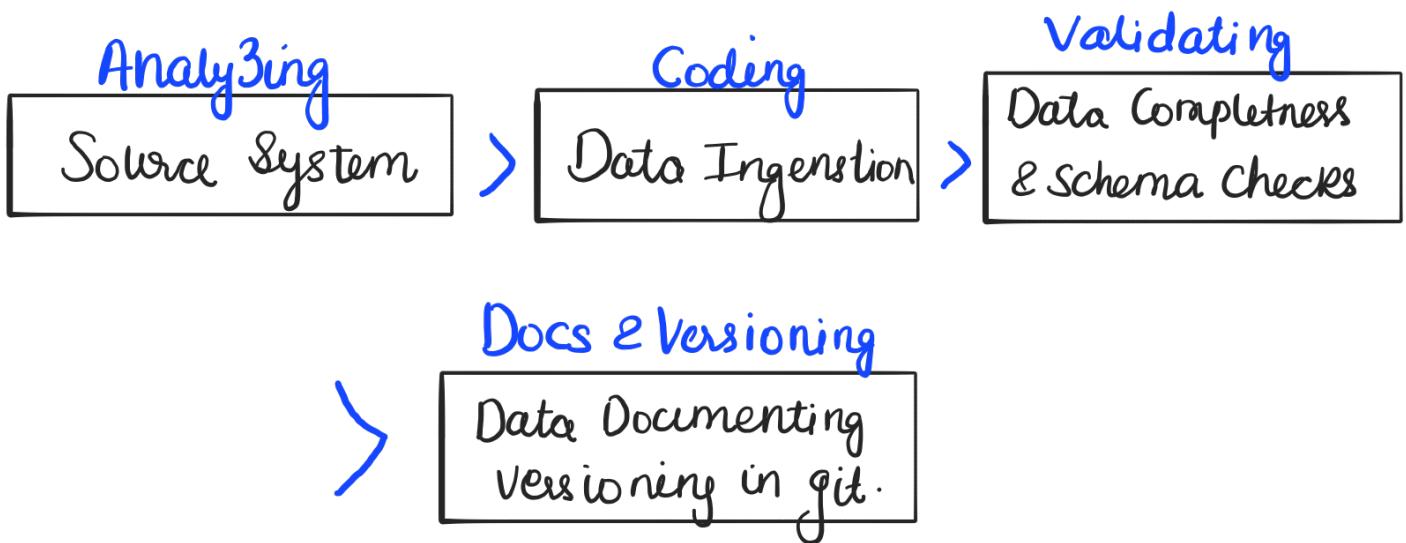


PROJECT NOTES

- > Initializing git hub and cloning the repo, or make a local git repo.
- > Setup Draw.io (For daigrams).

Bronze Layer



- > Business Context and Ownership
 - who owns the data ?
 - what business process does it supports.
 - System & data documentation
 - Data Model & Data catalog.

- > Architecture & Technology stack.
 - How is data stored
 - what are the integration Capabilities.

- > Extract & Load
 - Incremental Vs Full Load
 - Data Scope & Historical Needs
 - what is the expected size of the extracts
 - Are there any data volume limitations.
 - How to avoid impacting the source system's performance.
 - Authentication and authorizations.

1. Create TABLES (ODL scripts), for all the files,
CRM, exp.

If object-id ('name', 'U') is not null
DROP name;

CREATE TABLE name (

C1 DT,

C2 DT,

...
...);

2. TRUNCATE TABLE table_name

3. BULK INSERT table_name

FROM 'Path\filename'

WITH (

FIRSTROW = 2, (Row=1, Contains info about fields

FIELDTERMINATOR = ',' net data)

TABLOCK

);

