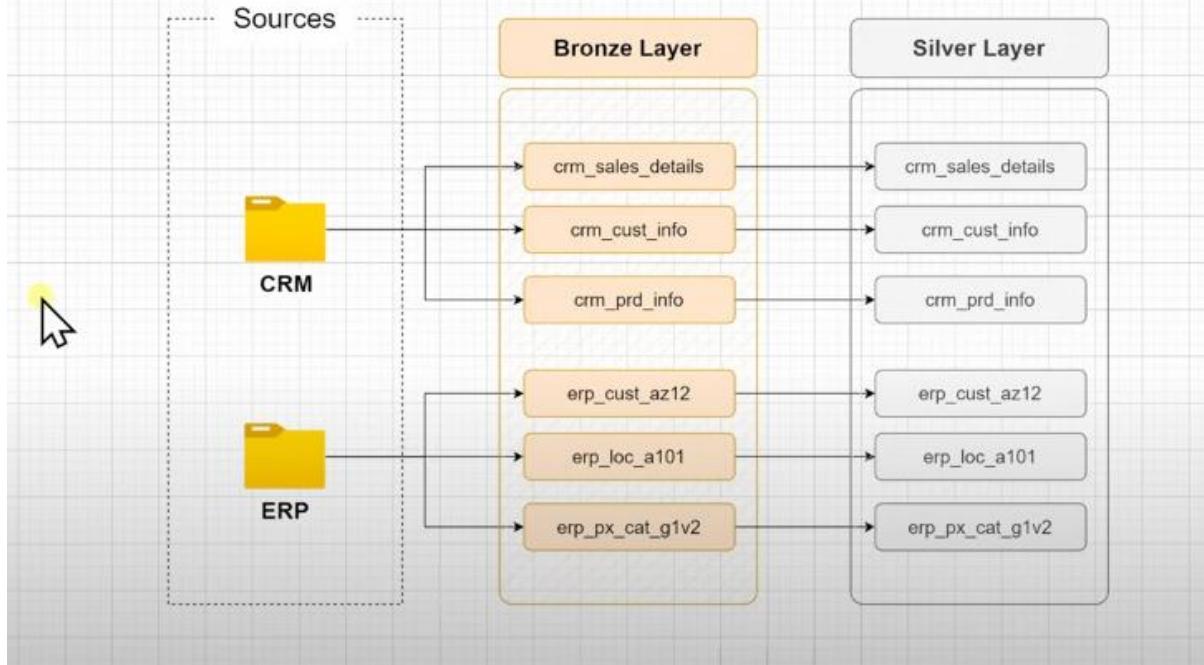
```

566  -- Data Standardization and Consistency
567
568  SELECT
569    DISTINCT
570    cat
571  FROM
572  bronze.erp_px_cat_g1v2;
573
574  SELECT
575    DISTINCT |
576    subcat
577  FROM bronze.erp_px_cat_g1v2;
578

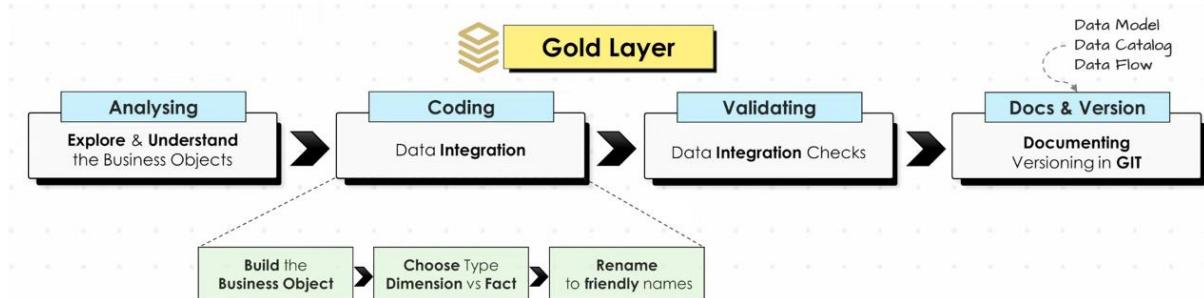
```

We don't need to do anything on this table as we've nice data quality for this table

# Data Flow Diagram



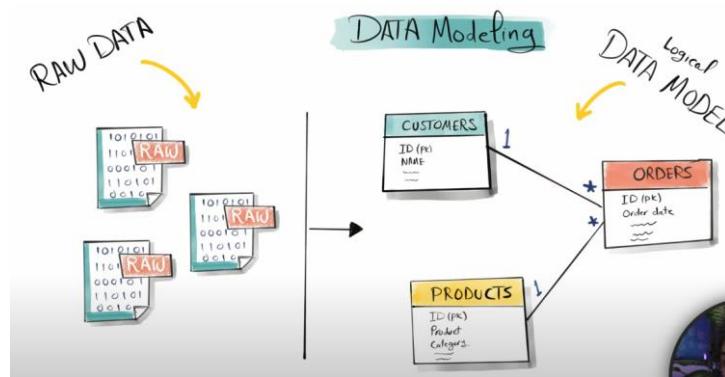
## Gold Layer



# Data Modelling

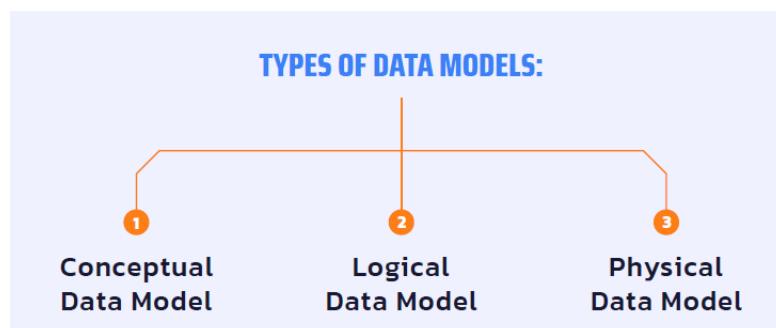
## Data Modelling:

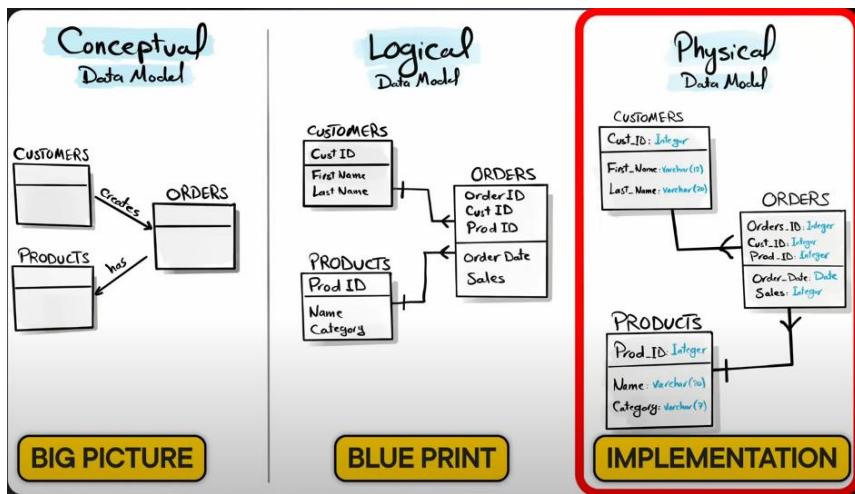
Data modelling is the process of creating a conceptual or logical representation of data structures, their relationships, and the rules governing them. It is widely used in designing databases, systems, or applications to organize and manage data effectively.



### A data model consists of 3 components

- **Entities:**  
These are the main things we want to store information about. For instance, in a business, an entity could be "Customer" or "Product."
- **Attributes:**  
These are the specific pieces of information about an entity. For a "Customer," attributes could include name, address, and phone number.
- **Relationships:**  
These define how entities are connected or related to each other. For example, a "Customer" can have a relationship with an "Order."





## 1 CONCEPTUAL DATA MODEL

A high-level view of what needs to be stored and how different entities relate to each other. It's like a bird's eye view.

## 2 LOGICAL DATA MODEL

More detailed than the conceptual model, specifying attributes and relationships. It's like a floor plan of a house.

## 3 PHYSICAL DATA MODEL

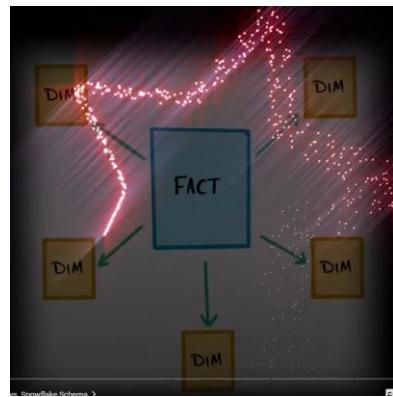
It specifies how the data will be stored, considering database technologies and constraints. It's like the actual construction of the house.

**Note:** In most cases, we work with conceptual or logical data models, as building a physical model requires significant effort. Additionally, tools like Databricks and Power BI allow us to create and manage physical data models efficiently. In our project, we primarily use a logical data model.

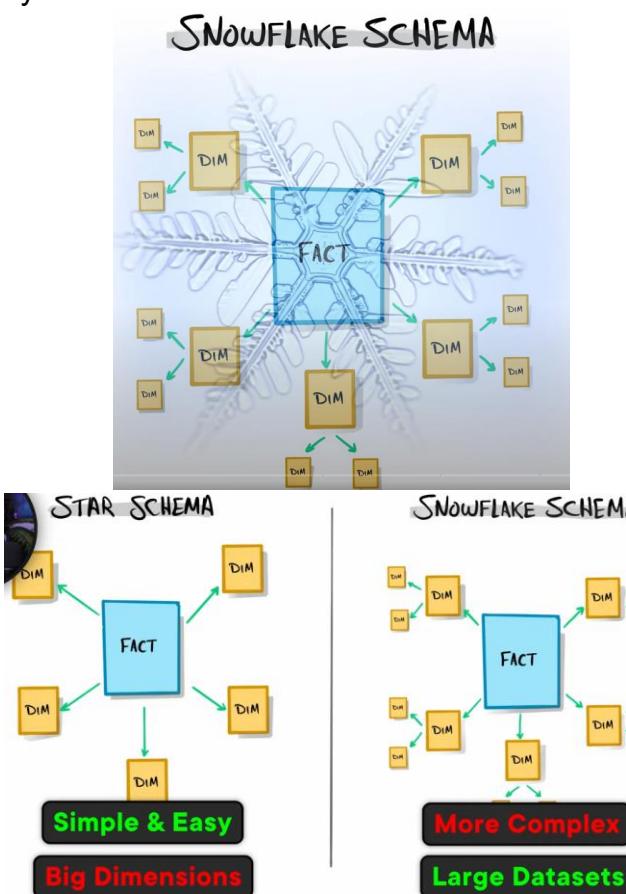
We've two special data models for,

- 1)optimized for reporting
- 2)Flexible
- Easy to understand

1) **Star schema:** In a star schema, the FACT table is centrally positioned and surrounded by DIMENSION (DIM) tables. The DIM tables store descriptive information, providing context to the data, while the FACT table contains measurable values such as quantity and sales.

