# KHO VÀ KHAI PHÁ DỮ LIỆU

## (DATA WAREHOUSE AND DATA MINING)

GV: Phan Đình Vấn

# DATA WAREHOUSE

- Định nghĩa kho dữ liệu
- Mục đích và ý nghĩa của kho dữ liệu
- Đặc tính của dw
- Demo ETL

# 1.1. ĐỊNH NGHĨA KHO DỮ LIỆU (DATA WAREHOUSE - DW)

■ Đầu 1990s, Bill Inmon đã đặt ra thuật ngữ kho dữ liệu:

➤ A data warehouse is a collection of

- Subject-oriented,

- Integrated,

- Time-variant,
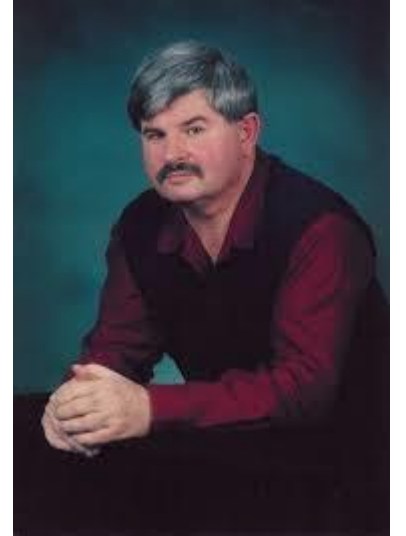
- Non-volatile

data in support of management decisions.



■ Ralph Kimball định nghĩa kho dữ liệu:

➤ A DW is a copy of transaction data specifically structured for query and analysis.

# 1.1. ĐỊNH NGHĨA DW

**Bill Inmon**, the Father of Data Warehousing

➤ William H. Inmon (born 1945) is an American computer scientist

➤ Bill Inmon first began to discuss the principles around the Data Warehouse.

➤ In 1992, Inmon published **Building the Data Warehouse**.

➤ In 2007, Inmon was named by Computerworld as one of the "Ten IT People Who Mattered in the Last 40 Years."

➤ Inmon's approach to Data Warehouse design focuses on a centralized data repository modeled to the third normal form. (Top-Down)

# 1.1. ĐỊNH NGHĨA DW

**Ralph Kimball**, other Father of Data Warehousing

➤ 1996, Ralph Kimball published **The Data Warehouse Toolkit**

➤ Kimball's early career in IT in the 1970s was highlighted by work as a key designer for the Xerox Star Workstation

➤ In the 1980s, he work with decision support systems as a Vice President for Metaphor Computer Systems.

➤ 1986, founded Red Brick Systems company with a full-fledged Data Warehouse application served as a major product.

➤ 1992, left Red Brick and start his own consultancy, Ralph Kimball Associates

➤ His well-regarded series of Data Warehouse Toolkit books soon followed.

- Web-based Data Warehousing
- ETL in a Data Warehousing environment,
- Microsoft-specific editions that cover SQL Server and the Microsoft Business Intelligence Toolset.

# 1.1. ĐỊNH NGHĨA DW

■ DW là tập các phương pháp, kỹ thuật và công cụ có thể kết hợp lại để cung cấp thông tin cho người dùng dựa trên việc tích hợp dữ liệu từ nhiều nguồn, nhiều môi trường khác nhau. (John Ladley)

■ Data Warehouse Technology is a set of methods, techniques and tools that can combine and support each other to provide information to the user based on integration data from different multiple sources/Environments. (John Ladley)

# 1.2. MỤC ĐÍCH VÀ Ý NGHĨA CỦA KHO DỮ LIỆU

- Sự bùng nổ của dữ liệu

- Sự phức tạp của dữ liệu

- Sự phân cấp và phân tán của hệ thống



- Vấn đề quản lý dữ liệu lịch sử

- Vấn đề lập báo cáo

- Vấn đề về quản lý và chia sẻ dữ liệu



- Hỗ trợ ra quyết định

# 1.2. MỤC ĐÍCH VÀ Ý NGHĨA CỦA KHO DỮ LIỆU

■ Kho dữ liệu là một cơ sở dữ liệu được thiết kế để hỗ trợ các hoạt động kinh doanh thông minh (BI)

➤ Giúp người dùng hiểu và nâng cao hiệu quả của tổ chức của họ.

➤ Nó được thiết kế để truy vấn và phân tích hơn là để xử lý giao dịch

➤ Thường chứa dữ liệu lịch sử được lấy từ dữ liệu giao dịch

➤ Kho dữ liệu tách biệt việc phân tích và các công việc giao dịch

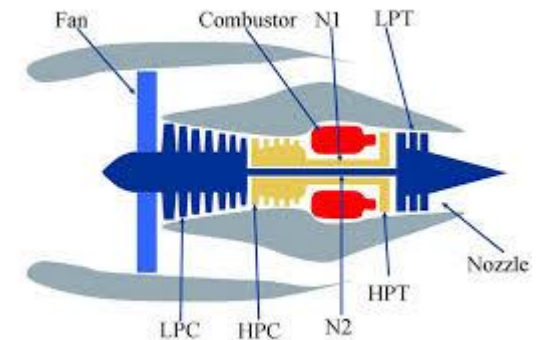➤ Cho phép hợp nhất dữ liệu từ nhiều nguồn.

