# DataCo Supply Chain Data Warehousing
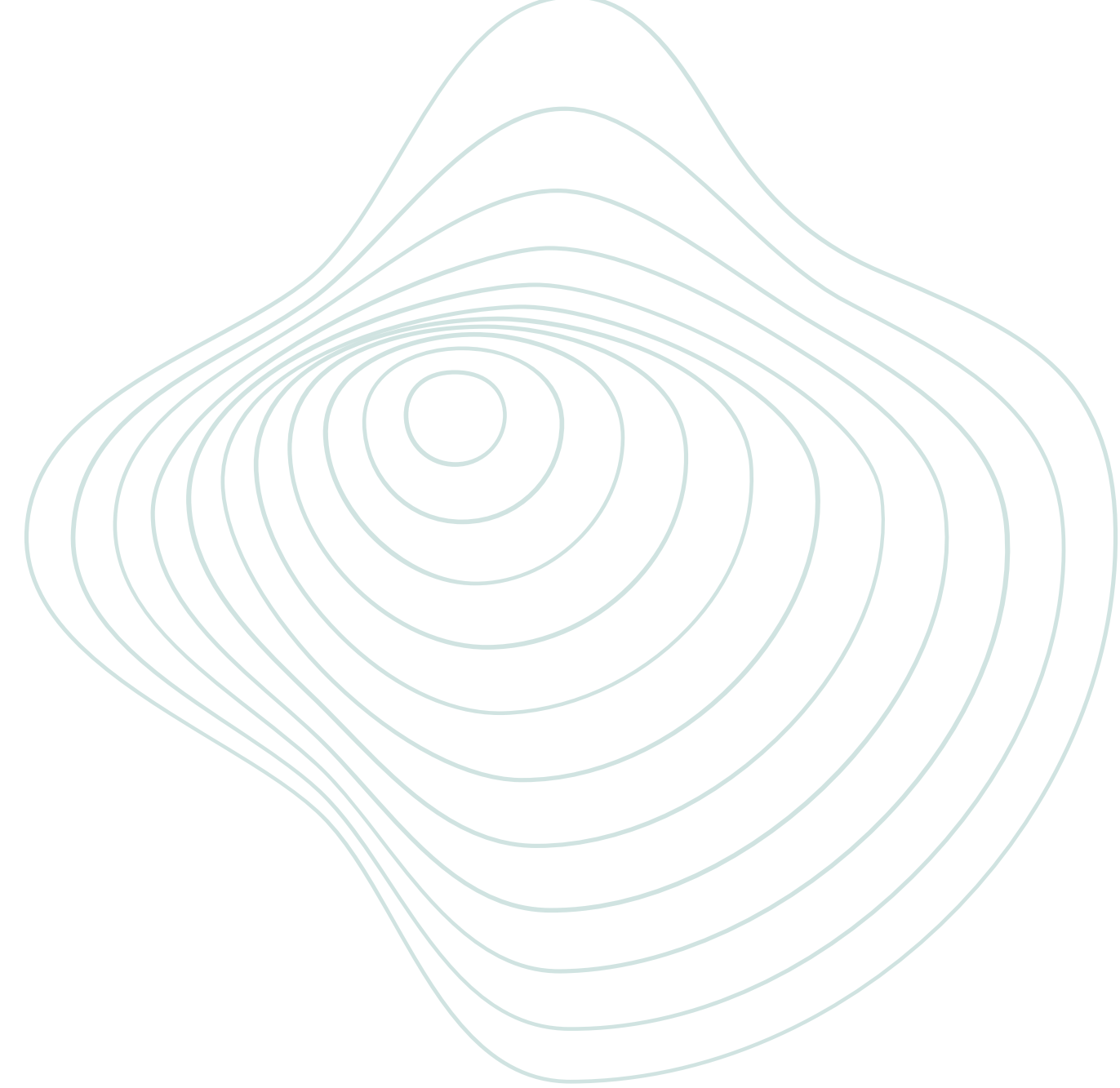
**FEUP - MECD - Data Warehouse**
**Middle presentation**