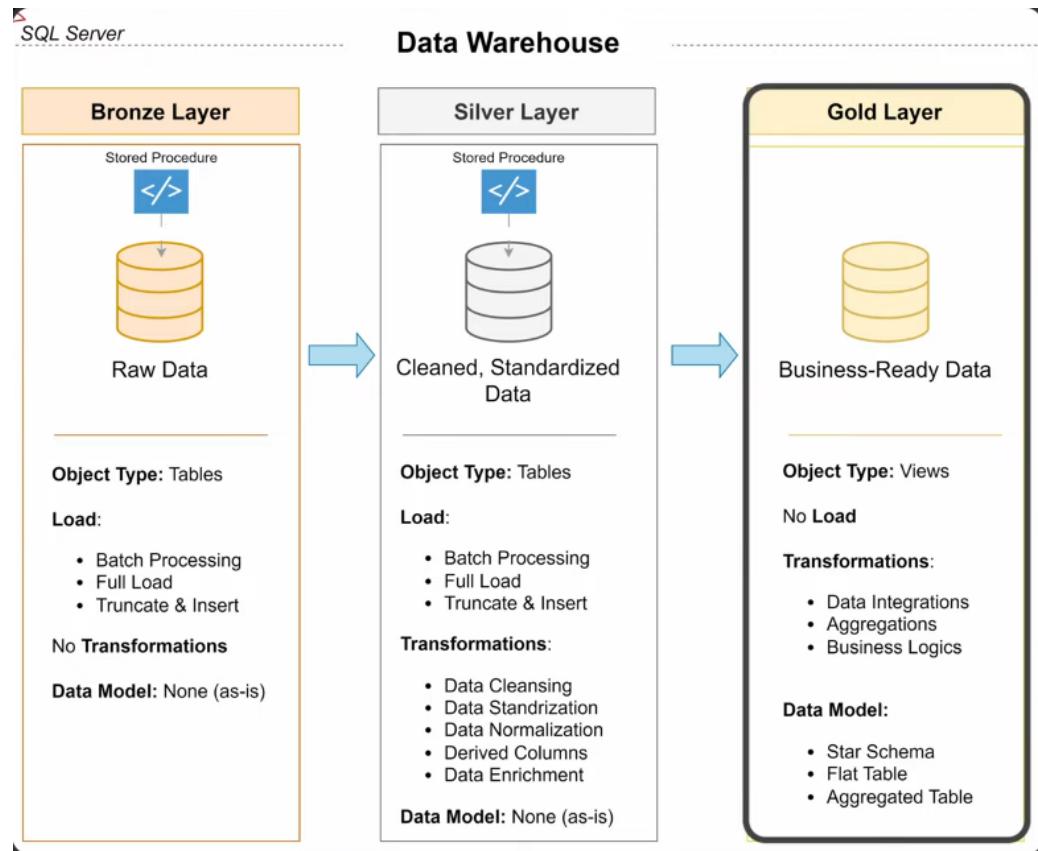


- 2) **Snowflake Schema:** This is like a star schema, but the DIMENSION (DIM) tables are further normalized to reduce data redundancy and optimize storage efficiency.



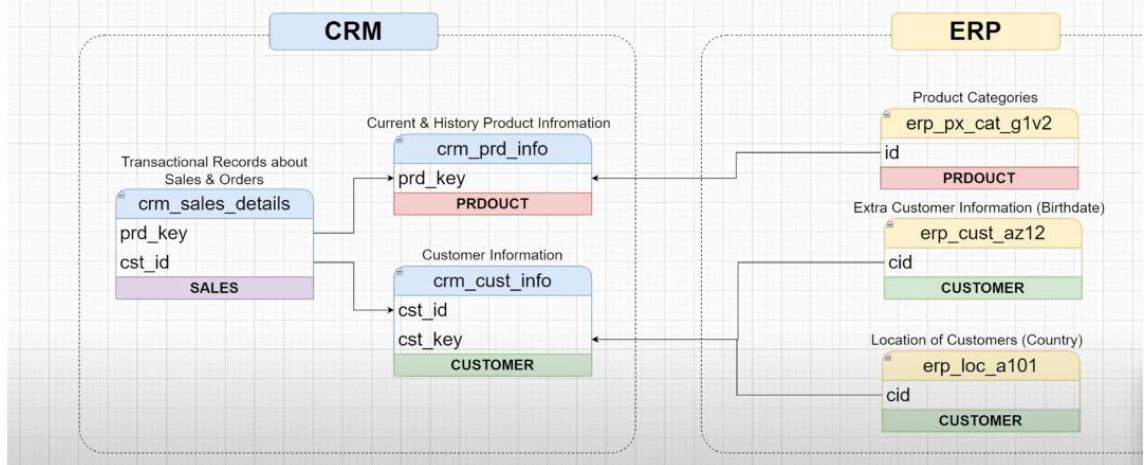
In our project, we've chosen star schema as it is easy for reporting and popular

DIMENSION	FACT
Descriptive information that give context to your data.  Who? What? Where?	Quantitative information that represents events  How much? How many?



### Integration Model

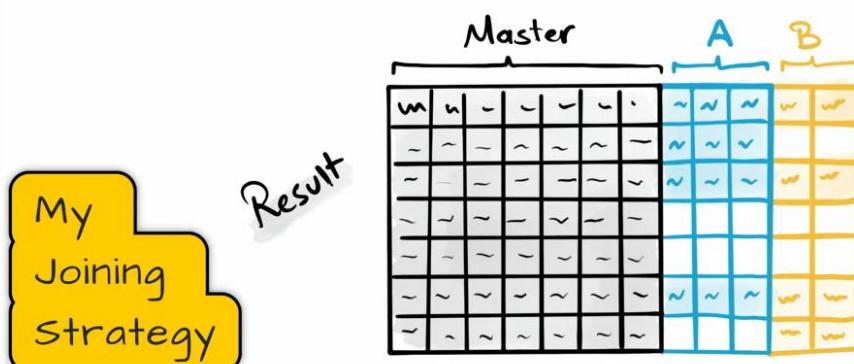
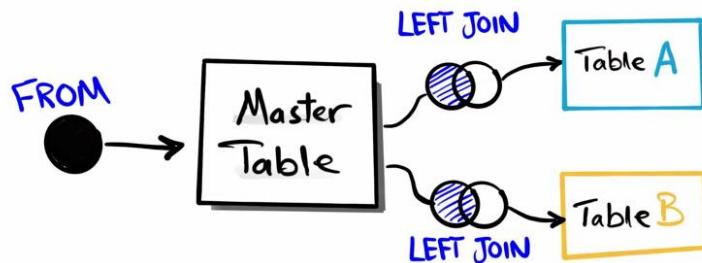
(how to tables are related)



## Build Gold Layer

### Create Dimension Customers

	Bronze Layer	Silver Layer	Gold Layer
Objective	Raw, unprocessed data as-is from sources	Clean & standardized data (Intermediate Layer) Prepare Data for Analysis	Business-Ready data Provide data to be consumed for reporting & Analytics
Object Type	Tables	Tables	Views
Load Method	Full Load (Truncate & Insert)	Full Load (Truncate & Insert)	None
Data Transformation	None (as-is)	<ul style="list-style-type: none"> <li>- Data Cleaning</li> <li>- Data Standardization</li> <li>- Data Normalization</li> <li>- Derived Columns</li> <li>- Data Enrichment</li> </ul>	<ul style="list-style-type: none"> <li>- Data Integration</li> <li>- Data Aggregation</li> <li>- Business Logic &amp; Rules</li> </ul>
Data Modeling	None (as-is)	None (as-is)	<ul style="list-style-type: none"> <li>- Start Schema</li> <li>- Aggregated Objects</li> <li>- Flat Tables</li> </ul>
Target Audience	- Data Engineers	<ul style="list-style-type: none"> <li>- Data Analysts</li> <li>- Data Engineers</li> </ul>	<ul style="list-style-type: none"> <li>- Data Analysts</li> <li>- Business Users</li> </ul> 



**NOTE:** It is preferable to use a LEFT JOIN instead of an INNER JOIN to ensure that we do not lose any data from the primary table, even if there are no matching records in the joined table.

**TIP** IN silver.emp\_cust\_az12 ca



After Joining table, check if any duplicates were introduced by the join logic

ON c1.cst\_key = ta.cst\_id

```

7 | SELECT
8 | ci.cst_id,
9 | ci.cst_key,
10| ci.cst_firstname,
11| ci.cst_lastname,
12| ci.cst_marital_status,
13| ci.cst_gndr,
14| ci.cst_create_date,
15| ca.bdate,
16| ca.gen,
17| la.cntry
18| FROM silver.crm_cust_info ci
19| LEFT JOIN silver.erp_cust_az12 ca
20| ON ci.cst_key = ca.cid
21| LEFT JOIN silver.erp_loc_a101 la
22| ON ci.cst_key = la.cid
23|

```

100 %

Results Messages

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date	bdate	gen	cntry
1 11000	AW00011000	Jon	Yang	Married	Male	2025-10-06	1971-10-06	Male	Australia
2 11001	AW00011001	Eugene	Huang	Single	Male	2025-10-06	1976-05-10	Male	Australia
3 11002	AW00011002	Ruben	Torres	Married	Male	2025-10-06	1971-02-09	Male	Australia
4 11003	AW00011003	Christy	Zhu	Single	Female	2025-10-06	1973-08-14	Female	Australia
5 11004	AW00011004	Elizabeth	Johnson	Single	Female	2025-10-06	1979-08-05	Female	Australia
6 11005	AW00011005	Julio	Ruiz	Single	Male	2025-10-06	1976-08-01	Male	Australia

```

3 | SELECT cst_id,COUNT(*) FROM
4 | (
5 | SELECT
6 | ci.cst_id,
7 | ci.cst_key,
8 | ci.cst_firstname,
9 | ci.cst_lastname,
10| ci.cst_marital_status,
11| ci.cst_gndr,
12| ci.cst_create_date,
13| ca.bdate,
14| ca.gen,
15| la.cntry
16| FROM silver.crm_cust_info ci
17| LEFT JOIN silver.erp_cust_az12 ca
18| ON ci.cst_key = ca.cid
19| LEFT JOIN silver.erp_loc_a101 la
20| ON ci.cst_key = la.cid
21| )t
22| GROUP BY cst_id
23| HAVING COUNT(*)>1

```

100 %

Results Messages

cst_id	(No column name)
--------	------------------

**We don't have any duplicates after joining multiple tables.**

cst_id	cst_key	cst_firstname	cst_lastname	cst_marital_status	cst_gndr	cst_create_date	bdate	gen	cntry
11000	AW00011000	Jon	Yang	Married	Male	2025-10-06	1971-10-06	Male	Australia
11001	AW00011001	Eugene	Huang	Single	Male	2025-10-06	1976-05-10	Male	Australia

Currently, we have multiple gender columns from different tables. To ensure consistency and avoid redundancy, we need to integrate the data and consolidate it into a single gender column.

```

29  SELECT
30  DISTINCT
31  ci.cst_gndr,
32  ca.gen
33  FROM silver.crm_cust_info ci
34  LEFT JOIN silver.erp_cust_az12 ca
35  ON      ci.cst_key = ca.cid
36  LEFT JOIN silver.erp_loc_a101 la
37  ON      ci.cst_key = la.cid
38  ORDER BY 1,2
    
```

The screenshot shows the SQL query in the top pane and its results in the bottom pane. The results table has two columns: 'cst\_gndr' and 'gen'. The 'gen' column contains several NULL values. A callout bubble points to one of these NULL values with the text: 'NULLs often come from joined tables! NULL will appear if SQL finds no match'.

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

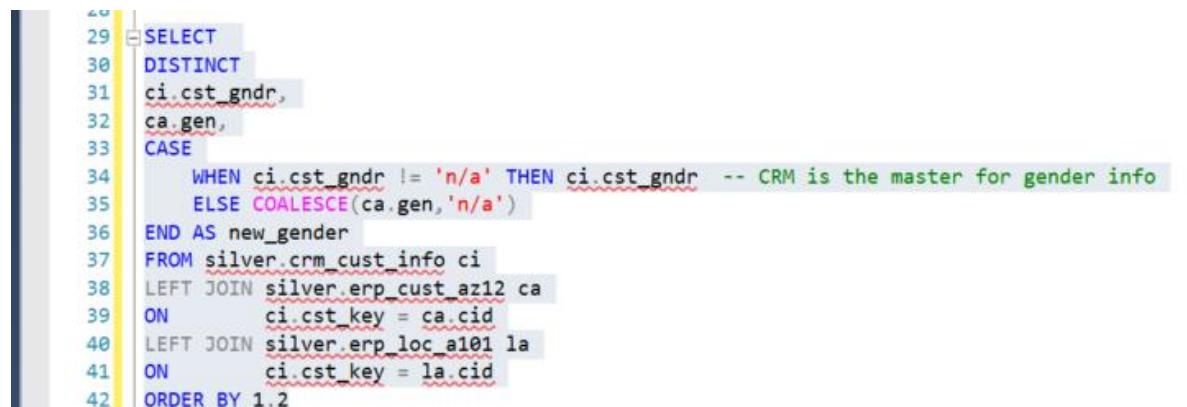
cst_gndr	gen
1	Female
2	Male
3	n/a
4	Female
5	Male
6	n/a
7	NULL

cst_gndr	gen
1	Female
2	Male
3	n/a
4	Female
5	Male

1) We found NULL values after joining the tables, even though NULLs were already handled in the silver table. These NULL values are not originating from the silver table but are instead a result of the join operation. When there are non-matching records between the tables, NULL values appear in the output.

