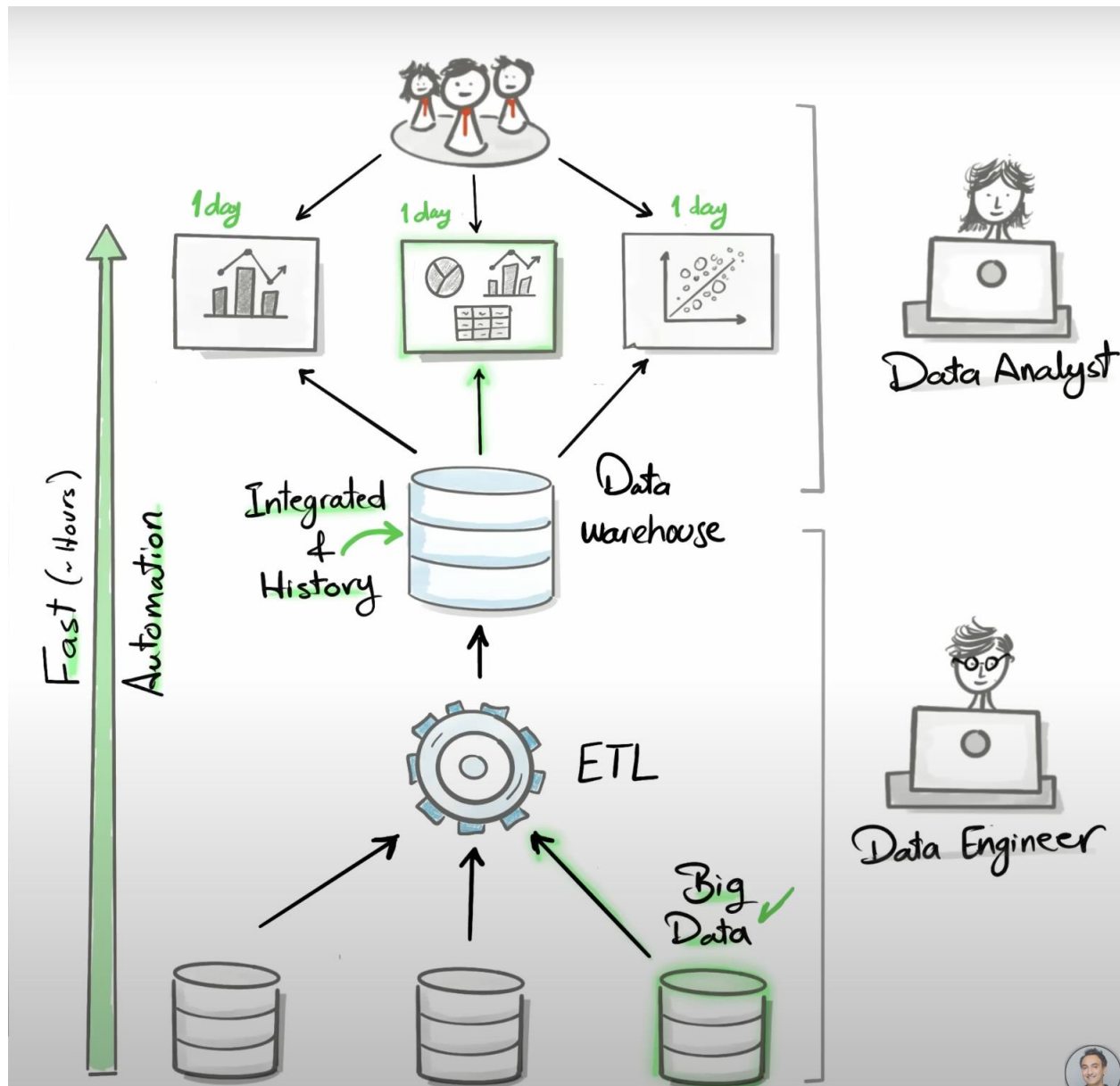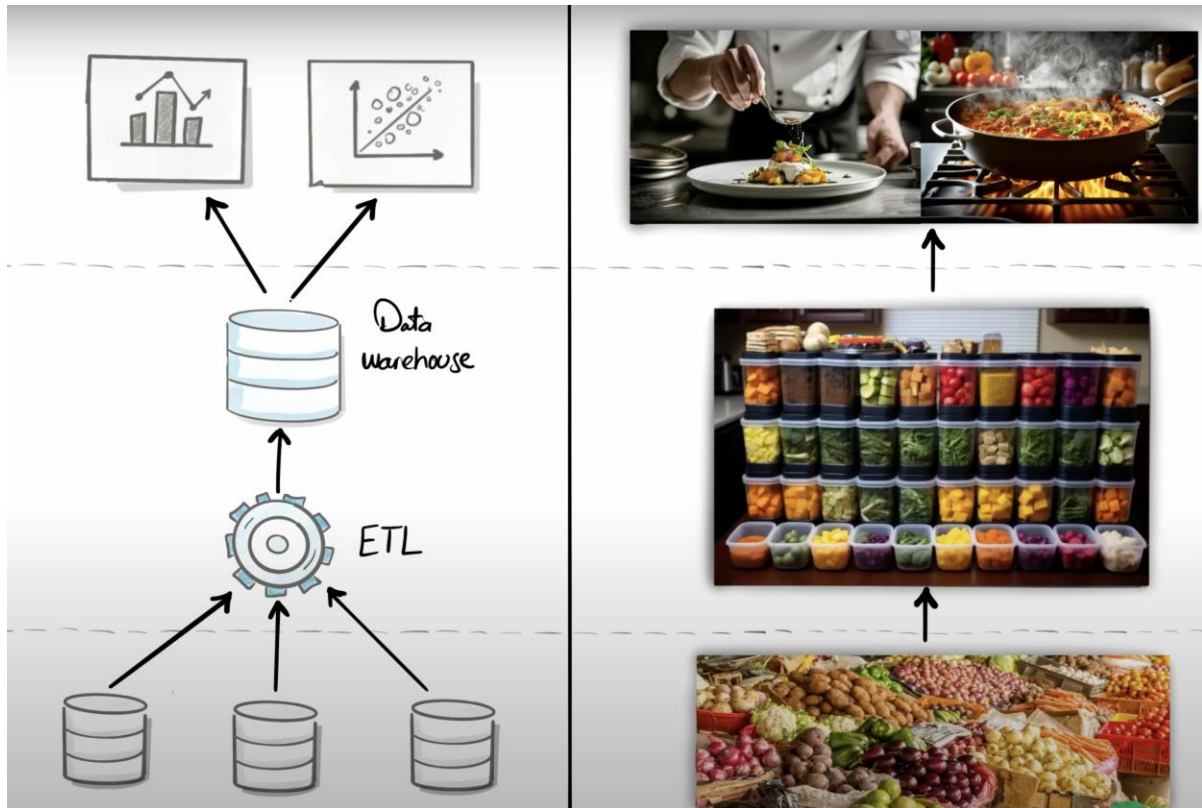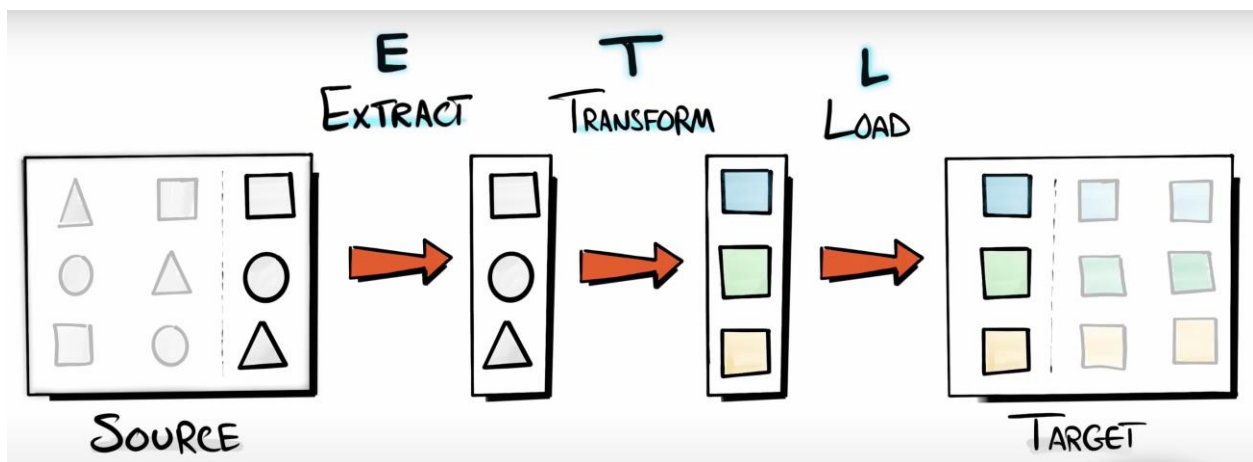# Data Warehouse project