Carlos Miguel Veloso
Cátia Teixeira
Luís Henriques
Rojan Aslani

# Introduction

## Assignment Goals

- To design a data warehouse, implement it, and exemplify its use
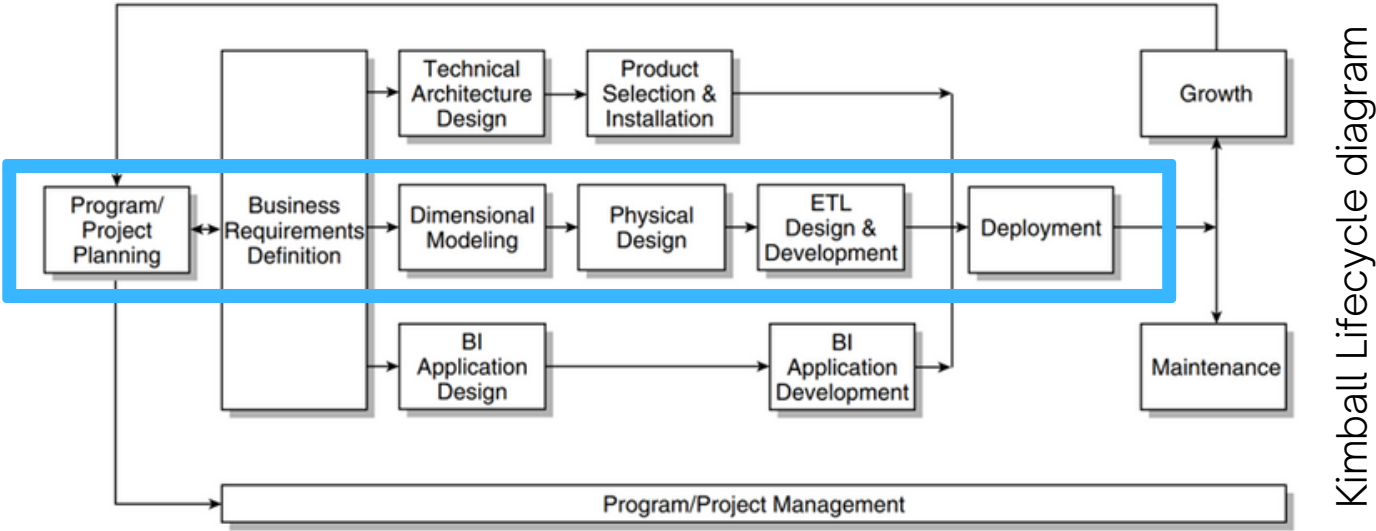
## Assignment Requirements

- The number of rows must be over 10000 with, at least, one additive measure
- There must be aggregated facts or snapshots with at least one semi-additive measure
- There must be at least 4 dimensions, one of which temporal, and some of them are common to both kinds of facts.

# Introduction



Kimball Lifecycle diagram

## Assignment steps

| Project phase | Tasks |
|---|---|
| Project Planning | • Timeline, general tasks definition and distribution<br>• Finding data |
| Business Requirements Definition | • Data understanding<br>• Scope definition and data filtering |
| Dimensional Modeling | • Relational model<br>• Entity relationship<br>• Bus Matrix<br>• Dimensional design |
| Physical Design | • Data warehouse implementation |
| ETL Design & Development | • ETL process definition<br>• Loading data do Postgres |
| Deployment | • Data analysis and business analytics |

# Source

Dataset source

Identify and profile operational data (OLTP) sources

# Dataset

**Online store transactions**

- Product data

- Financial data

- Sales and demand data

| | Original dataset | Reduced dataset |
|---|---|---|
| Columns | 53 | 47 |
| Rows | 180000+ | 27128 |
| Timespan | 2015 to 2018 | 2nd Semester of 2017 |

**kaggle**
- Data set was sourced from Kaggle platform

**POLITÉCNICO DE LEIRIA** | ESCOLA SUPERIOR DE TECNOLOGIA E GESTÃO
- Made available by Politécnico de Leiria

**DataCo**
- Data related to Datco Company

# Transactional schema

**Segment**
| | |
|---|---|
| PK | segment_id |
| | segment_desc |

**Customer**
| | |
|---|---|
| PK | customer_id |
| FK1 | city_id |
| FK2 | state_id |
| FK3 | country_id |
| FK4 | segment_id |
| | first_name |
| | last_name |
| | email |
| | zip_code |
| | latitude |
| | longitude |