2) In some cases, one column contains "Female" while another contains "Male," making it unclear which value to use. To resolve this, we need to confirm with the source administrator and choose the appropriate column. In our project, we are using **cst\_gndr** from **silver.crm\_cust\_info**.



```
29  SELECT
30    DISTINCT
31    ci.cst_gndr,
32    ca.gen,
33    CASE
34      WHEN ci.cst_gndr != 'n/a' THEN ci.cst_gndr -- CRM is the master for gender info
35      ELSE COALESCE(ca.gen, 'n/a')
36    END AS new_gender
37  FROM silver.crm_cust_info ci
38  LEFT JOIN silver.erp_cust_az12 ca
39  ON      ci.cst_key = ca.cid
40  LEFT JOIN silver.erp_loc_a101 la
41  ON      ci.cst_key = la.cid
42  ORDER BY 1,2
```



	cst_gndr	gen	new_gender
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

A B C

Rename columns to friendly, meaningful names

## General Principles

- Naming Conventions:** Use snake\_case, with lowercase letters and underscores ( \_ ) to separate words.
- Language:** Use English for all names.
- Avoid Reserved Words:** Do not use SQL reserved words as object names.

la.cntry AS country  
CRM silver.crm\_cust\_info ci

sort the columns into logical groups to improve readability

```

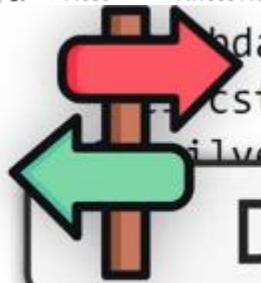
7  SELECT
8    ci.cst_id AS customer_id,
9    ci.cst_key AS customer_number,
10   ci.cst_firstname AS first_name,
11   ci.cst_lastname AS last_name,
12   la.cntry AS country,
13   ci.cst_marital_status AS marital_status,
14  CASE
15    WHEN ci.cst_gndr != 'n/a' THEN ci.cst_gndr -- CRM is the master for gender info
16    ELSE COALESCE(ca.gen, 'n/a')
17  END AS gender,
18  ca.bdate AS birthdate,
19  ci.cst_create_date AS create_date
20  FROM silver.crm_cust_info ci
21  LEFT JOIN silver.erp_cust_az12 ca
22  ON      ci.cst_key = ca.cid
23  LEFT JOIN silver.erp_loc_a101 la
24  ON      ci.cst_key = la.cid
25
26

```

100 %

Results Messages

	customer_id	customer_number	first_name	last_name	country	marital_status	gender	birthdate	create_date
79	11078	AW00011078	Gina	Martin	Canada	Single	Female	1985-01-06	2025-10-07
80	11079	AW00011079	Donald	Gonzalez	Australia	Single	Male	1970-03-07	2025-10-07
81	11080	AW00011080	Damien	Chander	Australia	Married	Male	1965-01-13	2025-10-07
82	11081	AW00011081	Savannah	Baker	United States	Married	Female	1972-01-21	2025-10-07
83	11082	AW00011082	Angela	Butler	United States	Single	Female	1972-02-01	2025-10-07
84	11083	AW00011083	Alyssa	Cox	United States	Married	Female	1977-03-11	2025-10-07
85	11084	AW00011084	Lucas	Phillips	United States	Single	Male	1963-03-12	2025-10-07
86	11085	AW00011085	Emily	Johnson	United States	Single	Female	1963-01-16	2025-10-07
87	11086	AW00011086	Ryan	Brown	United States	Married	Male	1963-06-22	2025-10-07



birthdate AS birthdate,  
 ... cst\_create\_date AS create\_date  
 silver.crm\_cust\_info ci

## Dimension vs Fact ???

We need to determine the object type—whether it is a Dimension or a Fact table. Since it contains descriptive data, our object type (customers) is classified as a Dimension Table.

When creating a new Dimension table, it is essential to have a primary key. While we can use the primary key from the source system, there are instances where a reliable primary key is unavailable. In such cases, we generate a primary key within the data warehouse, known as a **Surrogate Key**.

A **Surrogate Key** is not a business key; it holds no intrinsic meaning, and no one in the business refers to it. It is solely used to establish relationships within the data model, giving us more control over data integration without relying entirely on the source system.

And there are different ways on how to generate surrogate keys



## Surrogate Key

System-generated unique identifier assigned to each record in a table.

- DDL-based generation.
- Query-based using Window function (Row\_Number)

We've decided to create views as the objects in the gold layer.

```
6  CREATE VIEW gold.dim_customers AS
7  SELECT
8      ROW_NUMBER() OVER (ORDER BY cst_id) AS customer_key, ---we can use any column for creating the Surrogate Key
9      ci.cst_id AS customer_id,
10     ci.cst_key AS customer_number,
11     ci.cst_firstname AS first_name,
12     ci.cst_lastname AS last_name,
13     la.cntry AS country,
14     ci.cst_marital_status AS marital_status,
15     CASE
16         WHEN ci.cst_gndr != 'n/a' THEN ci.cst_gndr -- CRM is the master for gender info
17         ELSE COALESCE(ca.gen,'n/a')
18     END AS gender,
19     ca.bdate AS birthdate,
20     ci.cst_create_date AS create_date
21     FROM silver.crm_cust_info ci
22     LEFT JOIN silver.erp_cust_az12 ca
23     ON      ci.cst_key = ca.cid
24     LEFT JOIN silver.erp_loc_a101 la
25     ON      ci.cst_key = la.cid
26
27     SELECT * FROM gold.dim_customers
28
```

Results

	customer_key	customer_id	customer_number	first_name	last_name	country	marital_status	gender	birthdate	create_date
1	1	11000	AW00011000	Jon	Yang	Australia	Married	Male	1971-10-06	2025-10-06
2	2	11001	AW00011001	Eugene	Huang	Australia	Single	Male	1976-05-10	2025-10-06
3	3	11002	AW00011002	Ruben	Torres	Australia	Married	Male	1971-02-09	2025-10-06
4	4	11003	AW00011003	Christy	Zhu	Australia	Single	Female	1973-08-14	2025-10-06
5	5	11004	AW00011004	Elizabeth	Johnson	Australia	Single	Female	1979-08-05	2025-10-06
6	6	11005	AW00011005	Julio	Ruiz	Australia	Single	Male	1976-08-01	2025-10-06

# Build Gold Layer

## Create Dimension Products

```

72  SELECT *
73  FROM
74  silver.crm_prd_info pn
75  LEFT JOIN silver.emp_px_cat_g1v2 pc
76  ON pn.cat_id = pc.id
77
78
79
80
  
```

Results

prd_id	cat_id	prd_key	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt	dwh_create_date	id	cat	subcat	maintenance	dwh_create_date	
1	AC_BC	BC-M005	Mountain Bottle Cage	4	Mountain	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BC	Accessories	Bottles and Cages	No	2025-03-15 13:28:02.690000	
2	479	AC_BC	BC-R205	Road Bottle Cage	3	Road	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BC	Accessories	Bottles and Cages	No	2025-03-15 13:28:02.690000
3	477	AC_BC	WB-H098	Water Bottle - 30 oz.	2	Other Sales	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BC	Accessories	Bottles and Cages	No	2025-03-15 13:28:02.690000
4	483	AC_BR	RA-H123	Hitch Rack - 4-Bike	45	Other Sales	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BR	Accessories	Bike Racks	Yes	2025-03-15 13:28:02.690000
5	486	AC_BS	ST-1401	All-Purpose Bike Stand	59	Mountain	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BS	Accessories	Bike Stands	No	2025-03-15 13:28:02.690000
6	484	AC_CL	CL-9009	Bike Wash - Dissolver	3	Other Sales	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_CL	Accessories	Cleaners	Yes	2025-03-15 13:28:02.690000
7	485	AC_FE	FE-6654	Fender Set - Mountain	8	Mountain	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_FE	Accessories	Fenders	No	2025-03-15 13:28:02.690000
8	215	AC_HF	HI-LI509	Sport100 Helmet - Black	12	Other Sales	2011-07-01	2012-06-30	2025-03-15 13:28:02.403333	AC_HF	Accessories	Helmets	Yes	2025-03-15 13:28:02.690000

We want to retain only currently available products and filter out historical products. We can identify active products by checking if **prd\_end\_dt** is **NULL**, as this indicates they are still available.

```

72  SELECT *
73  FROM
74  silver.crm_prd_info pn
75  LEFT JOIN silver.emp_px_cat_g1v2 pc
76  ON pn.cat_id = pc.id
77  WHERE pn.prd_end_dt IS NULL --- Filter out historical data
78
79
80
  
```

Results

prd_id	cat_id	prd_key	prd_nm	prd_cost	prd_line	prd_start_dt	prd_end_dt	dwh_create_date	id	cat	subcat	maintenance	dwh_create_date	
1	478	AC_BC	BC-M005	Mountain Bottle Cage	4	Mountain	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BC	Accessories	Bottles and Cages	No	2025-03-15 13:28:02.690000
2	479	AC_BC	BC-R205	Road Bottle Cage	3	Road	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BC	Accessories	Bottles and Cages	No	2025-03-15 13:28:02.690000
3	477	AC_BC	WB-H098	Water Bottle - 30 oz.	2	Other Sales	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BC	Accessories	Bottles and Cages	No	2025-03-15 13:28:02.690000
4	483	AC_BR	RA-H123	Hitch Rack - 4-Bike	45	Other Sales	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BR	Accessories	Bike Racks	Yes	2025-03-15 13:28:02.690000
5	486	AC_BS	ST-1401	All-Purpose Bike Stand	59	Mountain	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_BS	Accessories	Bike Stands	No	2025-03-15 13:28:02.690000
6	484	AC_CL	CL-9009	Bike Wash - Dissolver	3	Other Sales	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_CL	Accessories	Cleaners	Yes	2025-03-15 13:28:02.690000
7	485	AC_FE	FE-6654	Fender Set - Mountain	8	Mountain	2013-07-01	NULL	2025-03-15 13:28:02.403333	AC_FE	Accessories	Fenders	No	2025-03-15 13:28:02.690000

No duplicates found for prd\_id

```

72  SELECT prd_id, COUNT(*)
73  FROM
74  (SELECT
75    pn.prd_id,
76    pn.cat_id,
77    pn.prd_key,
78    pn.prd_nm,
79    pn.prd_cost,
80    pn.prd_line,
81    pn.prd_start_dt,
82    pc.cat,
83    pc.subcat,
84    pc.maintenance
85  FROM
86  silver.crm_prd_info pn
87  LEFT JOIN silver.emp_px_cat_g1v2 pc
88  ON pn.cat_id = pc.id
89  WHERE pn.prd_end_dt IS NULL) --- Filter out historical data
90  GROUP BY prd_id
91  HAVING COUNT(*) > 1
92
  
```

Results

prd_id	(No column name)
478	

We don't have duplicate columns and no need to do any data integration  
Now, grouping similar data together

```

93 ----- grouping similar data together
94 SELECT
95   pn.prd_id,
96   pn.prd_key,
97   pn.prd_nm,
98   pn.cat_id,
99   pc.cat,
100  pc.subcat,
101  pn.prd_cost,
102  pn.prd_line,
103  pn.prd_start_dt,
104  pc.maintenance
105 FROM
106 silver.crm_prd_info pn
107 LEFT JOIN silver.erp_px_cat_g1v2 pc
108 ON pn.cat_id = pc.id
109 WHERE pn.prd_end_dt IS NULL
110
111
112

```

Results Messages

	prd_id	prd_key	prd_nm	cat_id	cat	subcat	prd_cost	prd_line	prd_start_dt	maintenance
1	478	BC-M005	Mountain Bottle Cage	AC_BC	Accessories	Bottles and Cages	4	Mountain	2013-07-01	No
2	479	BC-R205	Road Bottle Cage	AC_BC	Accessories	Bottles and Cages	3	Road	2013-07-01	No
3	477	WB-H098	Water Bottle - 30 oz.	AC_BC	Accessories	Bottles and Cages	2	Other Sales	2013-07-01	No