- Here we can pull the data into multiple sources and using ETL process we can store them in a Data warehouse.
- If we have a data warehouse in cloud platform, then we can easily handle even big data easily.
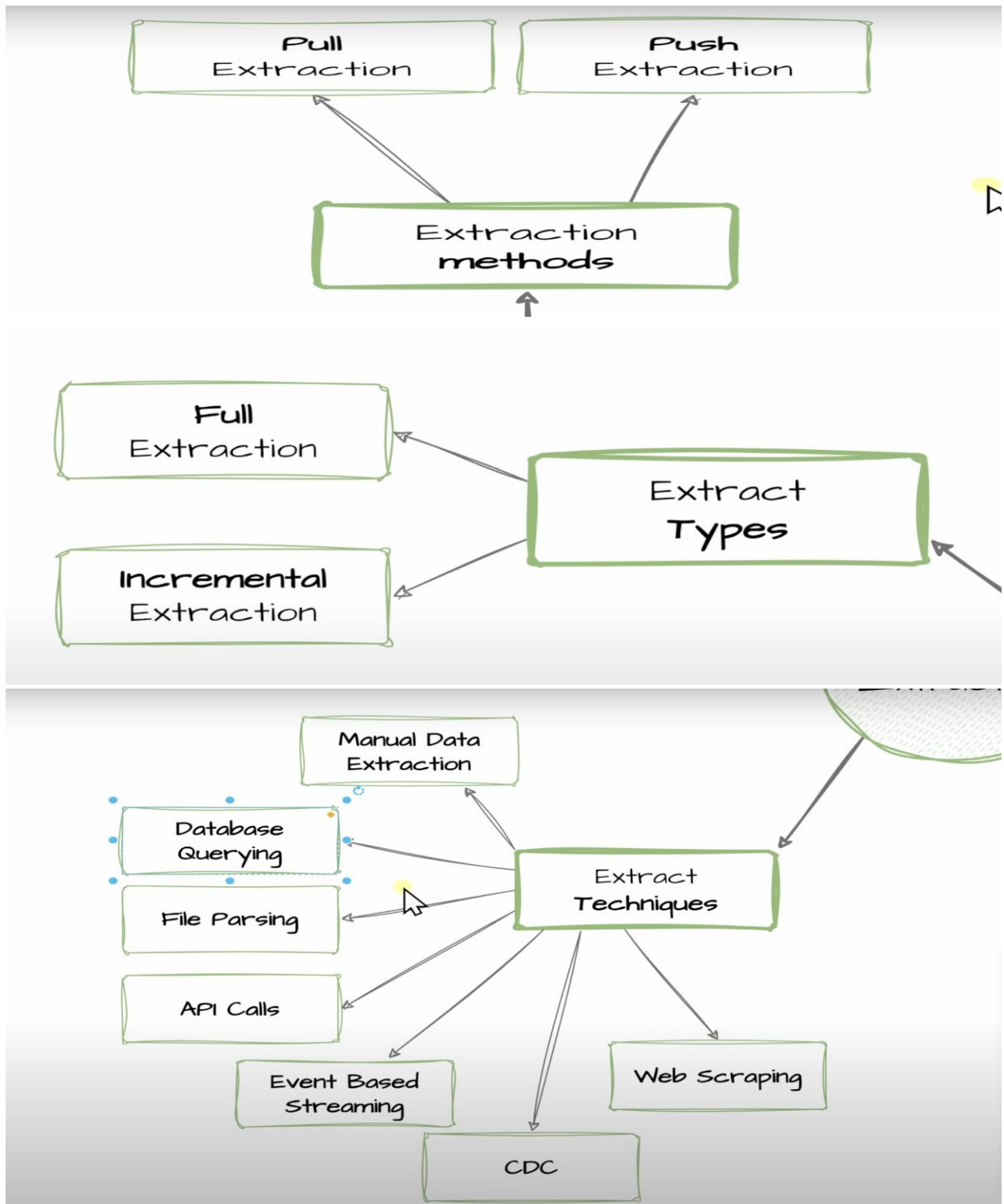
The example of a data warehouse:

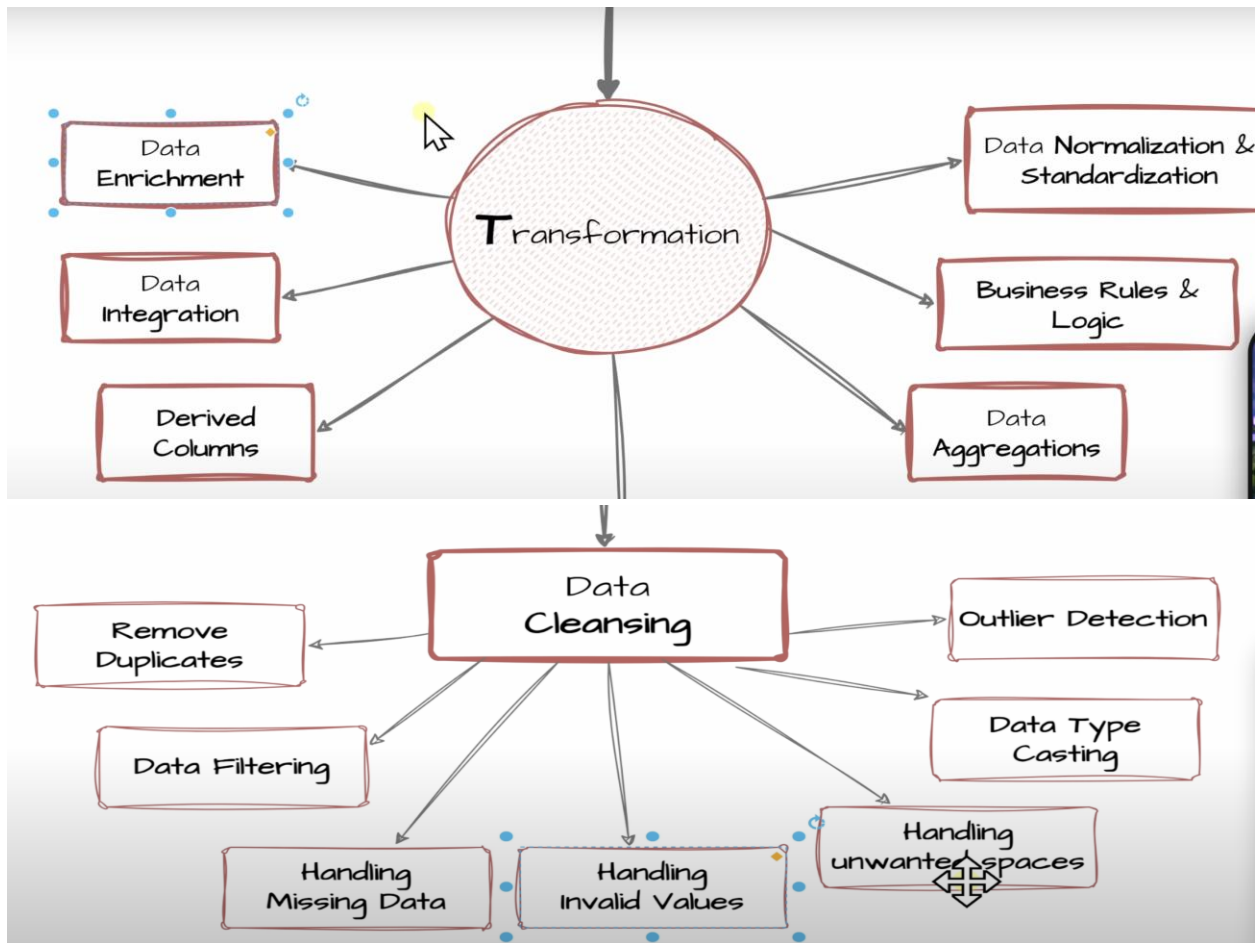Ans: The process of collecting data from to store it into a Datawarehouse is called ETL.