**Department**
| | |
|---|---|
| PK | department_id |
| | city_id |
| | state_id |
| | country_id |
| | region_id |
| | name |

**Category**
| | |
|---|---|
| PK | category_ld |
| | name |

| | |
|---|---|
| FK1 | order_id |
| FK2 | customer_id |
| | sales_per_customer |

**Product**
| | |
|---|---|
| PK | card_id |
| FK1 | category_id |
| | description |
| | image |
| | name |
| | price |
| | status |

**Order**
| | |
|---|---|
| PK | order_id |
| FK1 | city_id |
| FK2 | state_id |
| FK3 | country_id |
| FK4 | region_id |
| FK5 | customer_id |
| FK6 | product_card_id |
| FK7 | payment_type |
| | date |
| | item_total |
| | profit_per_order |
| | status |

**(Item)**
| | |
|---|---|
| PK | item_id |
| | discount |
| | discount_rate |
| | product_price |
| | profit_ratio |
| | quantity |
| | \sales |

**Shipping**
| | |
|---|---|
| PK | shipping_mode |

**(Delivery)**
| | |
|---|---|
| FK1 | order_id |
| FK2 | shipping_mode |
| | date |
| | days_shipping_real |
| | days_shipping_scheduled |
| | late_delivery_risk |
| | delivery_status |