A B C

Rename columns to friendly, meaningful names

```

93 ----- grouping similar data together and renaming the columns to appropriate name
94
95 SELECT
96   pn.prd_id AS product_id ,
97   pn.prd_key AS product_number,
98   pn.prd_nm AS product_name,
99   pn.cat_id AS category_id,
100  pc.cat AS category,
101  pc.subcat AS subcategory,
102  pn.prd_cost AS cost,
103  pn.prd_line AS product_line,
104  pn.prd_start_dt AS start_date,
105  pc.maintenance AS maintenance
106 FROM
107 silver.crm_prd_info pn
108 LEFT JOIN silver.erp_px_cat_g1v2 pc
109 ON pn.cat_id = pc.id
110 WHERE pn.prd_end_dt IS NULL
111
112

```

Results Messages

	product_id	product_number	product_name	category_id	category	subcategory	cost	product_line	start_date	maintenance
1	478	BC-M005	Mountain Bottle Cage	AC_BC	Accessories	Bottles and Cages	4	Mountain	2013-07-01	No
2	479	BC-R205	Road Bottle Cage	AC_BC	Accessories	Bottles and Cages	3	Road	2013-07-01	No
3	477	WB-H098	Water Bottle - 30 oz.	AC_BC	Accessories	Bottles and Cages	2	Other Sales	2013-07-01	No
4	483	RA-H123	Hitch Rack - 4-Bike	AC_BR	Accessories	Bike Racks	45	Other Sales	2013-07-01	Yes
5	486	ST-1401	All-Purpose Bike Stand	AC_BS	Accessories	Bike Stands	59	Mountain	2013-07-01	No
6	484	CL-9009	Bike Wash - Dissolver	AC_CL	Accessories	Cleaners	3	Other Sales	2013-07-01	Yes



Dimension vs Fact ???

We need to determine the object type—whether it is a Dimension or a Fact table. Since it contains descriptive data, our object type (products) is classified as a Dimension Table.

## Declaring product\_key as the surrogate key

```
94
93     ---- grouping similar data together and renaming the columns to appropriate name
94
95     SELECT
96         ROW_NUMBER() OVER(ORDER BY pn.prd_start_dt,pn.prd_key) AS product_key,
97         pn.prd_id AS product_id ,
98         pn.prd_key AS product_number,
99         pn.prd_nm AS product_name,
100        pn.cat_id AS category_id,
101        pc.cat AS category,
102        pc.subcat AS subcategory,
103        pn.prd_cost AS cost,
104        pn.prd_line AS product_line,
105        pn.prd_start_dt AS start_date,
106        pc.maintenance AS maintenance
107    FROM
108        silver.crm_prd_info pn
109    LEFT JOIN silver.erp_px_cat_g1v2 pc
110        ON pn.cat_id = pc.id
111    WHERE pn.prd_end_dt IS NULL
112
```

Results Messages

	product_key	product_id	product_number	product_name	category_id	category	subcategory	cost	product_line	start_date	maintenance
1	1	210	FR-R92B-58	HL Road Frame - Black- 58	CO_RF	Components	Road Frames	0	Road	2003-07-01	Yes
2	2	211	FR-R92R-58	HL Road Frame - Red- 58	CO_RF	Components	Road Frames	0	Road	2003-07-01	Yes
3	3	348	BK-M82B-38	Mountain-100 Black- 38	BI_MB	Bikes	Mountain Bikes	1898	Mountain	2011-07-01	Yes

## Creating VIEW

```
115     CREATE VIEW gold.dim_products AS
116     SELECT
117         ROW_NUMBER() OVER(ORDER BY pn.prd_start_dt,pn.prd_key) AS product_key,--if we use one column also it will work
118         pn.prd_id AS product_id ,
119         pn.prd_key AS product_number,
120         pn.prd_nm AS product_name,
121         pn.cat_id AS category_id,
122         pc.cat AS category,
123         pc.subcat AS subcategory,
124         pn.prd_cost AS cost,
125         pn.prd_line AS product_line,
126         pn.prd_start_dt AS start_date,
127         pc.maintenance AS maintenance
128     FROM
129        silver.crm_prd_info pn
130    LEFT JOIN silver.erp_px_cat_g1v2 pc
131        ON pn.cat_id = pc.id
132    WHERE pn.prd_end_dt IS NULL
```

Build Gold Layer

Create Fact Sales

Since it contains sales and quantity metrics, it qualifies as a **Fact Table**.

```
SELECT
sd.sls_ord_num,
sd.sls_prd_key,
sd.sls_cust_id,
sd.sls_order_dt,
sd.sls_ship_dt,
sd.sls_due_dt,
sd.sls_sales,
sd.sls_quantity,
sd.sls_price
FROM silver.crm_sales_details sd
```



KEYS			DATES			MEASURES		
sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_quantity	sls_price
SO43697	BK-R93R-62	21768	2010-12-29	2011-01-05	2011-01-10	3578	1	3578
SO43698	BK-M82S-44	283	2010-12-29	2011-01-05	2011-01-10	3400	1	3400

```
SELECT
sd.sls_ord_num,
sd.sls_prd_key,
sd.sls_cust_id,
sd.sls_order_dt,
sd.sls_ship_dt,
sd.sls_due_dt,
sd.sls_sales,
sd.sls_quantity,
sd.sls_price
FROM silver.crm_sales_details sd
```

### Building Fact

Use the dimension's surrogate keys instead of IDs  
to easily connect facts with dimensions

Original IDs								
	sls_ord_num	sls_prd_key	sls_cust_id	sls_order_dt	sls_ship_dt	sls_due_dt	sls_sales	sls_price
1	SO43697	BK-R93R-62	21768	2010-12-29	2011-01-05	2011-01-10	3578	1

ABC	Rename columns to friendly, meaningful names

```

158 -- Renaming the columns
159
160 SELECT
161   sd.sls_ord_num AS order_number,
162   --sd.sls_prd_key,
163   pr.product_key,      -- as we want use this surrogate key instead of sls_prd_key
164   --sd.sls_cust_id,
165   cu.customer_key,
166   sd.sls_order_dt AS order_date,
167   sd.sls_ship_dt AS shipping_date,
168   sd.sls_due_dt AS due_date,
169   sd.sls_sales AS sales_amount,
170   sd.sls_quantity AS quantity,
171   sd.sls_price AS price
172   FROM silver.crm_sales_details sd
173   LEFT JOIN gold.dim_products pr
174   ON      sd.sls_prd_key = pr.product_number
175   LEFT JOIN gold.dim_customers cu
176   ON      sd.sls_cust_id = cu.customer_id

```



sort the columns into logical groups to improve readability

```
sd.sls_due_dt AS due_date,  
sd.sls_sales AS sales_amount,  
sd.sls_quantity AS quantity,  
sd.sls_price price  
FROM silver.crm_sales_details sd  
LEFT JOIN gold.dim_products pr  
ON sd.sls_prd_key = pr.product_number  
LEFT JOIN gold.dim_customers cu  
ON sd.sls_cust_id = cu.customer_id
```

The columns are grouped into three categories:

- DIMENSION KEYS**: order\_number, product\_key, customer\_key
- DATES**: order\_date, shipping\_date, due\_date
- MEASURES**: sales\_amount, quantity, price



We don't need to order the columns here as we've them in correct order.

## Creating VIEW

```
181 CREATE VIEW gold.fact_sales AS  
182     SELECT  
183         sd.sls_ord_num AS order_number,  
184         --sd.sls_prd_key,  
185         pr.product_key,      -- as we want use this surrogate key instead of sls_prd_key  
186         --sd.sls_cust_id,  
187         cu.customer_key,  
188         sd.sls_order_dt AS order_date,  
189         sd.sls_ship_dt AS shipping_date,  
190         sd.sls_due_dt AS due_date,  
191         sd.sls_sales AS sales_amount,  
192         sd.sls_quantity AS quantity,  
193         sd.sls_price AS price  
194     FROM silver.crm_sales_details sd  
195     LEFT JOIN gold.dim_products pr  
196     ON      sd.sls_prd_key = pr.product_number  
197     LEFT JOIN gold.dim_customers cu  
198     ON      sd.sls_cust_id = cu.customer_id  
199  
200  
201     SELECT * FROM [gold].[fact_sales]
```

100 %

Results Messages

	order_number	product_key	customer_key	order_date	shipping_date	due_date	sales_amount	quantity	price
1	SO43697	20	10769	2010-12-29	2011-01-05	2011-01-10	3578	1	3578
2	SO43698	9	17390	2010-12-29	2011-01-05	2011-01-10	3400	1	3400
3	SO43699	9	14864	2010-12-29	2011-01-05	2011-01-10	3400	1	3400

## Fact Check

Check if all dimension tables can successfully join to the fact table

```
204 |
203 | ----Data Integrity check
204 |
205 | SELECT *
206 | FROM [gold].[fact_sales] f
207 | LEFT JOIN [gold].[dim_customers] c
208 | ON f.customer_key = c.customer_key
209 | LEFT JOIN [gold].[dim_products] p
210 | ON f.product_key = p.product_key
211 | WHERE c.customer_key IS NULL
212 |
213 |
```

100 % □ Results □ Messages  
order\_number product\_key customer\_key order\_date shipping\_date due\_date sales\_amount quantity price customer\_key customer\_id customer\_number first\_name last\_name country marital\_status gender birthdate create\_date produ

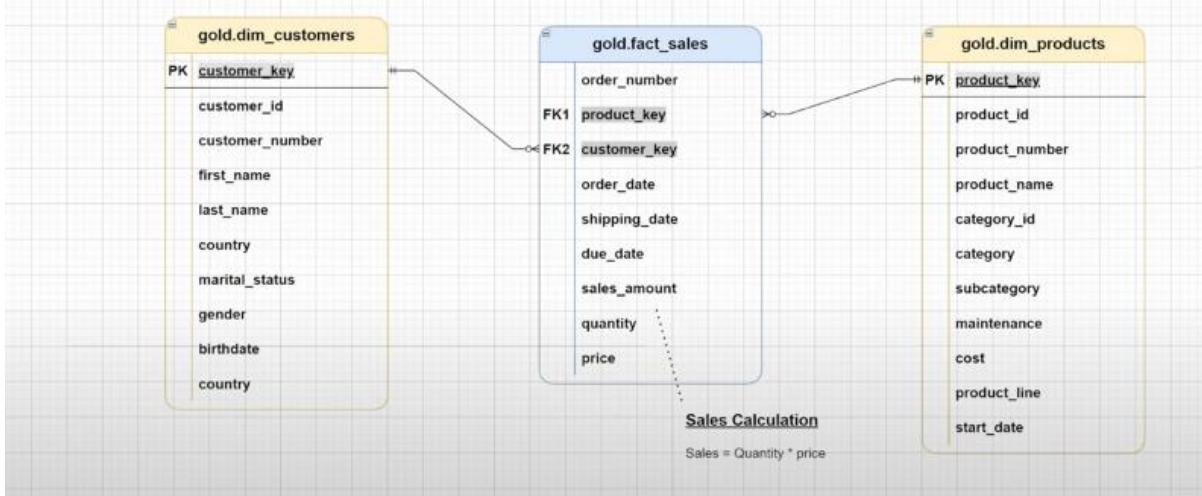
## STAR SCHEMA

Relationship  
in a star schema, the relationship between fact and dimensions is 1-to-many (1:N)