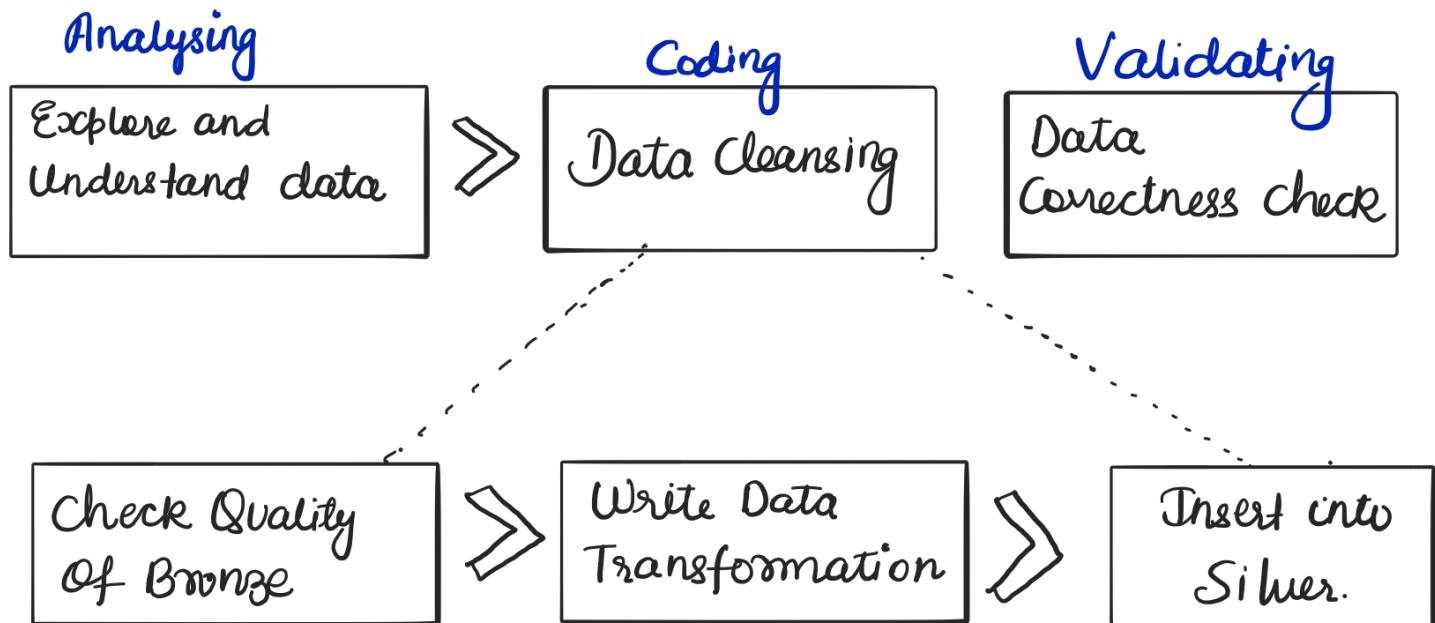
→ Delimiter (separator of values).

- > For data quality check all the fields, and count the no of rows.
- > Repeat for all the files.

For Load Scripts:

1. Make stored Procedure.
2. Use messages
3. Use Try & Catch.
4. Use Load Time. (for each table using variables).
5. Use Load time for each source or entire process.
6. Draw
7. Push Code. to git.

Silver Layer



1. Data exploration, try to create links between the table on draw.io
2. Create tables (DDL), just copy DDL from the bronze Layer and change the schema to Silver (ctrl+F).
add dwi-create_date datetime2 default getdate();

> Metadata Columns:

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

create_date : The records load timestamp.

update_date : The records last update timestamp.

source_system : The origin system of the record.

file_location : The file source of the record.

3. Check for duplicates and null values in the column you want to make a primary key.
(In bronze layer).

> Expectation : No Results.

> pick the value with the highest date
(or make a standard criteria).

4. Check for unwanted spaces in string values.

Expectation : No results.

(TRIM)

5. Data Standardization & Consistency :

(CASE), eg : M → Male | M - Married
 F → Female | S - Single

(UPPER).

6. Insert new data into the silver schema.

7. bronze. erp C filter out unmatched data after applying transformation, check for symbols.
(-), (_), replace if inconsistent.

8. Extract data two columns from 1 column.

SUBSTRING (C, s.v, Len(C)) .

9. Prod line (start_date, end_date, ambiguous),
decide a Rule, then apply Rule,
Start of 2nd is end of 1st. { Lead - 1 }
OR, (dateadd (-1) day) .

10. Cast datetime as date.

Derived Columns:

- > Create new columns based on calculations or transformation of existing ones.
- > Data enrichment - adding new data to the table.

Table CRM_Sales_Details:

1. Date quality check ≤ 0 OR length $\neq 8$ OR
 - > Current_date OR $<$ Start_date.
2. Can't cast int to date, so into varchar.
then date int \rightarrow Varchar \rightarrow date.
3. Sales order date cannot be greater than
Ship date & vice versa, same for due
dates.
4. Sales = quantity \times Price

Quantity $\neq 0, -ve, \text{NULLS}$

5. For price and quantity related issues,

- > Fix the source system
- > Fix the issue in data warehouse.
 - > Pick a Rule.
 - > Handle division by 0
 $\text{NULLIF}(\text{col}, 0)$.