Extract

## Extraction methods

Pull Extraction

Push Extraction

Extraction methods

## Extract Types

Full Extraction

Incremental Extraction

Extract Types

## Extract Techniques

Manual Data Extraction

Database Querying

File Parsing

API Calls

Event Based Streaming

CDC

Web Scraping

Extract Techniques

Transformation

## Transformation

- Data Enrichment
- Data Integration
- Derived Columns
- Data Normalization & Standardization
- Business Rules & Logic
- Data Aggregations

## Data Cleansing

- Remove Duplicates
- Data Filtering
- Handling Missing Data
- Handling Invalid Values
- Handling unwanted spaces
- Outlier Detection
- Data Type Casting

Load

## Processing Types

- Batch Processing
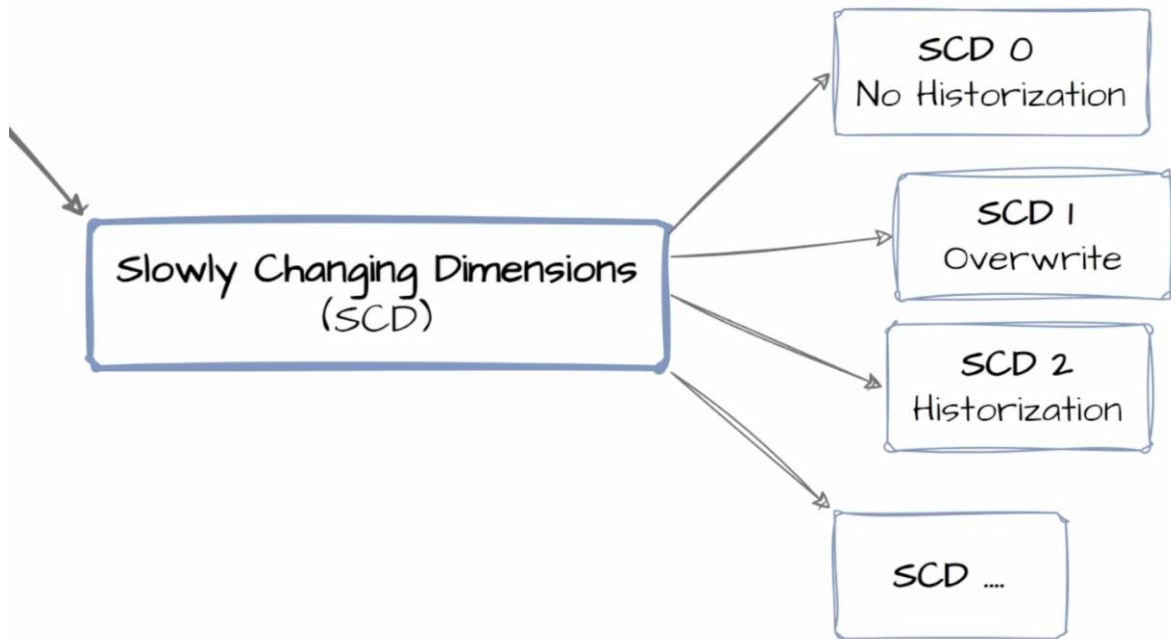- Stream Processing
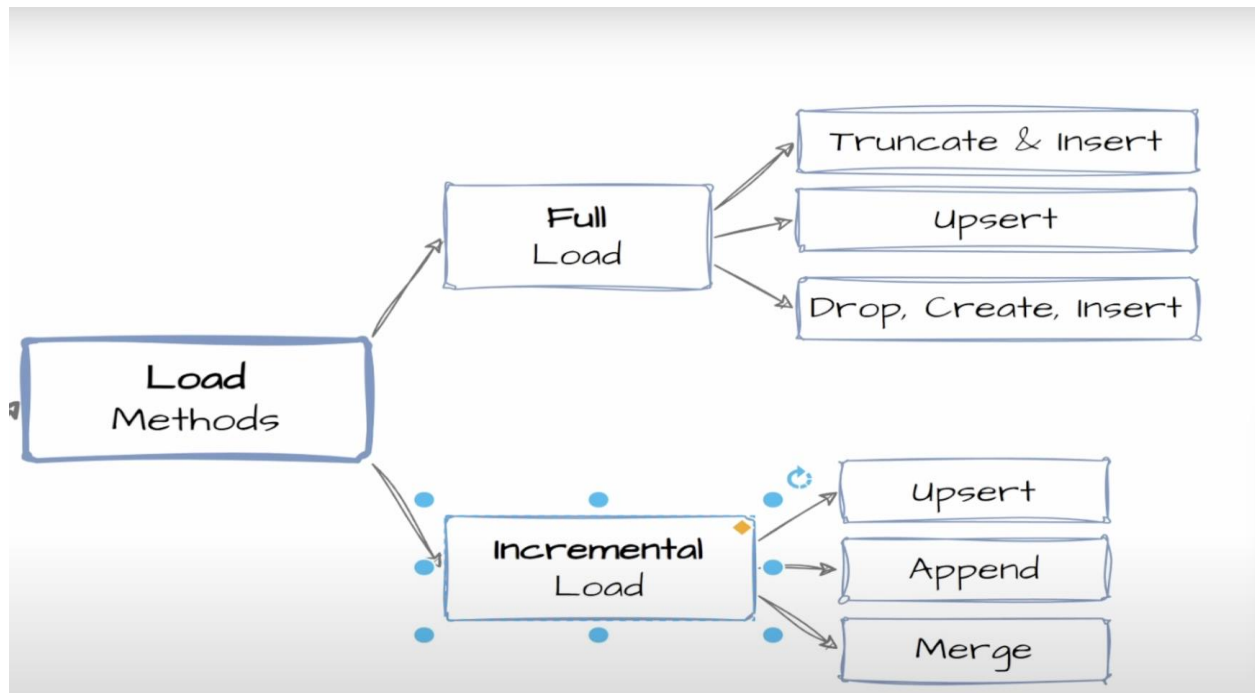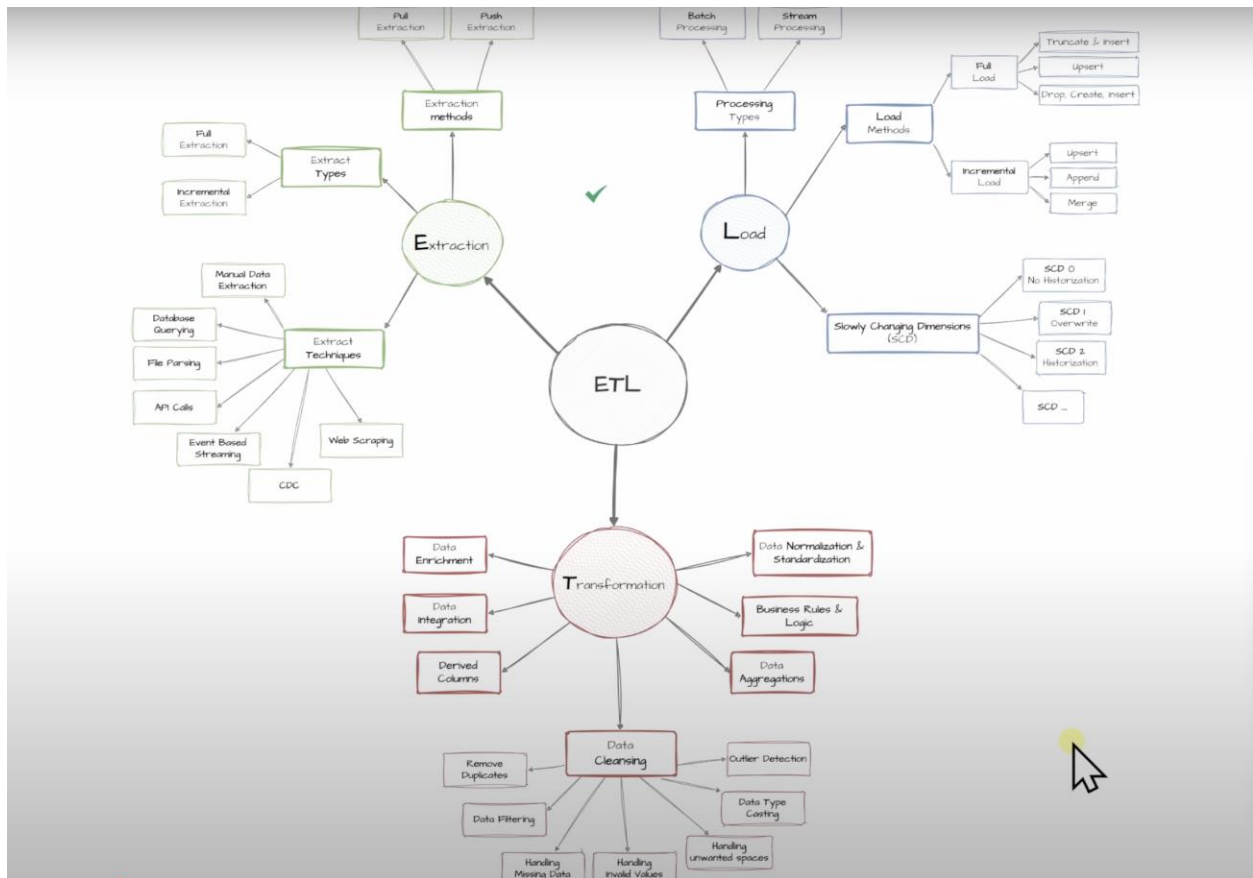
ETL:

- Here in this project, we are going to do Extraction methods such as pull and Extraction type as Full Extraction and Extraction technique as File Parsing.
- As part of transformation, we are going to do all methods of transformation.

**Project Tools**

- Notion: It helps to organize the ideas, plans,resources into a single place
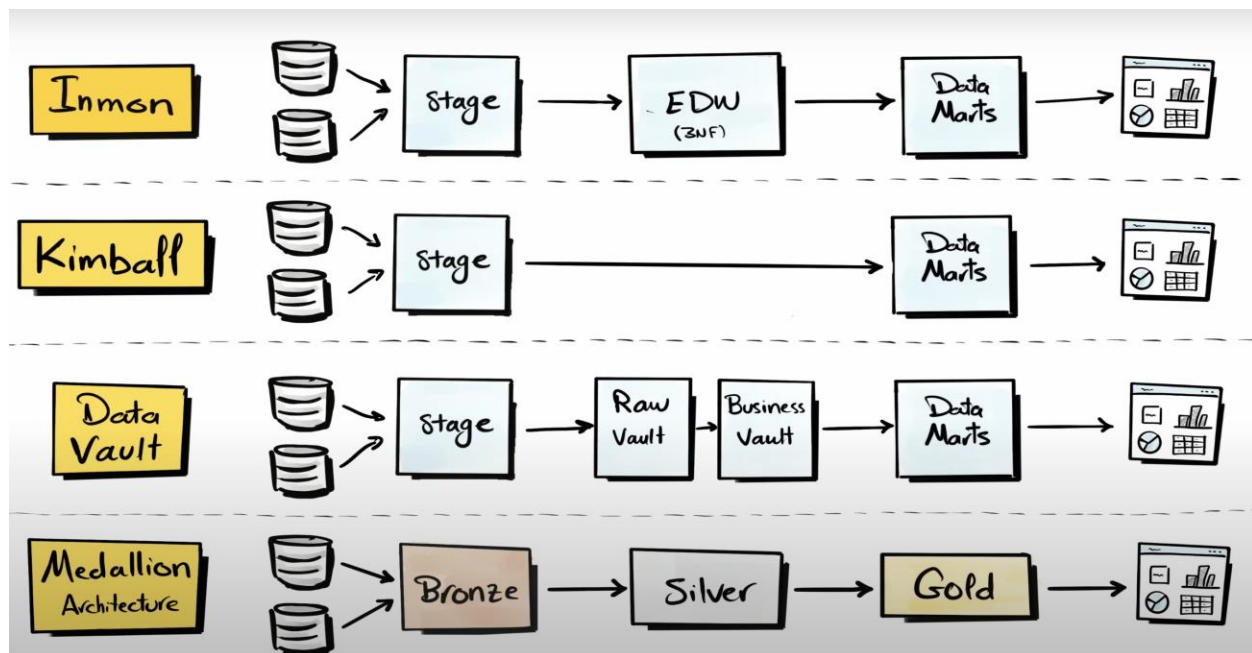
# Steps Of projects (Epics)

1) Requirement Analysis

- Before starting the work, we should be clear about the requirements of the business.
- We should be clear about the objective

2) <u>Design the Data achitechture</u>
- Choose the Data management approach
- Design rough layers of architecture
- Drawing the data architecture (There is no thumb rule for creating the architecture, according to the type of architecture you are choosing you can draw the diagram.

3) <u>Create a Project plan/Project Initialization</u>
- Create a detail Project task
- Define the project naming convention
- Create a Git Repo for the project
- Create Database and Schema.