### Many (Optional)

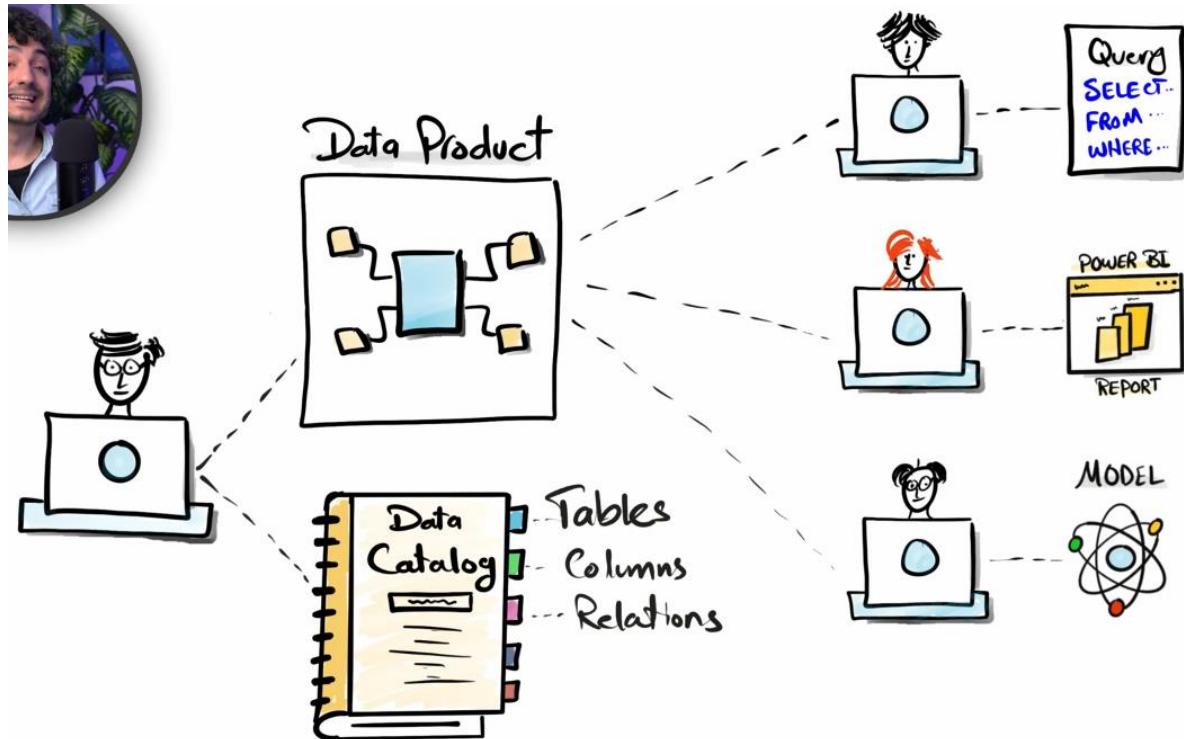
1. Customers who haven't placed any orders yet.
2. Customers who have placed only one order.
3. Customers who have placed multiple orders.

### Sales Data Mart (Star Schema)



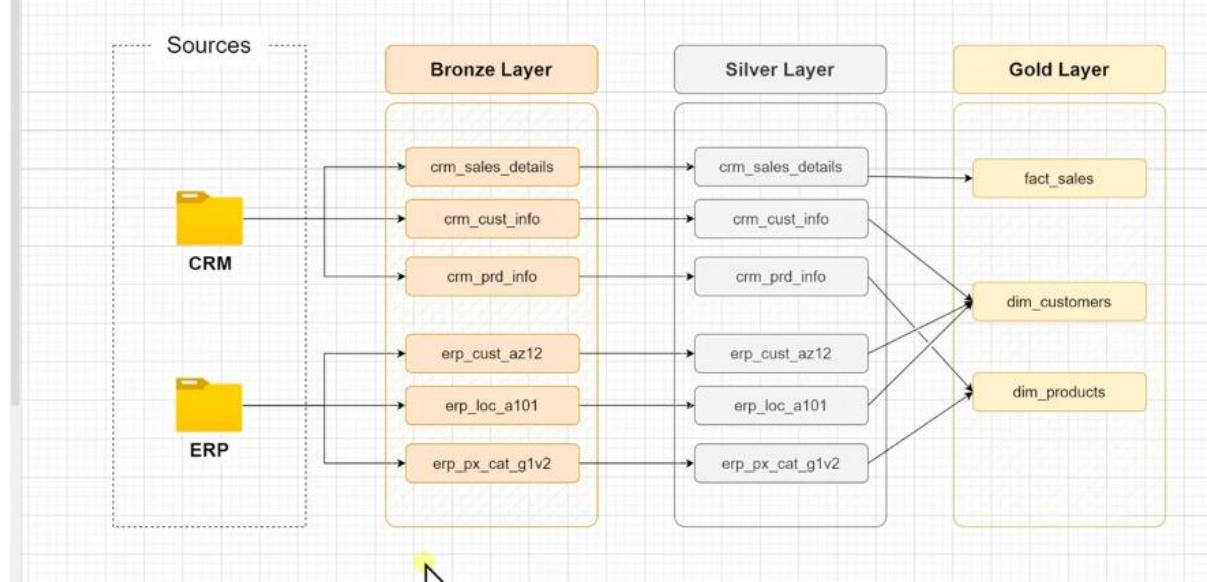
# Build Gold Layer

## Data Catalog

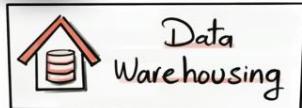


A Data Catalog is a centralized metadata repository that helps organizations efficiently discover, manage, and govern their data assets. It provides a structured inventory of data across various sources, enabling users to search, understand, and utilize data effectively.

## Data Flow Diagram



SQL Projects



Data Warehousing

"Organize, Structure, Prepare,"

- ETL/ELT Processing
- Data Architecture
- Data Integration
- Data Cleansing
- Data Load
- Data Modeling



Exploratory  
Data Analysis  
(EDA)

"Understand Data,"

- Basic Queries
- Data Profiling
- Simple Aggregations
- Subquery



Advanced  
Data Analytics

"Answer Business Questions."

- Complex Queries
- Window Functions
- CTE
- Subqueries
- Reports

## Exploratory Data Analysis (EDA)

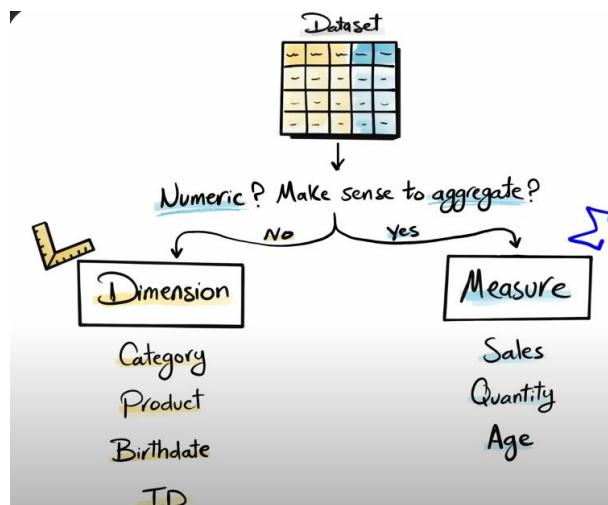
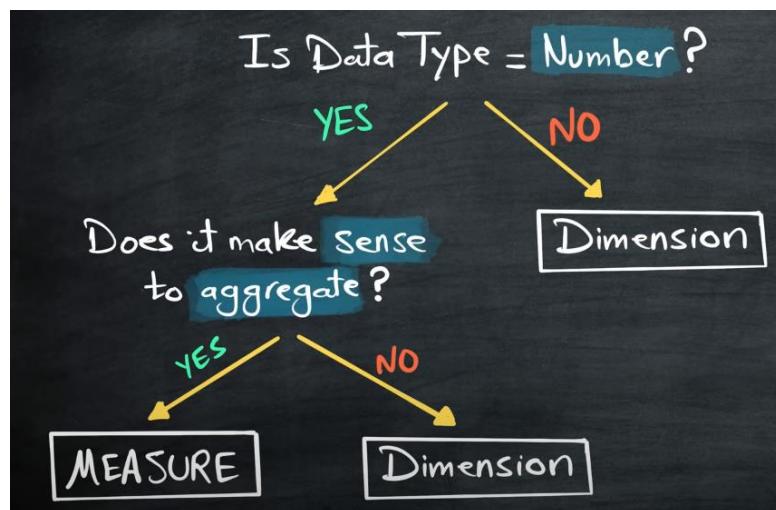
**Exploratory Data Analysis (EDA)** is the process of analysing and summarizing datasets to uncover patterns, detect anomalies, and extract insights using statistical and visualization techniques. It helps in understanding data distributions, relationships, and potential data quality issues before applying machine learning or further analysis.

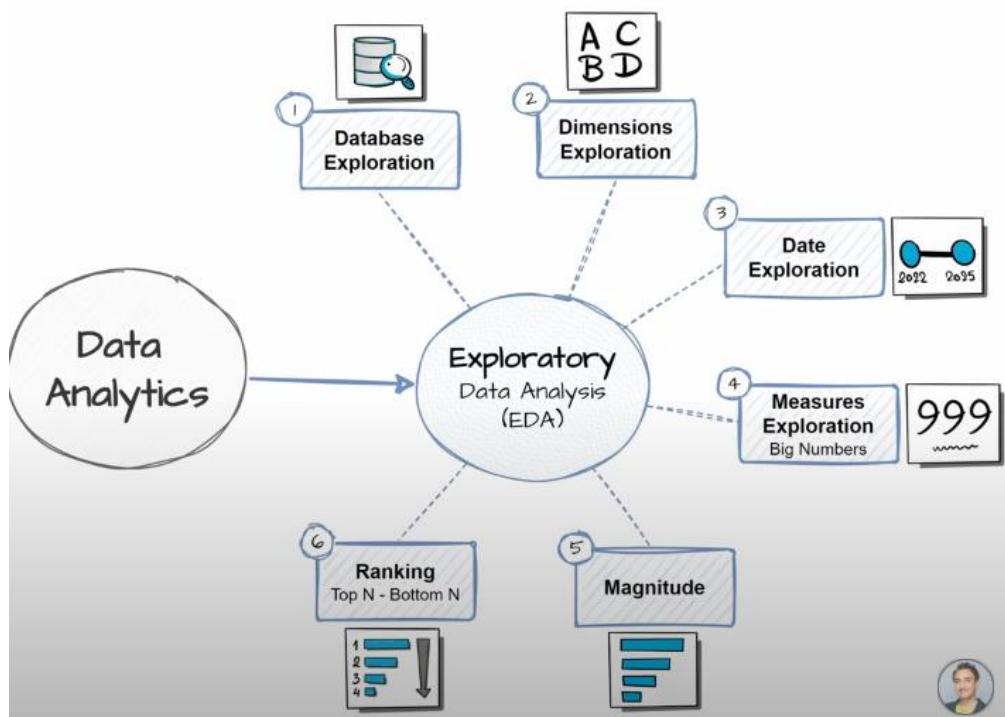
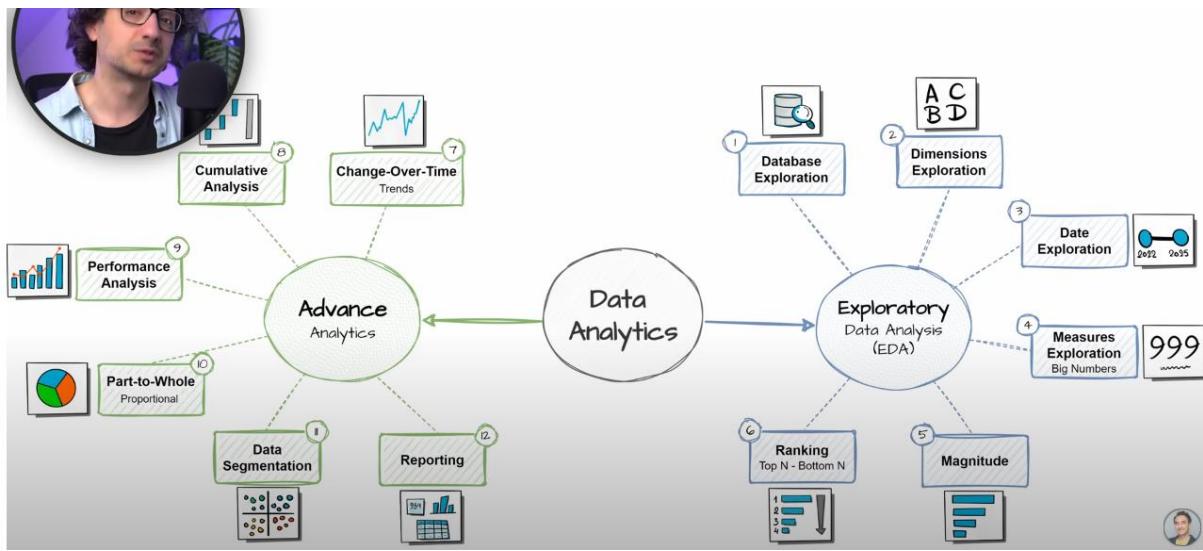
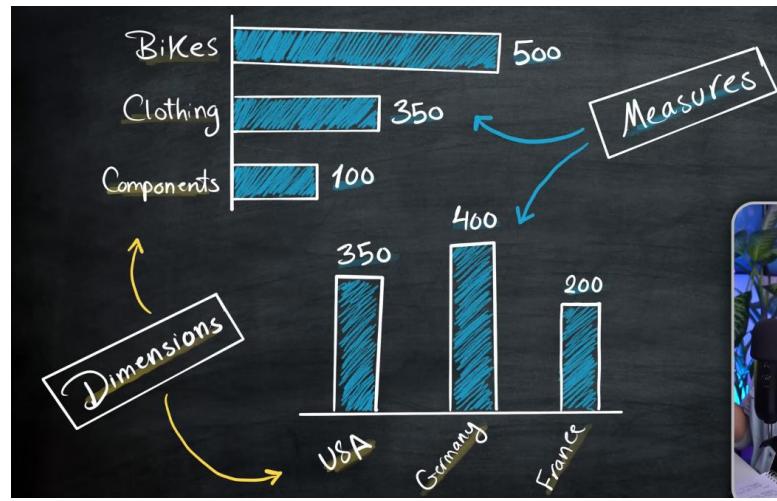
## Dimensions & Measures