Table : Bronze .erb - cert - qz12 :

- > Remove NAs from the cid using CASE.
- > Check D.O.B Range. ($\text{olddate} < > \text{geldate}$)
 - > geldate with NULL.
- > Cleanup gender column (first distinct) then convert it to Male, Female, NRA.
 - (as when $\text{IN}('F', 'FEMALE')$ - Female)

Table : exp - loc - cat

- clean cid, replace '-' with ''.
- Normalize the country.

Table : exp - px - cat - g1v2

- check cat for empty space , subcat, maintenance.
- check the distinct Cat, subcategory.
- Just load.

Create Procedure : (Silver. Load-Silver).

- ① Print
 - ② Truncate
 - ③ Print inserting
 - ④ print insertion finished.
- Print message for each Section 2 Step.
 - Implement error handling.
 - Print duration of each Step.
 - Print total duration.

Note: Maintain consistency across stored Procedures .

GOLD LAYER

Analyzing

Explore & Understand
the Business objects.

Coding

Data
Integration

Validating

Data Integration
checks

Docs & Version

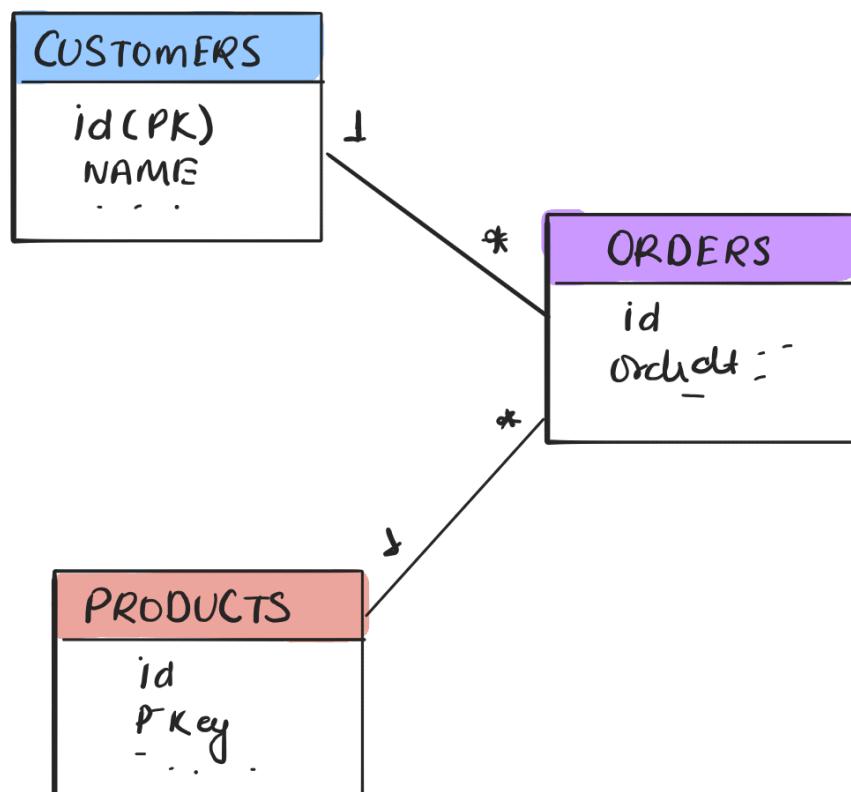
Documenting
Versioning in GIT

Build the
business object

Choose Type
Dimension Vs Fact

Rename to
Friendly names

Data Modelling:



Conceptual Model (BIG PICTURE)

- > Focus is only on entity (like what tables) and RIs between them.
- > No columns or stuff

Logical Data Model (BLUE PRINT)

- > Specify the details (columns) in these entities.
- > Also establish the relationship between the entities (how they are connected).