**Payment**
| | |
|---|---|
| PK | id |
| | type |

**City**
| | |
|---|---|
| PK | city_id |
| FK1 | state_id |
| | city |

**State**
| | |
|---|---|
| PK | state_id |
| FK1 | country_id |
| | state |

**Market**
| | |
|---|---|
| PK | market |

**Region**
| | |
|---|---|
| PK | region_id |
| FK1 | country_id |
| | region |

**Country**
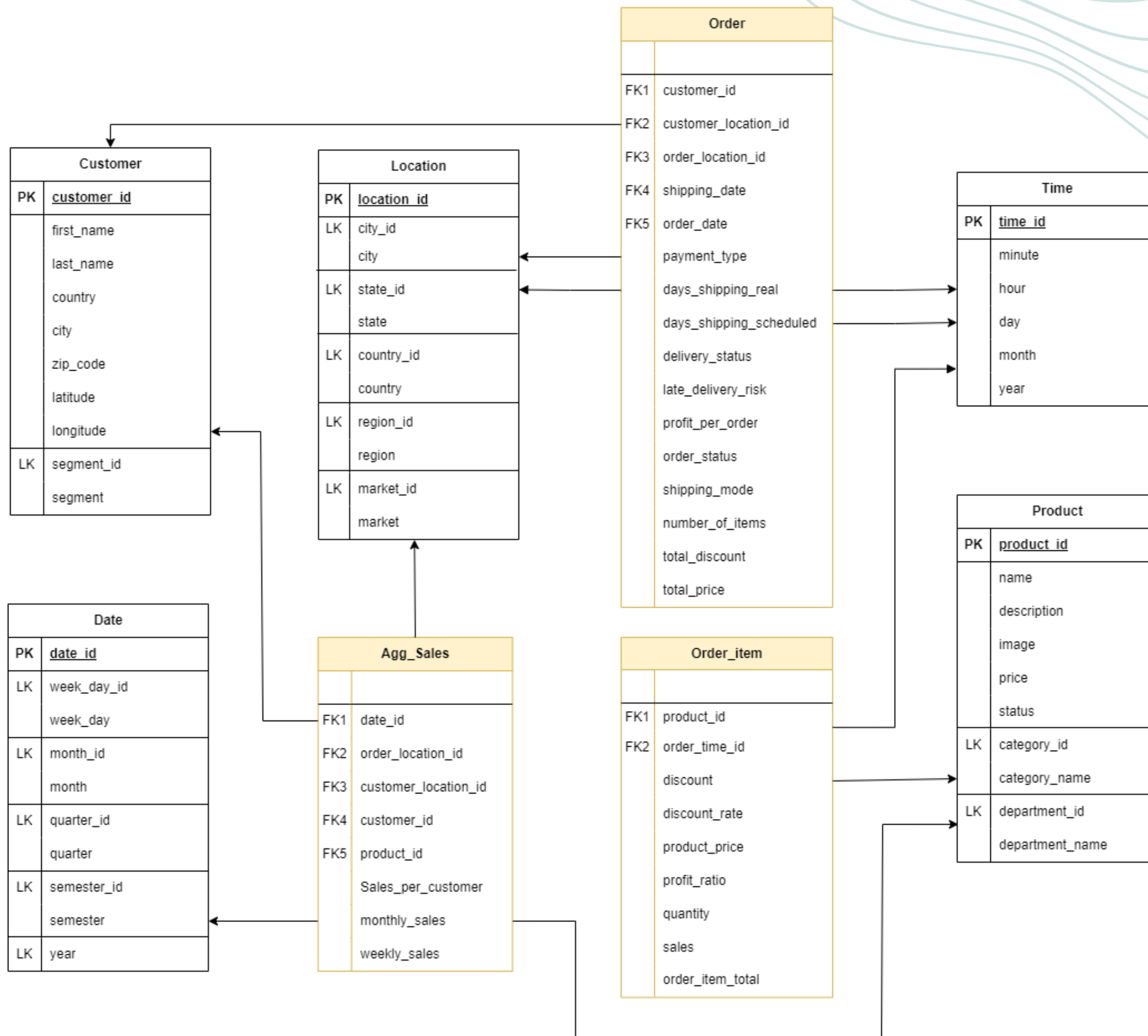| | |
|---|---|
| PK | country_id |
| | country |

6

# Dimensional Model

Develop a dimensional model that includes Dimensions and Facts

# Dimensional Model

Dimensional Bus Matrix

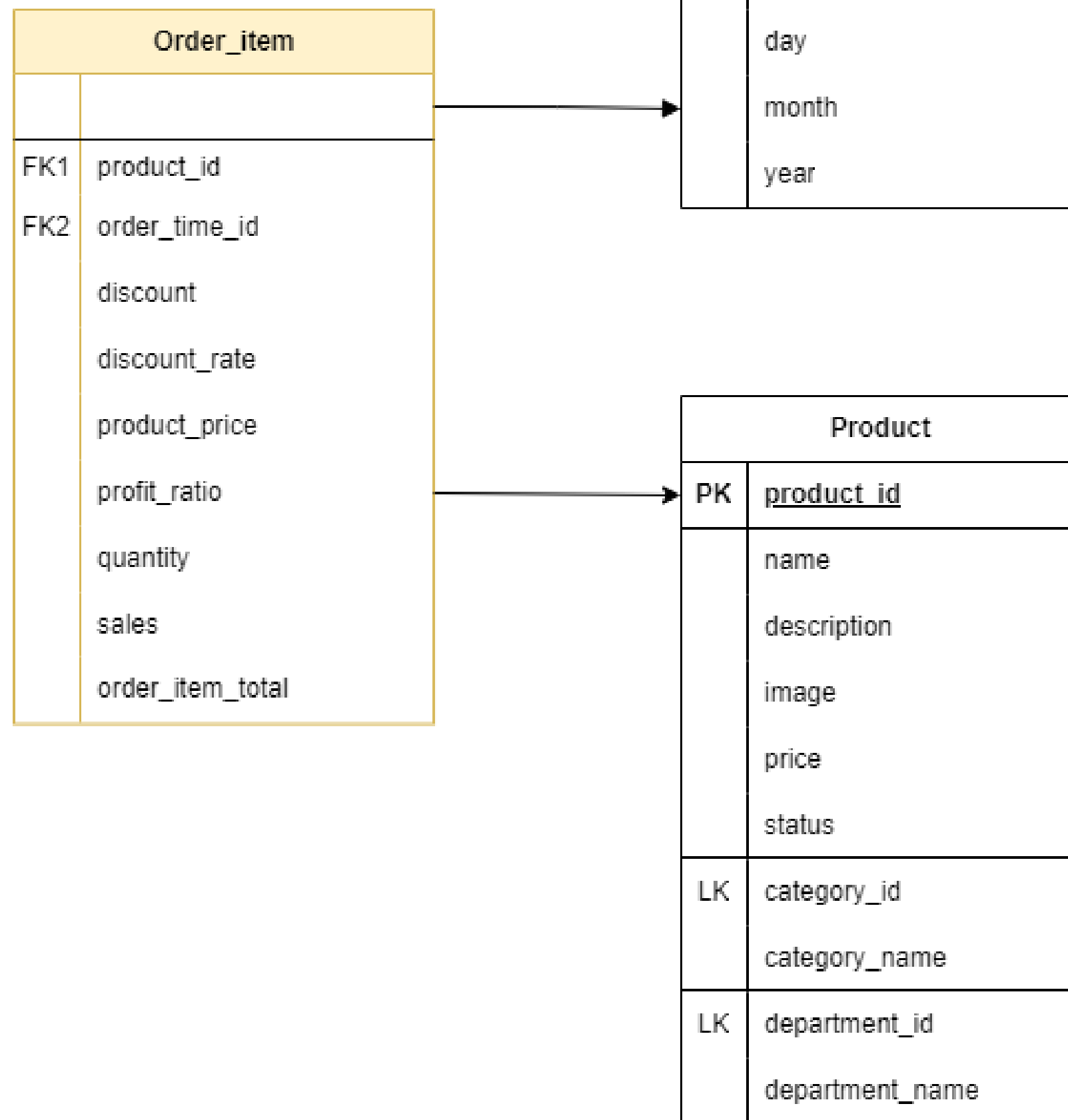| | | Dimensions | | | | |
|---|---|---|---|---|---|---|
| Stars (fact tables) | Granularity | Time | Date | Location | Product | Customer |
| Order | 1 / customer / date | x | | x | | x |
| Order items | 1 / product | x | | | x | |
| Sales (aggregation) | 1 / month | | x | x | X | x |