Any dataset is typically categorized into **dimensions** and **measures**:

- **Dimensions:** Categorical attributes that provide context to the data (e.g., Product Name, Customer ID, Region, Date). They help in grouping, filtering, and segmenting data.
- **Measures:** Numeric values that can be aggregated and analyzed (e.g., Sales, Quantity, Revenue, Profit). These represent the actual metrics used for calculations.

In a **star schema**, dimensions are stored in **dimension tables**, and measures are stored in **fact tables**.





# Exploratory Data Analysis

## Database Exploration

In **Database Exploration**, we analyze the structure of the database by identifying the number of tables, their relationships, and the types of columns present in each table. This helps in understanding the data model and preparing for further analysis.

The screenshot shows three vertically stacked sections of a SQL query interface. The top section displays a list of tables in the 'DataWarehouse' database. The middle section shows the detailed structure of the 'fact\_sales' table, including its columns and their properties. The bottom section is a repeat of the first two, likely demonstrating a search or a copy-paste operation.

```
1 USE DataWarehouse
2
3 -- Explore All objects in the Database
4
5 SELECT * FROM INFORMATION_SCHEMA.TABLES
```

TABLE_CATALOG	TABLE_SCHEMA	TABLE_NAME	TABLE_TYPE
1 DataWarehouse	silver	crm_sales_details	BASE TABLE
2 DataWarehouse	gold	dim_customers	VIEW
3 DataWarehouse	gold	dim_products	VIEW
4 DataWarehouse	gold	fact_sales	VIEW
5 DataWarehouse	gold	dim_customers_tb	BASE TABLE
6 DataWarehouse	gold	dim_products_tb	BASE TABLE
7 DataWarehouse	gold	fact_sales_tb	BASE TABLE
8 DataWarehouse	bronze	crm_cust_info	BASE TABLE
9 DataWarehouse	bronze	erp_cust_az12	BASE TABLE
10 DataWarehouse	bronze	erp_loc_a101	BASE TABLE
11 DataWarehouse	bronze	crm_prd_info	BASE TABLE
12 DataWarehouse	bronze	crm_sales_details	BASE TABLE
13 DataWarehouse	bronze	erp_px_cat_g1v2	BASE TABLE
14 DataWarehouse	silver	crm_cust_info	BASE TABLE
15 DataWarehouse	silver	erp_cust_az12	BASE TABLE
16 DataWarehouse	silver	erp_loc_a101	BASE TABLE
17 DataWarehouse	silver	erp_px_cat_g1v2	BASE TABLE
18 DataWarehouse	silver	crm_prd_info	BASE TABLE

```
7 -- Explore All columns in the table
8
9 SELECT * FROM INFORMATION_SCHEMA.COLUMNS
10 WHERE TABLE_NAME = 'fact_sales'
```

TABLE_CATALOG	TABLE_SCHEMA	TABLE_NAME	COLUMN_NAME	ORDINAL_POSITION	COLUMN_TYPE
1 DataWarehouse	gold	fact_sales	order_number	1	NULL
2 DataWarehouse	gold	fact_sales	product_key	2	NULL
3 DataWarehouse	gold	fact_sales	customer_key	3	NULL
4 DataWarehouse	gold	fact_sales	order_date	4	NULL
5 DataWarehouse	gold	fact_sales	shipping_date	5	NULL
6 DataWarehouse	gold	fact_sales	due_date	6	NULL
7 DataWarehouse	gold	fact_sales	sales_amount	7	NULL
8 DataWarehouse	gold	fact_sales	quantity	8	NULL
9 DataWarehouse	gold	fact_sales	price	9	NULL

```
1 USE DataWarehouse
2
3 -- Explore All objects in the Database
4
5 SELECT * FROM INFORMATION_SCHEMA.TABLES
6
7 -- Explore All columns in the table
8
9 SELECT * FROM INFORMATION_SCHEMA.COLUMNS
10 WHERE TABLE_NAME = 'fact_sales'
```

# Exploratory Data Analysis

## Dimensions Exploration



A C  
B D

Dimensions Exploration

Identifying the unique values (or categories) in each dimension.

Recognizing how data might be grouped or segmented,  
which is useful for later analysis.



A C  
B D

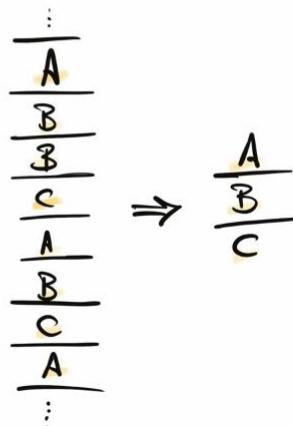
Dimensions Exploration

DISTINCT [Dimension]

DISTINCT Country

DISTINCT Category

DISTINCT Product



```
20 -- Explore All Countries our customers come from.  
21  
22 SELECT DISTINCT country FROM gold.dim_customers  
23 ORDER BY 1
```

100 %

Results Messages

	country
1	Australia
2	Canada
3	France
4	Germany
5	n/a
6	United Kingdom
7	United States

```

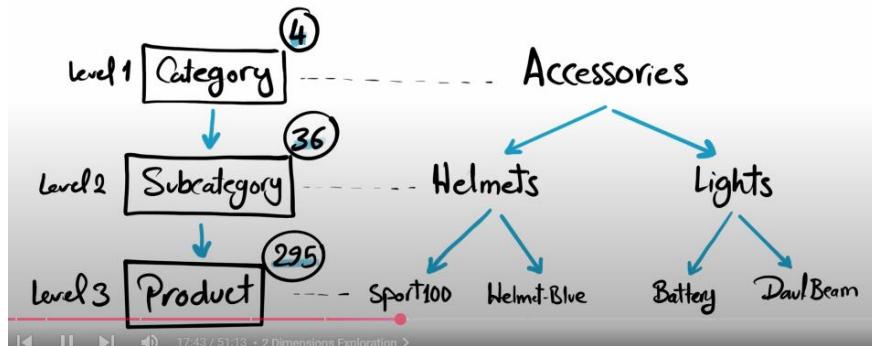
25 -- Explore All Categories "The Major Divisions"
26
27 SELECT DISTINCT category, subcategory, product_name FROM gold.dim_products
28 ORDER BY 1,2,3

```

100 %

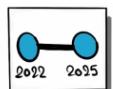
Results Messages

	category	subcategory	product_name
28	Accessories	Tires and Tubes	LL Mountain Tire
29	Accessories	Tires and Tubes	LL Road Tire
30	Accessories	Tires and Tubes	ML Mountain Tire
31	Accessories	Tires and Tubes	ML Road Tire
32	Accessories	Tires and Tubes	Mountain Tire Tube
33	Accessories	Tires and Tubes	Patch Kit/8 Patches
34	Accessories	Tires and Tubes	Road Tire Tube
35	Accessories	Tires and Tubes	Touring Tire
36	Accessories	Tires and Tubes	Touring Tire Tube
37	Bikes	Mountain Bikes	Mountain-100 Black- 38
38	Bikes	Mountain Bikes	Mountain-100 Black- 42
39	Bikes	Mountain Bikes	Mountain-100 Black- 44
40	Bikes	Mountain Bikes	Mountain-100 Black- 48
41	Bikes	Mountain Bikes	Mountain-100 Silver- 38
42	Bikes	Mountain Bikes	Mountain-100 Silver- 42



## Exploratory Data Analysis

# Date Exploration



Date Exploration

Identify the earliest and latest dates (boundaries).

Understand the scope of data and the timespan.



## Date Exploration

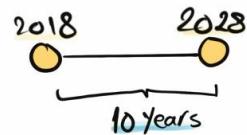
### MIN/MAX [Date Dimension]

MIN Order\_date

MAX Create\_date

MIN Birthdate

2019
2020
2018
2018
2022
2023
2023
2023
2023
2022
2022



DATEDIFF

```
55  
34 -- Find the date of the first and last order  
35  
36 SELECT  
37 MIN(order_date) AS first_order_date ,  
38 MAX(order_date) AS last_order_date  
39 FROM gold.dim_sales_tb
```

100 %

	first_order_date	last_order_date
1	2010-12-29	2014-01-28

```
41 -- How many years of sales are available  
42  
43 SELECT  
44 MIN(order_date) AS first_order_date ,  
45 MAX(order_date) AS last_order_date ,  
46 DATEDIFF(MONTH,MIN(order_date), MAX(order_date)) AS order_range_months  
47 FROM gold.dim_sales_tb
```

100 %

	first_order_date	last_order_date	order_range_months
1	2010-12-29	2014-01-28	37

```
49 -- Find the youngest and oldest customer  
50  
51 SELECT  
52 MIN(birthdate) AS oldest_birthdate ,  
53 DATEDIFF(YEAR, MIN(birthdate), GETDATE() ) AS oldest_age ,  
54 MAX(birthdate) AS youngest_birthdate ,  
55 DATEDIFF(YEAR, MAX(birthdate), GETDATE() ) AS youngest_age  
56 FROM  
57 gold.dim_customers_tb
```

100 %

	oldest_birthdate	oldest_age	youngest_birthdate	youngest_age
1	1916-02-10	109	1986-06-25	39

# Exploratory Data Analysis

## Measures Exploration



999

Measures Exploration

Calculate the key metric of the business (Big Numbers)

- Highest Level of Aggregation | Lowest Level of Details -



999

Measures Exploration

$\sum$  [Measure]

SUM (Sales)

AVG (Price)

SUM (Quantity)

10
20
50
30
10
80
30
10

}

BIG Number

240

Key Metric

```
-- Find the total sales
SELECT SUM(sales_amount) AS total_sales FROM gold.fact_sales_tb

-- Find how many items are sold
SELECT SUM(quantity) AS total_quantity FROM gold.fact_sales_tb

-- Find the average selling price
SELECT AVG(price) AS avg_price FROM gold.fact_sales_tb

-- Find the total no of orders
SELECT COUNT(DISTINCT order_number) AS total_orders FROM gold.fact_sales_tb

-- Find the total no of products
SELECT COUNT(product_key) AS total_products FROM gold.dim_products_tb

-- Find the total no of customers
SELECT COUNT(customer_key) AS total_customers FROM gold.dim_customers_tb

-- Find the total no of customers that has placed an order

SELECT COUNT(DISTINCT customer_key) FROM gold.fact_sales_tb
```

```

86 -- Generate a report that shows all key metrics of the business
87
88 SELECT 'total_sales' AS measure_name, SUM(sales_amount) AS measure_value FROM gold.fact_sales_tb
89 UNION ALL
90 SELECT 'total_quantity' AS measure_name, SUM(quantity) AS measure_value FROM gold.fact_sales_tb
91 UNION ALL
92 SELECT 'avg_price' AS measure_name, AVG(price) AS measure_value FROM gold.fact_sales_tb
93 UNION ALL
94 SELECT 'total_no_of_orders' AS measure_name, COUNT(DISTINCT order_number) AS measure_value FROM gold.fact_sales_tb
95 UNION ALL
96 SELECT 'total_no_of_products' AS measure_name, COUNT(product_key) AS measure_value FROM gold.dim_products_tb
97 UNION ALL
98 SELECT 'total_no_of_customers' AS measure_name, COUNT(customer_key) AS measure_value FROM gold.dim_customers_tb
99
100

```

100 %

Results Messages

measure_name	measure_value
total_sales	29356250
total_quantity	60423
avg_price	486
total_no_of_orders	27659
total_no_of_products	295
total_no_of_customers	18484

# Exploratory Data Analysis

## Magnitude Analysis



Magnitude

Compare the measure values by categories.