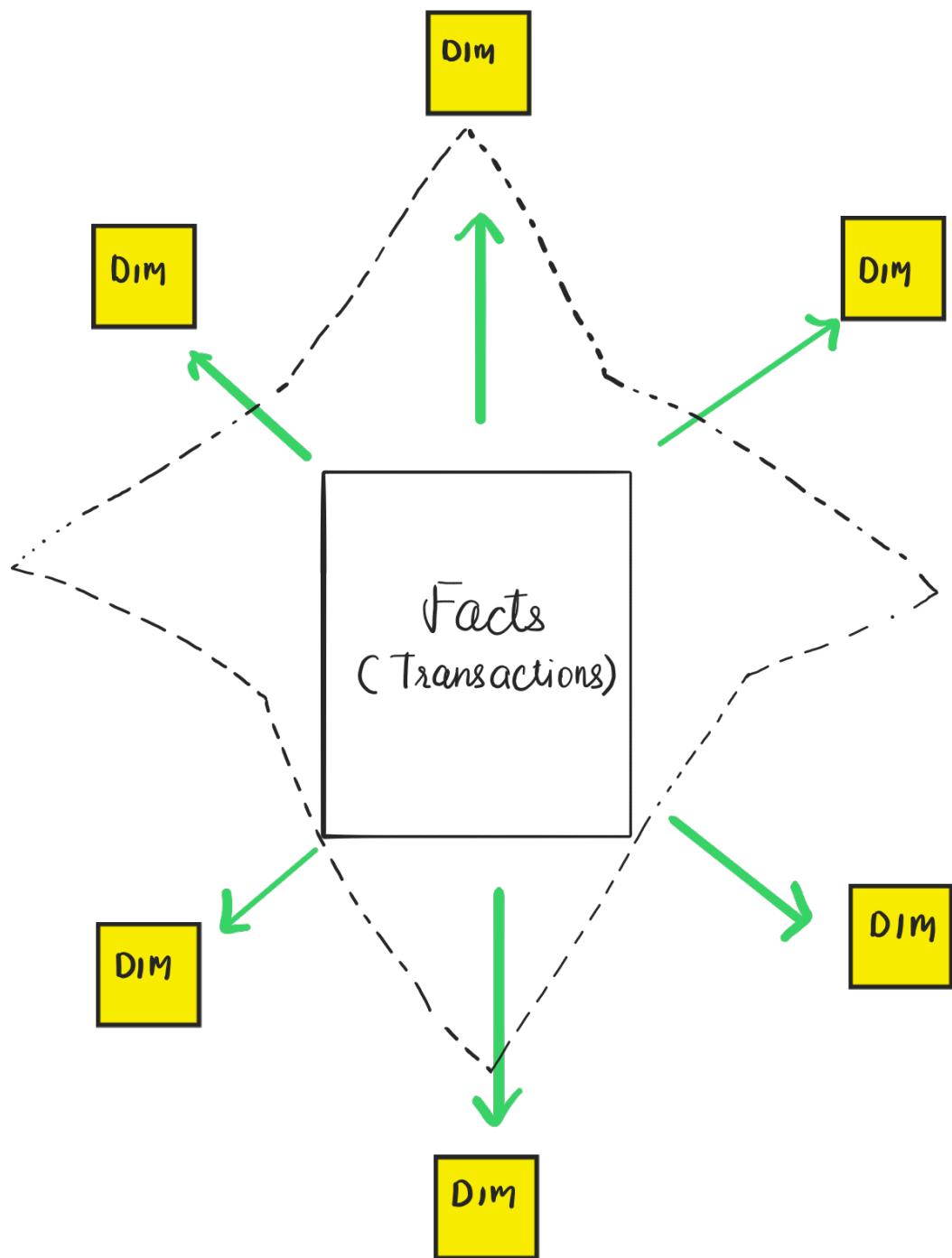
Physical Data Model (IMPLEMENTATION)

- > How the data will be stored in the database
- > Column Types.

NOTE: For this project we will build the logical Data Models.