Dimensional Model

9

# Order star

| Order | |
|---|---|
| FK1 | customer_id |
| FK2 | customer_location_id |
| FK3 | order_location_id |
| FK4 | shipping_date |
| FK5 | order_date |
| | payment_type |
| | days_shipping_real |
| | days_shipping_scheduled |
| | delivery_status |
| | late_delivery_risk |
| | profit_per_order |
| | order_status |
| | shipping_mode |
| | number_of_items |
| | total_discount |
| | total_price |

| Customer | |
|---|---|
| PK | customer_id |
| | first_name |
| | last_name |
| | country |
| | city |
| | zip_code |
| | latitude |
| | longitude |
| LK | segment_id |
| | segment |

| Location | |
|---|---|
| PK | location_id |
| LK | city_id |
| | city |
| LK | state_id |
| | state |
| LK | country_id |
| | country |
| LK | region_id |
| | region |
| LK | market_id |
| | market |

| Time | |
|---|---|
| PK | time_id |
| | minute |
| | hour |
| | day |
| | month |
| | year |

10

Order item star

**Customer**

| PK | customer_id |
|----|----|
| | first_name |
| | last_name |
| | country |
| | city |
| | zip_code |
| | latitude |
| | longitude |
| LK | segment_id |
| | segment |

**Location**

| PK | location_id |
|----|----|
| LK | city_id |
| | city |
| LK | state_id |
| | state |
| LK | country_id |
| | country |
| LK | region_id |
| | region |
| LK | market_id |
| | market |

**Date**

| PK | date_id |
|----|----|
| LK | week_day_id |
| | week_day |
| LK | month_id |
| | month |
| LK | quarter_id |
| | quarter |
| LK | semester_id |
| | semester |
| LK | year |

**Agg_Sales**

| | |
|----|----|
| FK1 | date_id |
| FK2 | order_location_id |
| FK3 | customer_location_id |
| FK4 | customer_id |
| FK5 | product_id |
| | Sales_per_customer |
| | monthly_sales |
| | weekly_sales |