It helps us understand the importance of different categories.



Magnitude

$\Sigma$  [Measure] By [Dimension]

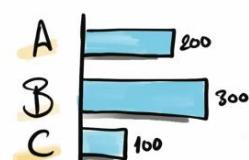
Total Sales By Country

Total Quantity By Category

Average Price By Product

Total Orders By Customer

	600
A	200
B	300
C	100



```
-- Find total customers by countries  
-- Find total customers by gender  
-- Find total products by category  
-- What is the average costs in each category?  
-- What is the total revenue generated for each category?  
-- Find total revenue is generated by each customer  
-- What is the distribution of sold items across countries?
```

## Exploratory Data Analysis

### Ranking Analysis



Order the values of dimensions by measure.

Top N performers | Bottom N Performers



Rank [Dimension] By  $\sum$  [Measure]

Rank Countries By Total Sales

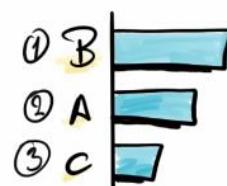
Top 5 Products By Quantity

Bottom 3 Customers By Total Orders

Top RANK()  
DENSE\_RANK()  
ROW\_NUMBER()

①	B	300
②	A	200

↓



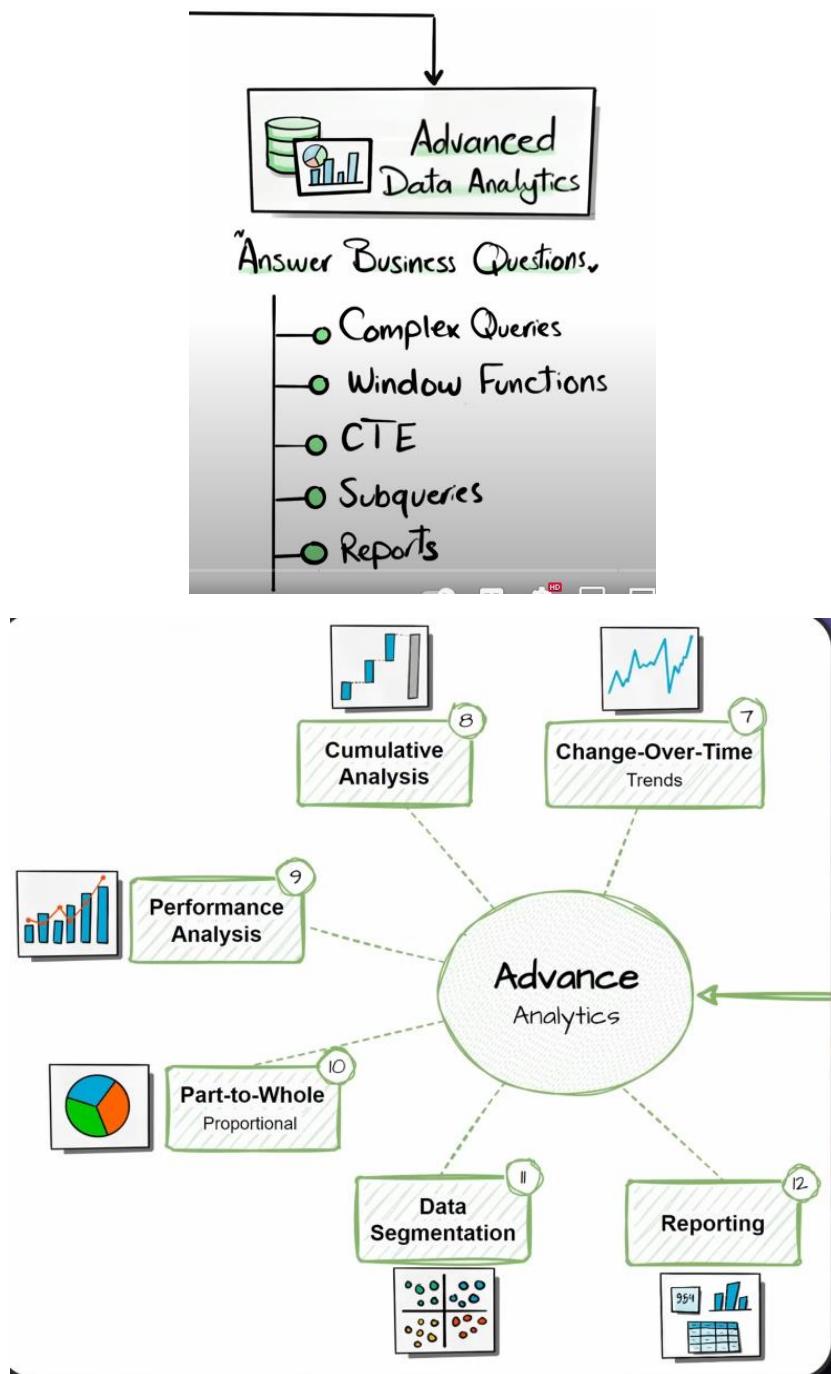
---

```
-- Which 5 products generate the highest revenue?
```

---

```
-- What are the 5 worst-performing products in terms of sales?
```

# Advance Data Analytics (EDA)



# Advanced Analytics Project

## Changes Over Time Analysis



Change - Over - Time Trends

Analyze how a measure evolves over time.

Helps track trends and identify seasonality in your data.



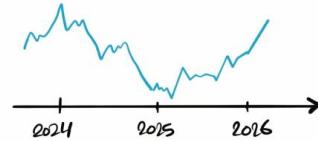
Change - Over - Time Trends

$\sum$  [Measure] By [Date Dimension]

Total Sales By Year

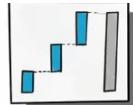
Average Cost By Month

2024	300
2025	100
2026	200



```
17  
18 -- Analyse sales performance over time  
19 -- -- Quick Date Functions  
20  
21  
22 SELECT  
23     YEAR(order_date) AS order_year,  
24     MONTH(order_date) AS order_month,  
25     SUM(sales_amount) AS total_sales,  
26     SUM(quantity) AS total_quantity,  
27     COUNT(DISTINCT customer_key) AS total_customers,  
28     COUNT(DISTINCT order_number) AS total_orders  
29     FROM  
30     gold.fact_sales_tb  
31     WHERE order_date IS NOT NULL  
32     GROUP BY YEAR(order_date),  
33                  MONTH(order_date)  
34     ORDER BY YEAR(order_date),  
35                  MONTH(order_date)
```

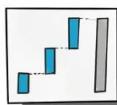
# Cumulative Analysis



Cumulative Analysis

Aggregate the data progressively over time.

Helps to understand whether our business is growing or declining.



Cumulative Analysis

$\sum$  [Cumulative Measure] By [Date Dimension]

Running Total Sales By Year

Moving Average of Sales By Month

2024	300	300
2025	100	400
2026	200	600

Cumulative



## WINDOW FUNCTIONS

```
-- Calculate the total sales per year
-- and the running total of sales over time
SELECT
    order_year,
    total_sales,
    SUM(total_sales)    OVER (ORDER BY order_year) AS running_total_sales,
    AVG(average_sales)  OVER (ORDER BY order_year) AS moving_average_sales
FROM
    (
        SELECT
            DATETRUNC(YEAR,order_date) AS order_year,
            SUM(sales_amount) AS total_sales,
            AVG(sales_amount) AS average_sales
        FROM gold.fact_sales_tb
        WHERE order_date IS NOT NULL
        GROUP BY DATETRUNC(YEAR,order_date) ) t
```

# Performance Analysis



## Performance Analysis

Comparing the current value to a target value.

Helps measure success and compare performance.



## Performance Analysis

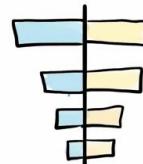
$\text{Current [Measure]} - \text{Target [Measure]}$

Current Sales - Average Sales

Current Year Sales - Previous Year Sales

Current Sales - lowest Sales

	Current	-	Target (Avg)	Performance
A	200	-	200	0
B	300	-	200	100
C	100	-	200	-100



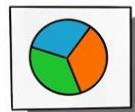
## WINDOW FUNCTIONS

```

/*
Analyze the yearly performance of products by comparing their sales
to both the average sales performance of the product and the previous year's sales */
WITH yearly_product_sales AS (
  SELECT
    YEAR(f.order_date)          AS order_year,
    p.product_name               AS product_name,
    SUM(f.sales_amount)          AS current_year_sales,
    COUNT(DISTINCT order_number) AS current_year_orders
  FROM
    gold.fact_sales_tb f
  LEFT JOIN gold.dim_products_tb p
  ON
    f.product_key = p.product_key
  WHERE f.order_date IS NOT NULL
  GROUP BY YEAR(f.order_date),
           p.product_name
)
SELECT
  order_year,
  product_name,
  current_year_sales,
  current_year_orders,
  AVG(current_year_sales) OVER (PARTITION BY product_name) AS avg_sale_price,
  current_year_sales - AVG(current_year_sales) OVER(PARTITION BY product_name) AS diff_avg_price,
  CASE
    WHEN current_year_sales - AVG(current_year_sales) OVER(PARTITION BY product_name) > 0 THEN 'Above Avg'
    WHEN current_year_sales - AVG(current_year_sales) OVER(PARTITION BY product_name) < 0 THEN 'Below Avg'
    ELSE 'Avg'
  END AS avg_change,
  -- Year-over-Year Analysis
  LAG(current_year_sales) OVER (PARTITION BY product_name ORDER BY order_year) AS prev_yr_sales,
  current_year_sales - LAG(current_year_sales) OVER (PARTITION BY product_name ORDER BY order_year) AS diff_prev_yr_price,
  CASE
    WHEN current_year_sales - LAG(current_year_sales) OVER (PARTITION BY product_name ORDER BY order_year) > 0 THEN 'Increase'
    WHEN current_year_sales - LAG(current_year_sales) OVER (PARTITION BY product_name ORDER BY order_year) < 0 THEN 'Decrease'
    ELSE 'No Change'
  END AS prev_yr_change
FROM yearly_product_sales
ORDER BY product_name,
         order_year ;

```

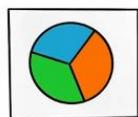
## Part-To-Whole Analysis



Part-to-Whole

Proportional Analysis

Analyze how an individual part is performing compared to the overall, allowing us to understand which category has the greatest impact on the business.



Part-to-Whole

Proportional Analysis

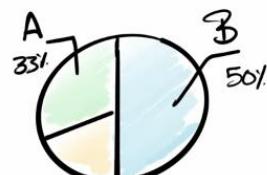
$([\text{Measure}] / \text{Total [Measure]}) * 100$  By [Dimension]

$(\text{Sales} / \text{Total Sales}) * 100$  By Category

$(\text{Quantity} / \text{Total Quantity}) * 100$  By Country

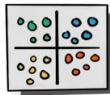


A	200	33%
B	300	50%
C	100	17%



```
-- Which categories contribute the most to overall sales?
WITH category_sales AS (
    SELECT
        p.category AS category,
        SUM(f.sales_amount) AS sales_amount
    FROM
        gold.fact_sales_tb f
    LEFT JOIN gold.dim_products_tb p
    ON
        f.product_key = p.product_key
    WHERE f.order_date IS NOT NULL
    GROUP BY p.category )
SELECT
    category,
    sales_amount,
    SUM(sales_amount) OVER() AS overall_sales,
    ROUND((CAST(sales_amount AS FLOAT)/ SUM(sales_amount) OVER())*100, 2 ) AS percentage_of_total
FROM category_sales
ORDER BY percentage_of_total DESC;
```

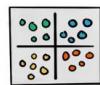
# Data Segmentation



Data Segmentation

Group the data based on a specific range.