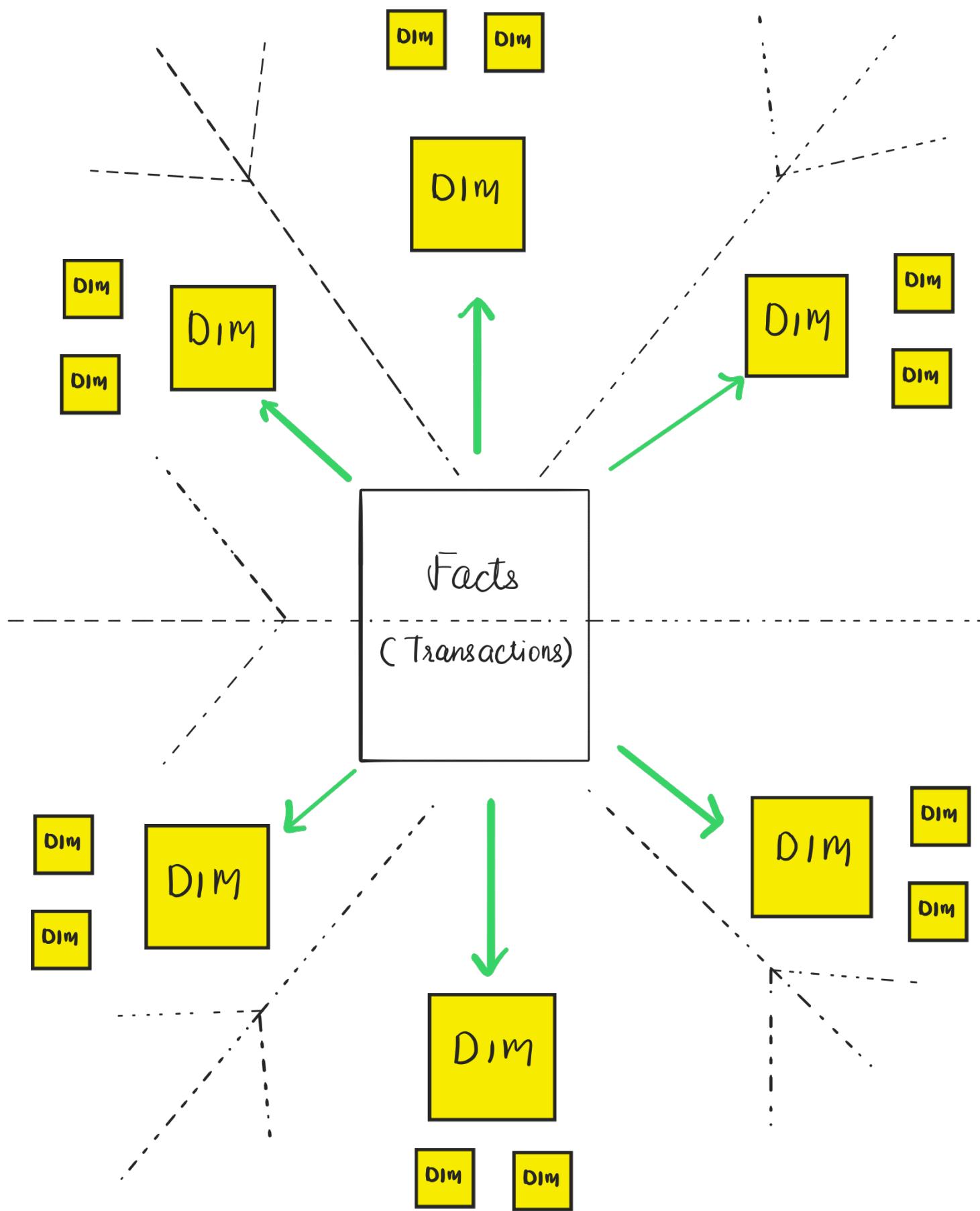
Schemas

Star Schema :



Dimensions, are description about
Transactions

Snow Flake Schema



→ Dimensions are Split into smaller dimension.

Dimensions

> Descriptive information that gives context to your data. (who, what were).

Facts

> Quantitative information that represents events (How much, how many).

Explore Business objects:

> fact Tables, (like sales, products, customers)

Create dimension Customer

> Join all tables that have cust info. and then create customer view that has all the necessary details about customer available

> If same data comes from multiple columns. then we have to do data integration.
(ASK system export). only for matching data.
(USE CASE handle NULLS) { Take any one source as master }.

> Rename columns (meaning full name)

> Assign a primary key.

Surrogate Keys : System generated unique identifier assigned to each record in table.

> DDL Based generations.

> Query-based using window function
(Row-Number).

→ In this project we have used query-based window function.

Create Dimension Product 8-

- > Join all tables with the product info.
- > Filter out historical data ,only current products.
(Connection is done using prod-key not the prod-id)
- > Check for duplicacy .
- > Arrange the columns.
- > Rename Columns.
- > Create surrogate key
- > Create view.

Create Fact Sales :

- > Facts connects multiple dimensions.
- > Use the dimension's surrogate key instead of IDs to easily connect facts with dimensions.
- > Join dim_sales_details with our views.
- > Rename column, arrange them.
- > Create view.

Steps after Creating Facts & Dimensions :

- > Check Foreign Key integrity (Right table p.k is NULL).
- > Draw data models.(gold).
- > Create Data Catalog.
- > Extend Data Flow diagram

Gotchas :

- If DDL was created earlier , and we have derived columns then update DDL for that table.
- Check for data types in DDL after performing operations in the silver Layer.
- After every transformation Run it through the quality check scripts.