4) <u>Buliding bronze Layer</u>
- Understand/Analyze the source system
- Coding: data ingestion
- Validating: Data completeness and Schema test
- Document: Draw the data flow
- Commit the code into Git

5) <u>Building Silver Layer</u>
- Analyze: Explore and understand the data
- Coding: Data cleaning (Check Quality, Transformation, Insert into silver
- Validating: Data Correctness checks
- Docs: Documenting and versioning to git

6) <u>Building Gold Layer</u>
- Analyzing: Explore and understand the business logic
- Coding: Data integration (Build the business object, choose table type (dimension, measure), convert the names into the friendly names)
- Validation: Data quality checks
- Docs & Version: Documenting and versioning to git

## ==Data Architecture:==

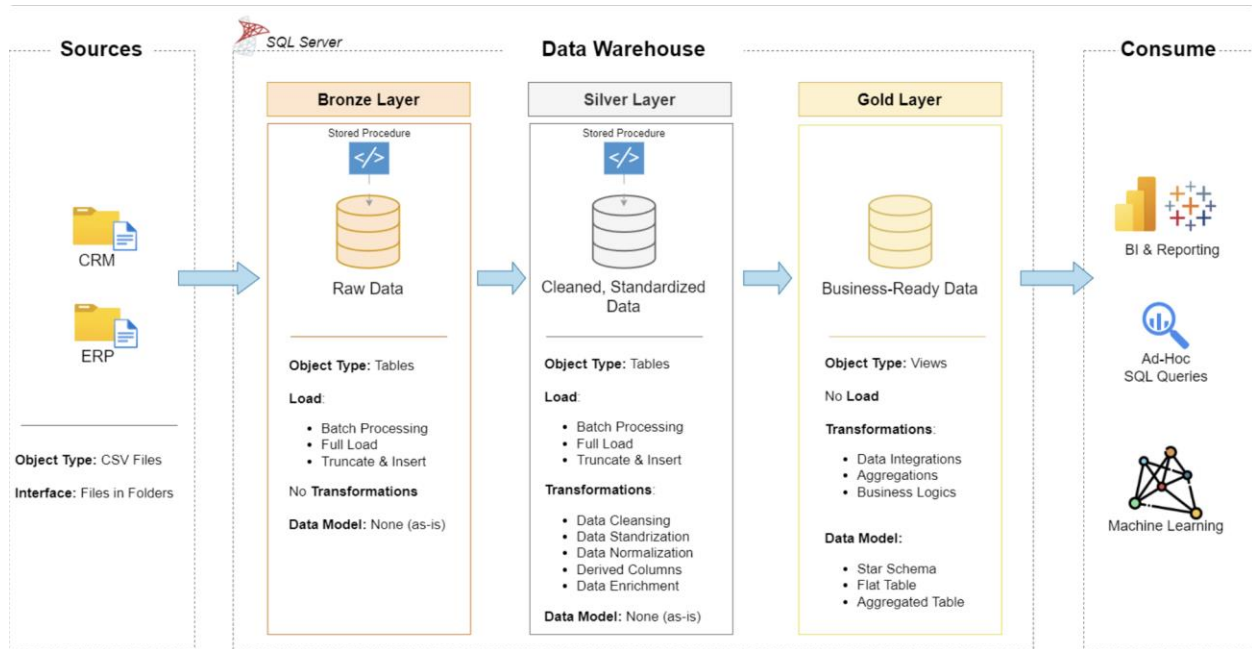Below is the diff type of Data architecture:

- Here we are selecting Medallion as the type of architecture, so below are the details of Medallion type.

Medallion Approach:

| | 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) | - Data Cleaning<br>- Data Standardization<br>- Data Normalization<br>- Derived Columns<br>- Data Enrichment | - Data Integration<br>- Data Aggregation<br>- Business Logic & Rules |
| Data Modeling | None (as-is) | None (as-is) | - Start Schema<br>- Aggregated Objects<br>- Flat Tables |
| Target Audience | - Data Engineers | - Data Analysts<br>- Data Engineers | - Data Analysts<br>- Business Users |

## High-level Architecture

## Project Initialization

- After creating the project task, we can go for proper naming convention of each objects (Database, Schema, Tables, Store procedure).
- There are certain rules to name the project (Camel Case, snake case, etc.)
- Bronze and Silver level:

  - 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 Level

  - 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.

- Other Naming convention

## 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.

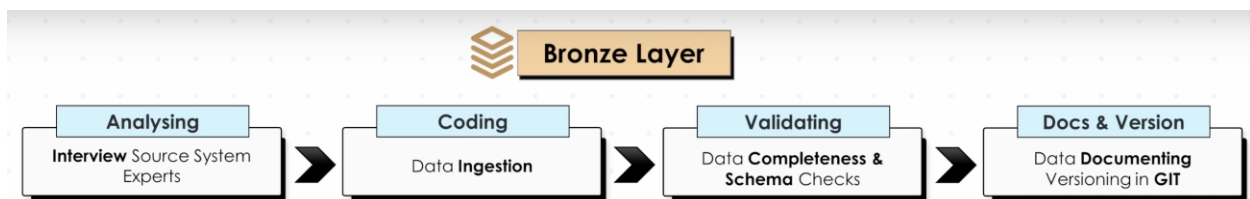## Stored Procedure

- 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.

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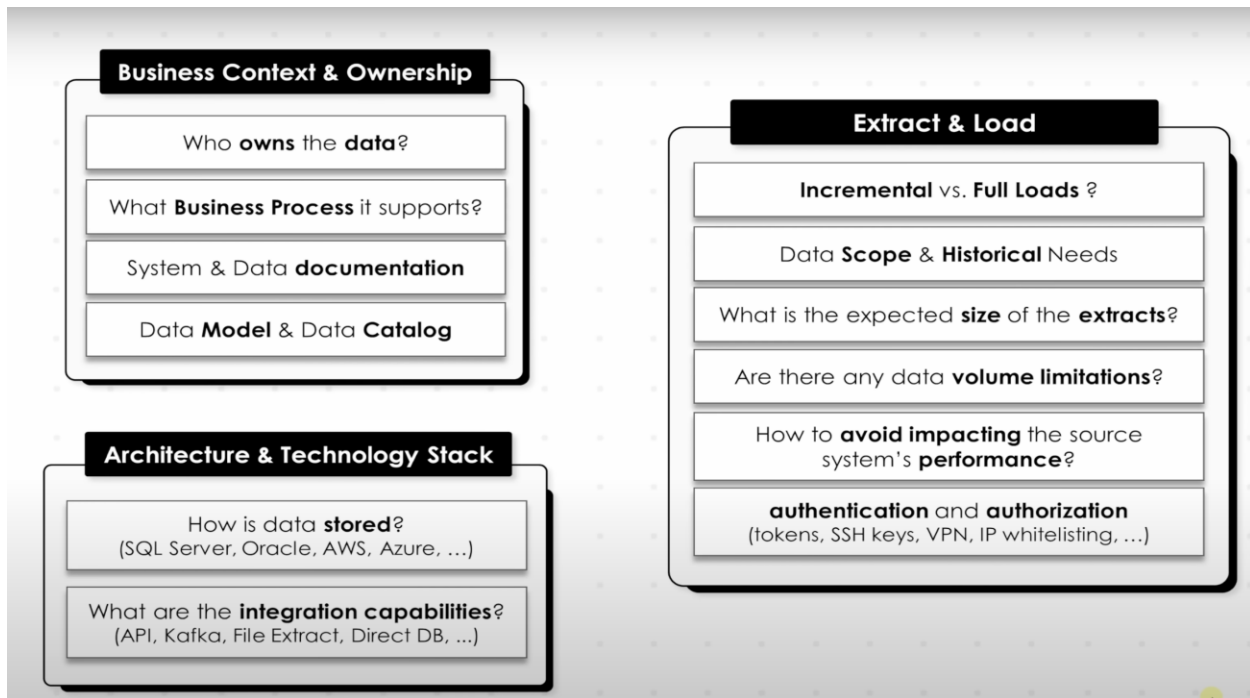
- All primary keys in dimension tables must use the suffix `_key`.
- `<table_name>_key`
  - `<table_name>` : Refers to the name of the table or entity the key belongs to.
  - `_key` : A suffix indicating that this column is a surrogate key.
  - Example: `customer_key` → Surrogate key in the `dim_customers`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

- Create a proper git folder and prepare the repo structure for the project.
- Create Database in the any DB tool(MS SQL, etc. )
- Create schemas (Gold, Silver, Bronze) under Database.
- Now we should have an empty database and schemas.