Helps understand the correlation between two measures.



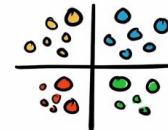
Data Segmentation

[Measure] By [Measure]

Total Products By Sales Range

Total Customers By Age

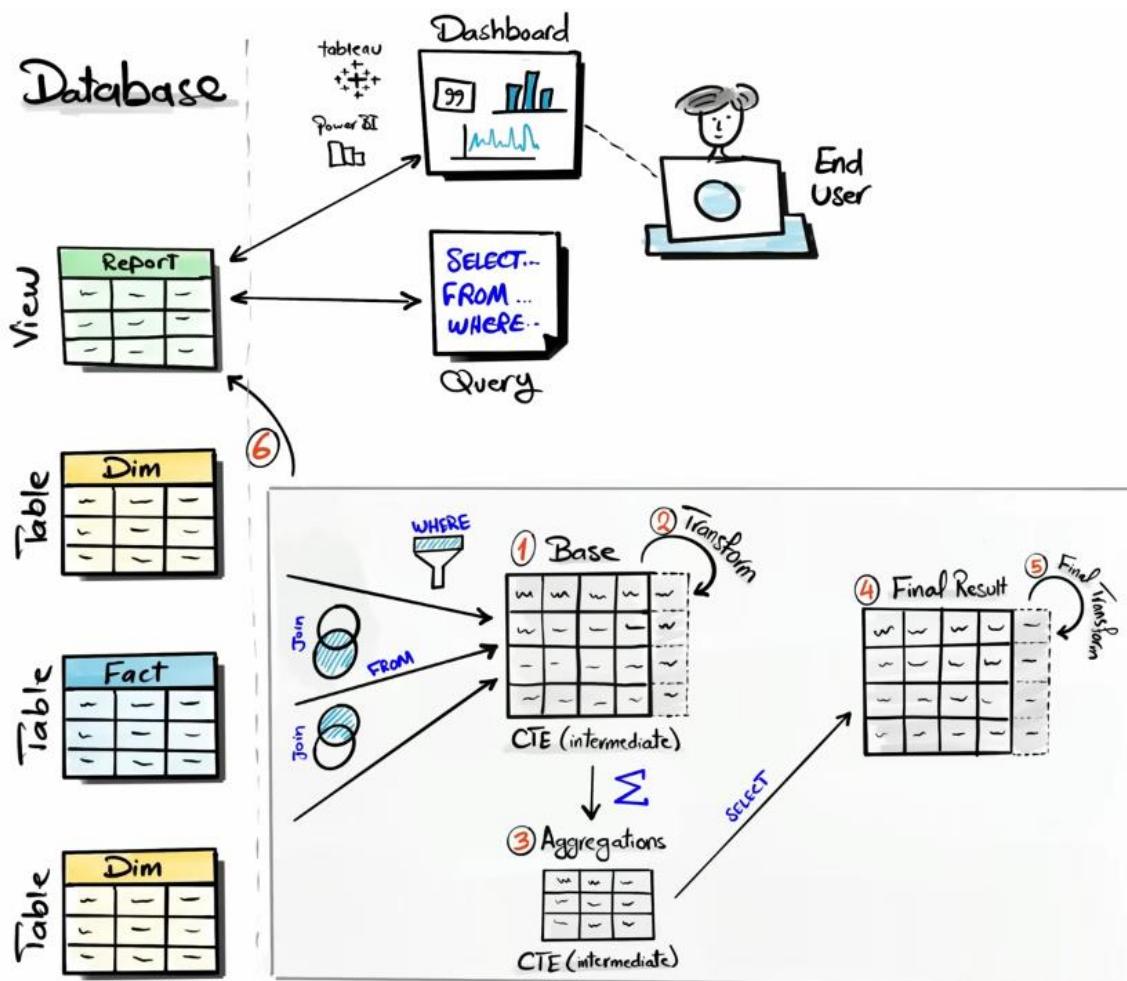
		Categorize	
3	50	Low	7
4	100	Medium	6
5	150		
1	200	Large	15
10	250		
5	300		



CASE WHEN STATEMENT

```
/*Group customers into three segments based on their spending behavior:  
 - VIP: Customers with at least 12 months of history and spending more than €5,000.  
 - Regular: Customers with at least 12 months of history but spending €5,000 or less.  
 - New: Customers with a lifespan less than 12 months.  
And find the total number of customers by each group*/  
  
WITH customer_spending AS (  
    SELECT  
        customer_key,  
        SUM(sales_amount) AS total_spending,  
        MIN(order_date) AS min_order_date,  
        MAX(order_date) AS max_order_date,  
        DATEDIFF(MONTH, MIN(order_date), MAX(order_date)) AS lifespan  
    FROM  
    gold.fact_sales_tb  
    GROUP BY customer_key )  
SELECT  
customer_segment,  
COUNT(DISTINCT customer_key) AS total_customers  
FROM  
(SELECT  
    customer_key,  
    total_spending,  
    min_order_date,  
    max_order_date,  
    lifespan,  
    CASE  
        WHEN lifespan >= 12 and total_spending > 5000 THEN 'VIP'  
        WHEN lifespan >= 12 and total_spending <= 5000 THEN 'Regular'  
        ELSE 'New'  
    END AS customer_segment  
    FROM  
    customer_spending ) t  
    GROUP BY customer_segment  
    ORDER BY total_customers DESC ;
```

# Build Customer Report



```
/*
=====
Customer Report
=====
Purpose:
- This report consolidates key customer metrics and behaviors

Highlights:
1. Gathers essential fields such as names, ages, and transaction details.
2. Segments customers into categories (VIP, Regular, New) and age groups.
3. Aggregates customer-level metrics:
- total orders
- total sales
- total quantity purchased
- total products
- lifespan (in months)
4. Calculates valuable KPIs:
- recency (months since last order)
- average order value
- average monthly spend
=====
*/
```

```

IF OBJECT_ID('gold.report_customers', 'V') IS NOT NULL
    DROP VIEW gold.report_customers;
GO

CREATE VIEW gold.report_customers AS
WITH base_query AS
/*
1) Base Query: Retrieves core columns from tables
*/
(
SELECT
    f.order_number,
    f.product_key,
    f.order_date,
    f.sales_amount,
    f.quantity,
    c.customer_key,
    c.customer_number,
    CONCAT(c.first_name, ' ', c.last_name) AS customer_name,
    c.birthdate,
    DATEDIFF(YEAR, c.birthdate, GETDATE()) AS age
FROM
    gold.fact_sales_tb f
LEFT JOIN gold.dim_customers_tb c
ON      f.customer_key = c.customer_key
),
/*
2) Customer Aggregations: Summarizes key metrics at the customer level
*/
/*
2) Customer Aggregations: Summarizes key metrics at the customer level
*/
customer_aggregation AS
(
SELECT
    customer_key,
    customer_number,
    customer_name,
    age,
    COUNT(DISTINCT order_number) AS total_orders,
    SUM(sales_amount) AS total_sales,
    SUM(quantity) AS total_quantity,
    COUNT(DISTINCT product_key) AS total_products,
    DATEDIFF(MONTH, MIN(order_date), MAX(order_date)) AS lifespan,
    MAX(order_date) AS last_order_date
FROM
    base_query
GROUP BY
    customer_key,
    customer_number,
    customer_name,
    age
)
/*
3) Calculations: Calculating valuable KPI's
*/

```

```
/*
3) Calculations: Calculating valuable KPI's
*/
SELECT
    customer_key,
    customer_number,
    customer_name,
    age,
    CASE
        WHEN age < 20 THEN 'Under 20'
        WHEN age between 20 and 29 THEN '20-29'
        WHEN age between 30 and 39 THEN '30-39'
        WHEN age between 40 and 49 THEN '40-49'
        ELSE '50 and above'
    END AS age_group,
    CASE
        WHEN lifespan >= 12 and total_sales > 5000 THEN 'VIP'
        WHEN lifespan >= 12 and total_sales <= 5000 THEN 'Regular'
        ELSE 'NEW'
    END AS customer_segment,
    last_order_date,
    total_orders,
    total_sales,
    total_quantity,
    total_products,
    lifespan,
    -- recency
    DATEDIFF(month, last_order_date, GETDATE()) AS recency,
    -- Compute average order value (AVO)
    CASE WHEN total_sales = 0 THEN 0
        ELSE total_sales / total_orders
    END AS avg_order_value,
    -- Compute average monthly spend
    CASE WHEN lifespan = 0 THEN total_sales
        ELSE total_sales / lifespan
    END AS avg_monthly_spend
FROM
    customer_aggregation;
```

# Build Product Report

```
/*
=====
Product Report
=====
Purpose:
    - This report consolidates key product metrics and behaviors.

Highlights:
    1. Gathers essential fields such as product name, category, subcategory, and cost.
    2. Segments products by revenue to identify High-Performers, Mid-Range, or Low-Performers.
    3. Aggregates product-level metrics:
        - total orders
        - total sales
        - total quantity sold
        - total customers (unique)
        - lifespan (in months)
    4. Calculates valuable KPIs:
        - recency (months since last sale)
        - average order revenue (AOR)
        - average monthly revenue
=====
*/
USE DataWarehouse

IF OBJECT_ID('gold.report_products', 'V') IS NOT NULL
    DROP VIEW gold.report_products;
GO

CREATE VIEW gold.report_products AS
WITH base_query AS (
/*
1) Base Query: Retrieves core columns from fact_sales and dim_products
*/
    SELECT
        p.product_key,
        p.product_number,
        p.category,
        p.product_name,
        p.subcategory,
        p.cost,
        f.order_number,
        f.customer_key,
        f.sales_amount,
        f.quantity,
        f.order_date
    FROM
        gold.fact_sales_tb f
    LEFT JOIN gold.dim_products_tb p
        ON      f.product_key = p.product_key
    WHERE f.order_date IS NOT NULL),
/*
2) Product Aggregations: Summarizes key metrics at the product level
*/

```

```

/*
2) Product Aggregations: Summarizes key metrics at the product level
*/
product_aggregations AS(
SELECT
    product_key,
    product_number,
    category,
    product_name,
    subcategory,
    cost,
    COUNT(DISTINCT order_number) AS total_orders,
    SUM(sales_amount) AS total_sales,
    SUM(quantity) AS total_quantity,
    COUNT(customer_key) AS total_customers,
    MAX(order_date) AS last_order_date,
    DATEDIFF(MONTH, MIN(order_date), MAX(order_date)) AS lifespan,
    ROUND(AVG(CAST(sales_amount AS FLOAT) / NULLIF(quantity, 0)),1) AS avg_selling_price
FROM base_query
GROUP BY
    product_key,
    product_number,
    category,
    product_name,
    subcategory,
    cost)
/*
3) Calculations: Calculating valuable KPI's
*/
/*
3) Calculations: Calculating valuable KPI's
*/
SELECT
    product_key,
    product_number,
    category,
    product_name,
    subcategory,
    cost,
    total_orders,
    total_sales,
    avg_selling_price,
    total_quantity ,
    total_customers,
    last_order_date,
    lifespan,
    -- recency
    DATEDIFF(MONTH, last_order_date, GETDATE() ) AS recency_in_months,
    -- product segment
    CASE
        WHEN total_sales > 50000 THEN 'High-Performer'
        WHEN total_sales >= 10000 THEN 'Mid-Range'
        ELSE 'Low-Performer'
    END AS product_segment,
    -- Average Order Revenue (AOR)
    CASE
        WHEN total_orders = 0 THEN 0
        ELSE total_sales / total_orders
    END AS avg_order_revenue,
    -- Average Monthly Revenue
    CASE
        WHEN lifespan = 0 THEN total_sales
        ELSE total_sales / lifespan
    END AS avg_monthly_revenue
FROM
    product_aggregations;
--Execution
SELECT * FROM gold.report_products

```