**Product**

| PK | product_id |
|----|----|
| | name |
| | description |
| | image |
| | price |
| | status |
| LK | category_id |
| | category_name |
| LK | department_id |
| | department_name |

# Sales Aggregation

12

# DM Implementation + ETL

Implement the dimensional model in an appropriate database system

# Implementation

- Communicate with database
- Creating dimension and fact tables
- ETL process
  - Extract
  - Transform (calculate fields and process data)
  - Load

# Input file

**Manual correction of the data fields**



15

# Load

**Slowly changing dimensions - allows easy future updates**

Dimensions made using data from CSV file:

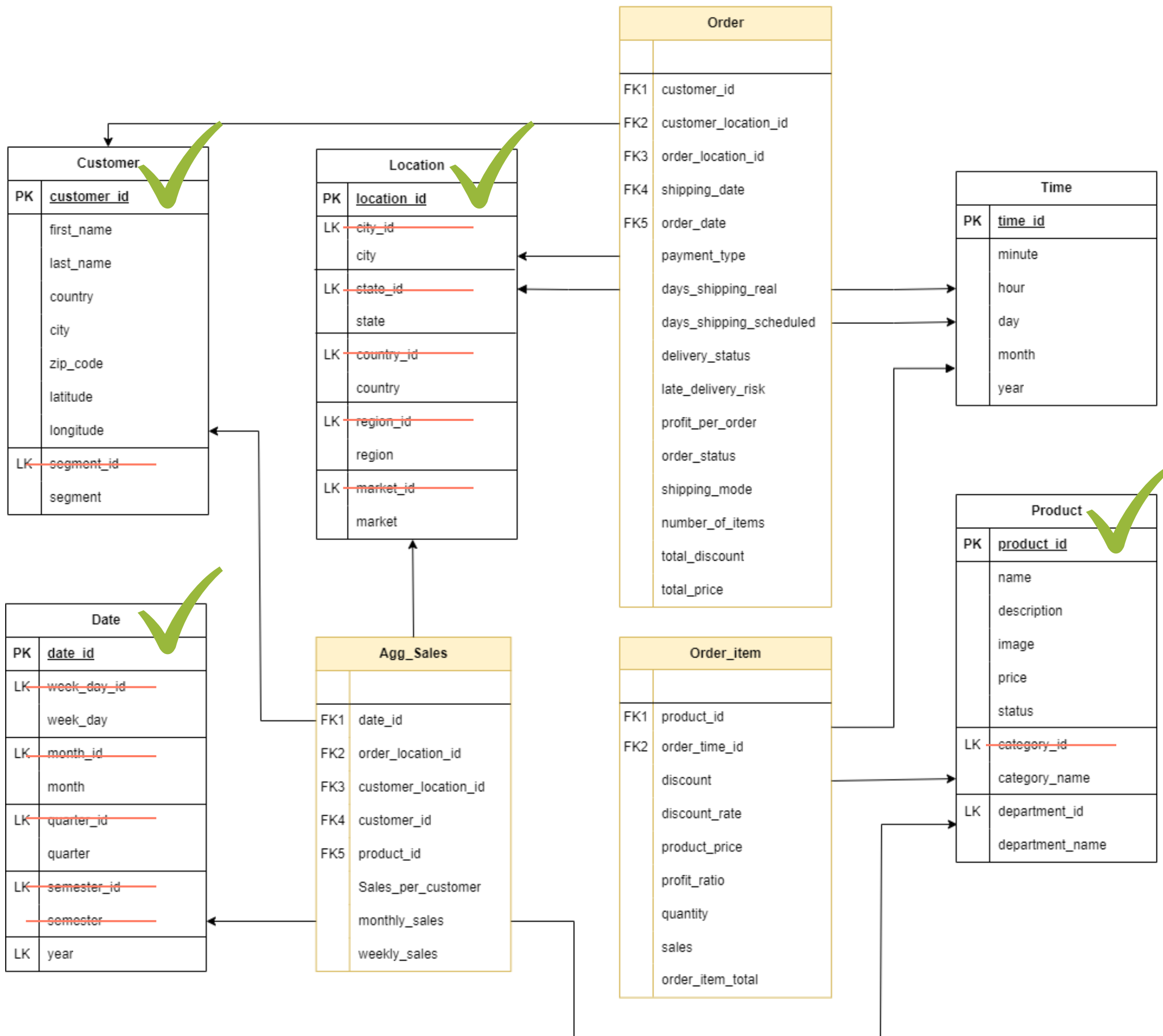- Customer_dim
- Product_dim
- Location_dim