## Building Bronze layer



1. One should have the proper understating about below question on the source system

2. Create tables according to the source system data and load data



(Above query for creating table)

- Create the same DDL for all the available tables, then verify all the source system tables that have been created in the database.

- Then By bulk import method we can load multiple data from the source system to the database.
- Write the query for all the table

```
TRUNCATE TABLE bronze.erp_loc_a101;
BULK INSERT bronze.erp_loc_a101
FROM 'C:\sql\dwh_project\datasets\source_erp\loc_a101.csv'
WITH (
    FIRSTROW = 2,
    FIELDTERMINATOR = ',',
    TABLOCK
);
```

(Sample query for loading the data)

3. If you need to need to run the code on daily basis or frequently then better to write a store Procedure for this.

(Sample for store procedure)

```
CREATE OR ALTER PROCEDURE bronze.load_bronze AS
BEGIN
    TRUNCATE TABLE bronze.crm_cust_info;
    BULK INSERT bronze.crm_cust_info
    FROM 'C:\sql\dwh_project\datasets\source_crm\cust_info.csv'
    WITH (
        FIRSTROW = 2,
        FIELDTERMINATOR = ',',
        TABLOCK
    );
```
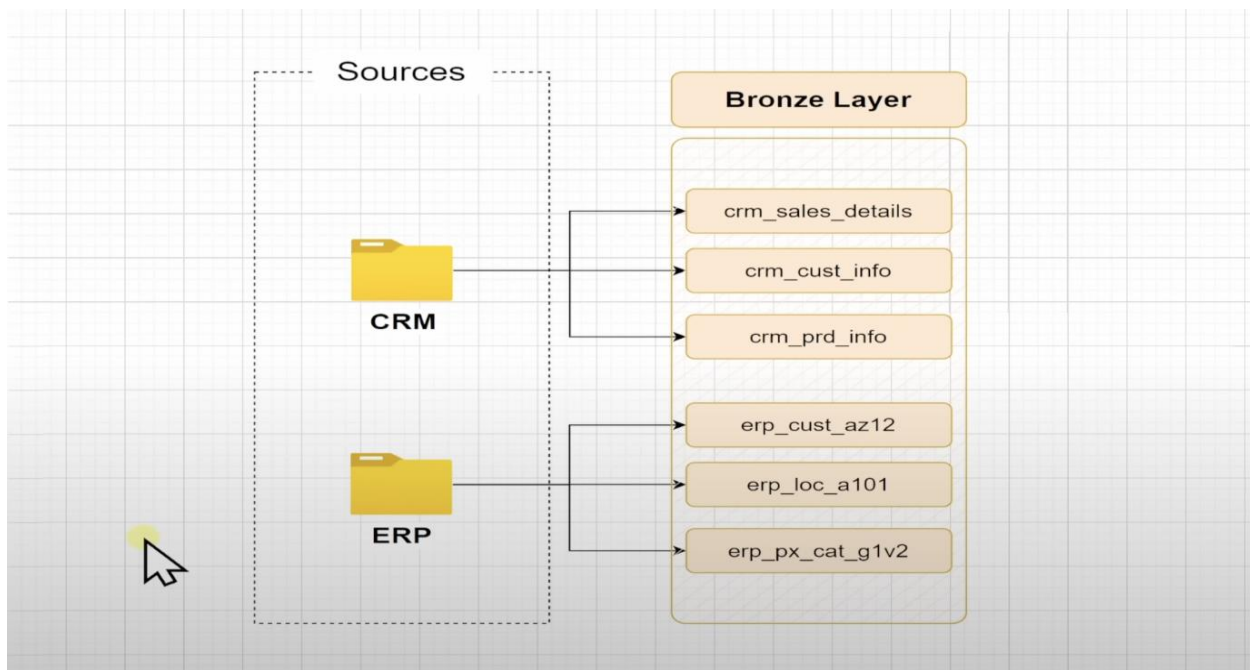
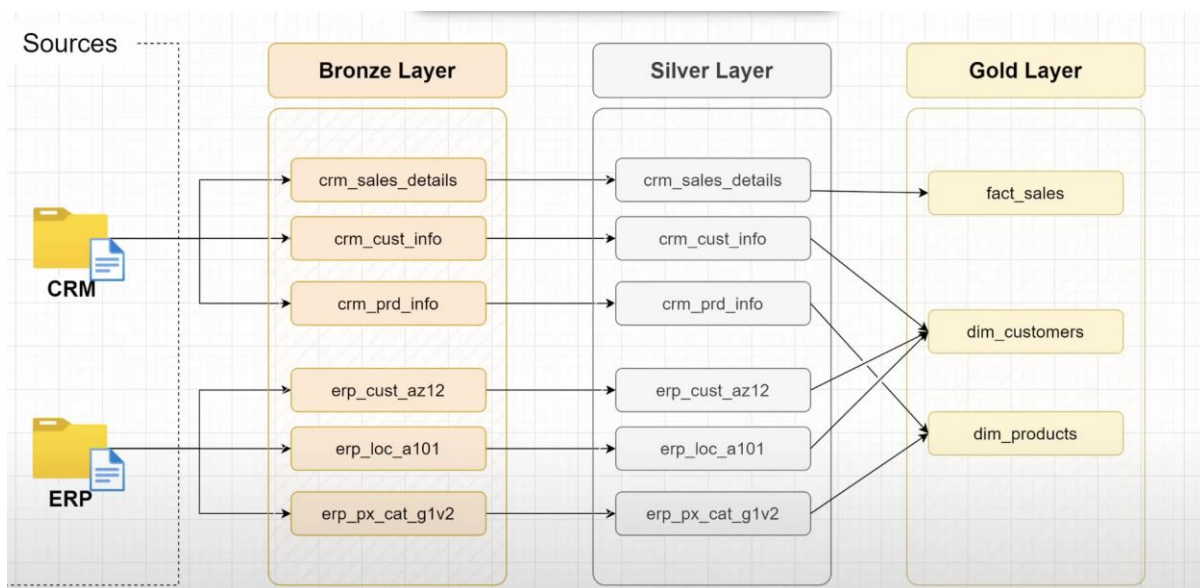Then we need to run only the below query, no need to run all the table query frequently
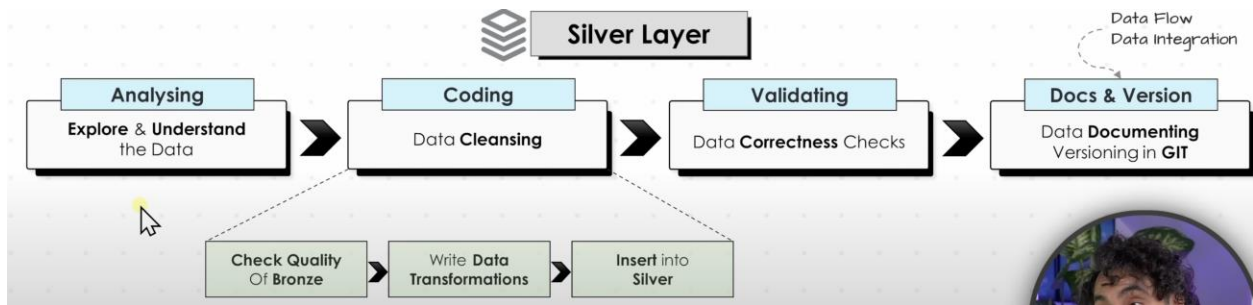
```
EXEC bronze.load_bronze
```

4. Data Flow

(Data flow for all the layers)



**Building Silver layer**

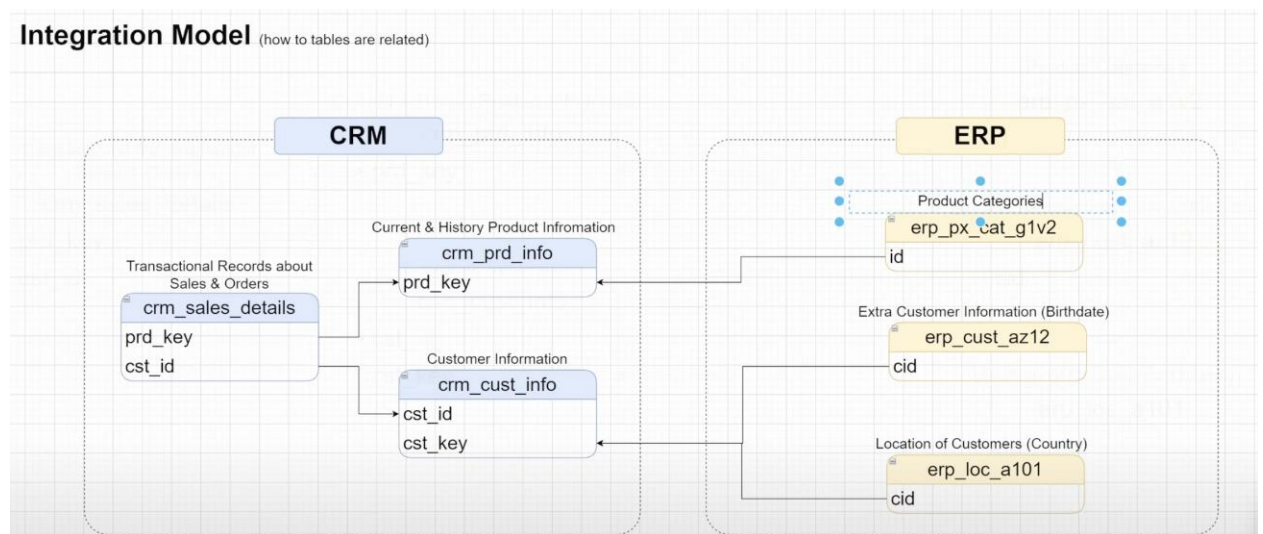1. Analyzing (Explore & understand the data)

- Query to see the data of all the table

(Like this we can see the data)

```
SELECT TOP 1000 * FROM bronze.crm_cust_info
SELECT TOP 1000 * FROM bronze.crm_prd_info
SELECT TOP 1000 * FROM bronze.crm_sales_details
```

- Create data model like this



2.Coding (Data cleansing)

- Writing Store procedure for silver layer to load the table in silver layer

(below is the sample query for the one of the table)

```sql
IF OBJECT_ID ('silver.crm_cust_info', 'U') IS NOT NULL
    DROP TABLE silver.crm_cust_info;
CREATE TABLE silver.crm_cust_info (
    cst_id INT,
    cst_key NVARCHAR(50),
    cst_firstname NVARCHAR(50),
    cst_lastname NVARCHAR(50),
    cst_material_status NVARCHAR(50),
    cst_gndr NVARCHAR(50),
    cst_create_date DATE,
    dwh_create_date DATETIME2 DEFAULT GETDATE()
);
```

- Then we can go start cleaning the data from the bronze layer and then add in the silver layer.
    1. Check for Nulls or duplicates in primary key
    2. Replace Invalid value
    3. Check for unwanted space
    4. Data standardization and consistency
    5. Convert the low cardinality column into meaningful column
    6. Convert to date time function
- We must perform these above steps for all the table present in the source
- After cleaning the data, we can perform the data transform in the tables
- Then we can load the data into the silver layer. (Here we are doing fool load)
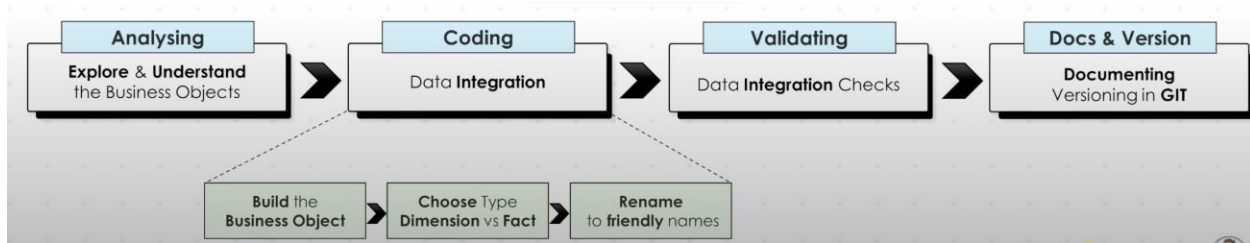  (Use truncate before running the insert query)

## 3. Validating (data correctness check)

- While writing the store procedure for silver layer, make sure we are keeping all our code in the try catch block
- We need to give proper print statement to check all the tables get inserted or not
- We can also add a couple of methods like time, etc. to check the total time required to execute the procedure.
- We must use error handling as well

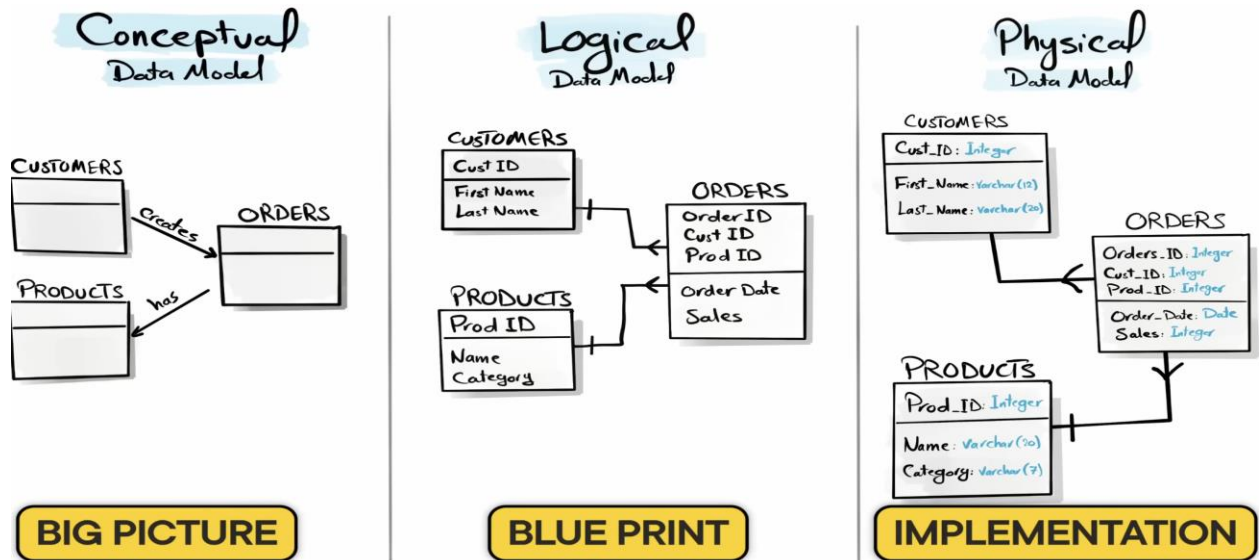## 4. Doc & Version (Data documenting and pushing into git)