DataCo CSV File Input       Dimension lookup/update
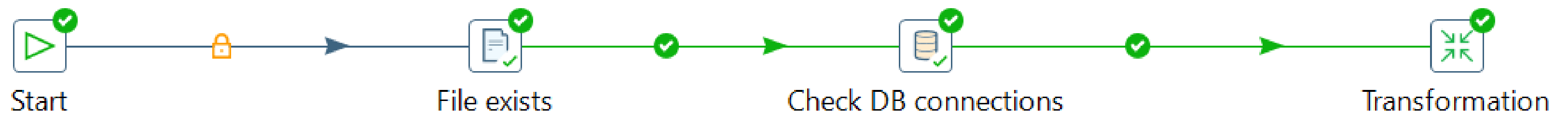
# Load

Dimensions rows generated in Pentaho:

- Date_dim

Done!

18

# Job scheduling

Schedule a set of ETL jobs to perform incremental data loads as new data is added to the operational (OLTP) systems and/or as operational data is changed.



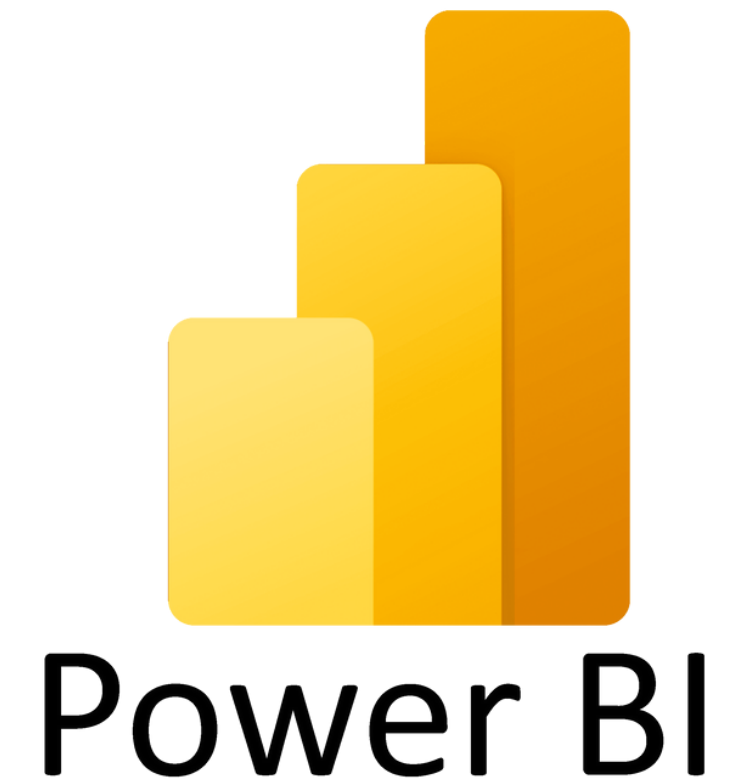| Start | File exists | Check DB connections | Transformation |

# Querying & Data Analysis

Develop business analytics reports, dashboards, or other user interfaces for the data warehouse.

# Querying



- Cubes
- Rollups

# Data Analytics



- Data Insights
- Visualization
- Transform & Clean data

# Conclusions

# Conclusions

**Datawarehouse VS Operational system**

- Scalability

- Data Integration

- Performance: complex analytical queries

- Historical Analysis

# Future works

- Finish the ETL process
- Querying and Analytics
- Use data to feed models and make forecasting systems to predict supply chain management

# DataCo Supply Chain Data Warehousing

**FEUP - MECD - Data Warehouse**
**Middle presentation**

Carlos Miguel Veloso
Cátia Teixeira
Luís Henriques
Rojan Aslani

Questions?

# LEVELS