- Create a file under script to upload the Store procedure query
- Create a folder/file name as test to upload all the data cleaning or transformation query that we wrote to clean the tables.
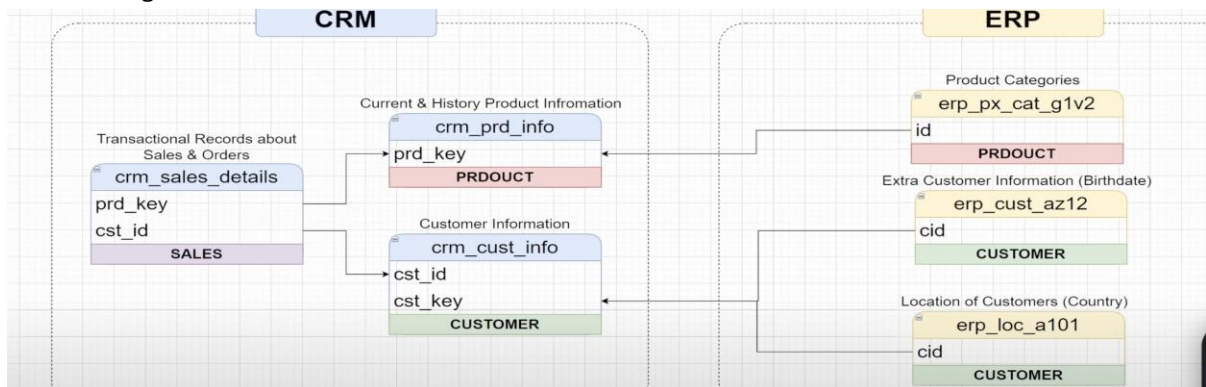
**Building Gold layer**

# 1. Analyzing: Explore & Understand Business



- 
- Above is the type of data model present, but here we are using Logical layer model,
- We are using Star schema here



- 
  Here we are adding the labels at the bottom of the tables to make it easy to understand the tables' content.

## 2&3. Data Integration and Validation
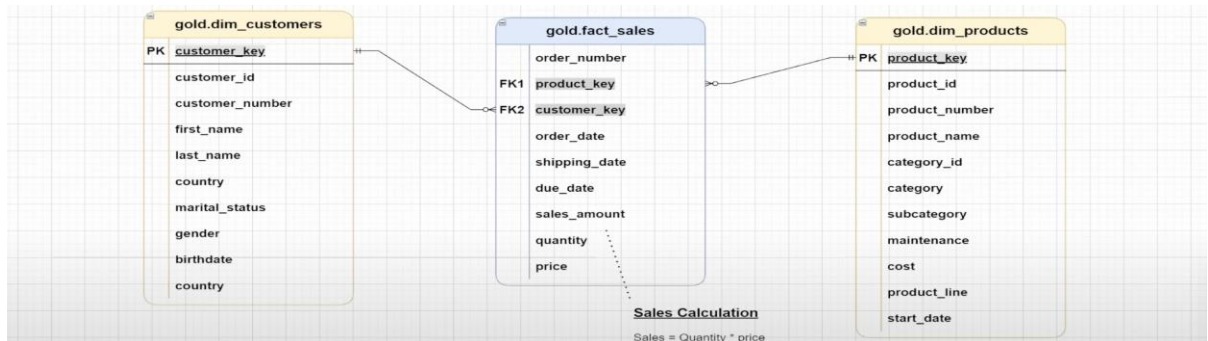
- Here we are collecting the customer information from two source system (erp, crm)

```
SELECT
    ci.cst_id,
    ci.cst_key,
    ci.cst_firstname,
    ci.cst_lastname,
    ci.cst_marital_status,
    CASE WHEN ci.cst_gndr != 'n/a' THEN ci.cst_gndr  -- CRM is the Master for gend
         ELSE COALESCE(ca.gen, 'n/a')
    END AS new_gen,
    ci.cst_create_date,
    ca.bdate,
    la.cntry
FROM silver.crm_cust_info ci
LEFT JOIN silver.erp_cust_az12 ca
ON        ci.cst_key = ca.cid
LEFT JOIN silver.erp_loc_a101 la
```

- Convert the column name to some easy and understandable name
- Make the proper order of the column.
- Find out whether it is a dimension or a fact table. Then start creating the dimension table and the fact table.
- Make sure we have some primary key present, if not then create it on basis of existing column.
- Then create a view of the gold level.
- Create a data model as per the table (Dimension & Fact) that has been created
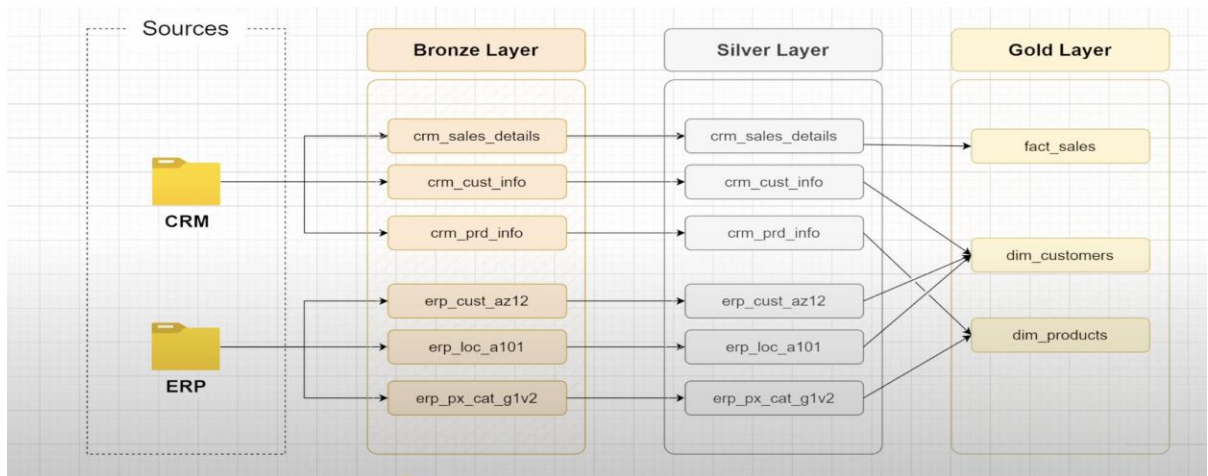


- Create a data catalogue (Catalogue describes the details of the data like column info, connection type, etc.)

### 1. gold.dim_customers

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

| Column Name | Data Type | Description |
| --- | --- | --- |
| customer_key | INT | Surrogate key uniquely identifying each customer record in the dimension table. |
| customer_id | INT | Unique numerical identifier assigned to each customer. |
| customer_number | NVARCHAR(50) | Alphanumeric identifier representing the customer, used for tracking and referencing. |
| first_name | NVARCHAR(50) | The customer's first name, as recorded in the system. |
| last_name | NVARCHAR(50) | The customer's last name or family name. |
| country | NVARCHAR(50) | The country of residence for the customer (e.g., 'Australia'). |
| marital_status | NVARCHAR(50) | The marital status of the customer (e.g., 'Married', 'Single'). |
| gender | NVARCHAR(50) | The gender of the customer (e.g., 'Male', 'Female', 'n/a'). |
| birthdate | DATE | The date of birth of the customer, formatted as YYYY-MM-DD (e.g., 1971-10-06). |
| create_date | DATE | The date and time when the customer record was created in the system |

- Draw the final data flow diagram

- Then commit all the script to the Git