# 1.2. MỤC ĐÍCH VÀ Ý NGHĨA CỦA KHO DỮ LIỆU

■ DW cũng là cơ sở dữ liệu quan hệ, môi trường kho dữ liệu có thể bao gồm:

➤ Extraction, transportation, transformation, and loading (ETL)

➤ Các công cụ phân tích thống kê, báo cáo, khai thác dữ liệu

➤ Các ứng dụng khác quản lý quá trình thu thập dữ liệu, chuyển thành thông tin hữu ích, cung cấp nó cho người dùng.

# 1.2. MỤC ĐÍCH VÀ Ý NGHĨA CỦA KHO DỮ LIỆU

■ Để đạt được mục tiêu nâng cao kinh doanh thông minh (trí tuệ kinh doanh), kho dữ liệu cần được thu thập từ nhiều nguồn.
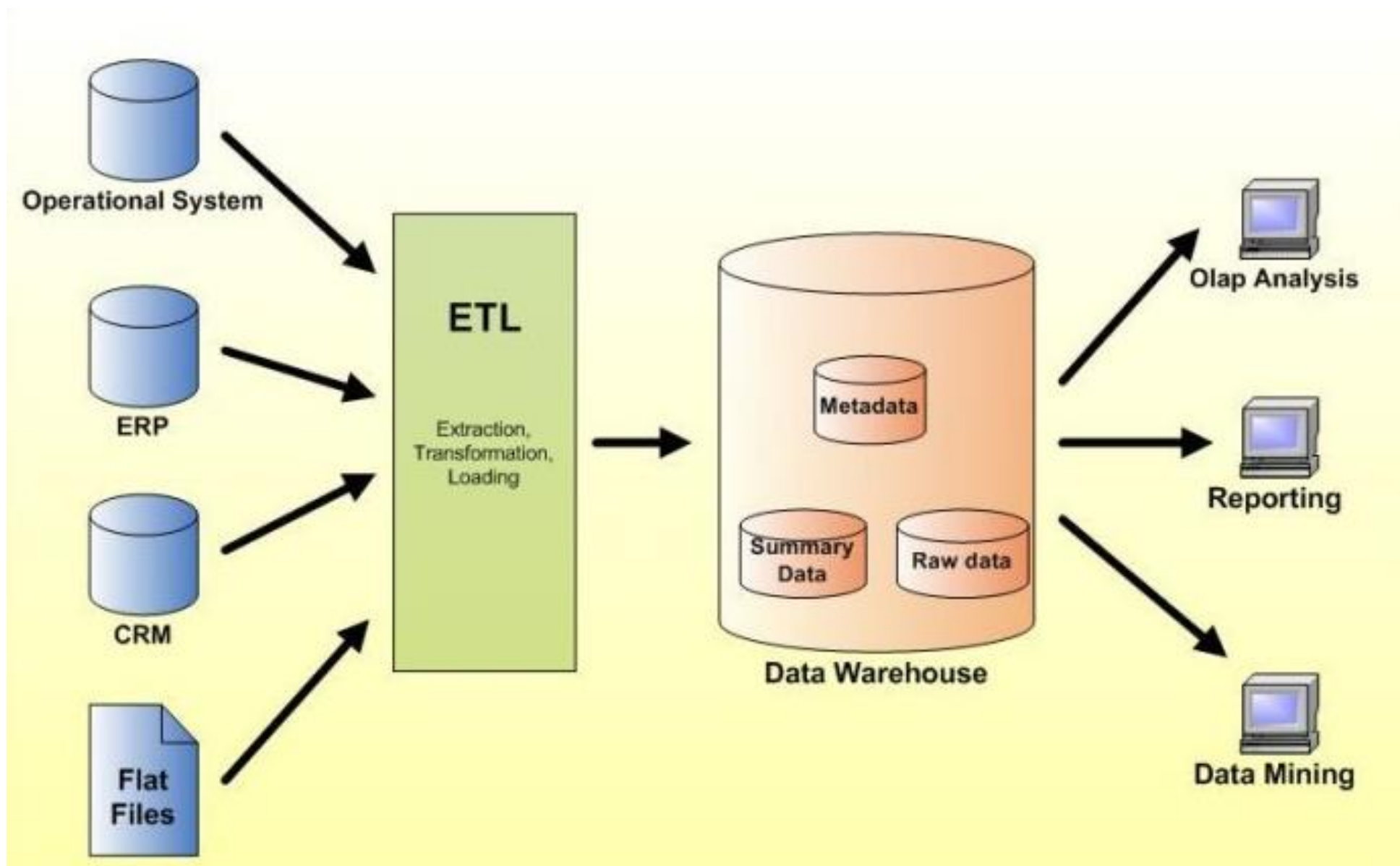
- ➤ Từ nội bộ
- ➤ Từ các phần mềm (đã mua)
- ➤ Từ bên thứ ba của các công ty cung cấp dữ liệu
- ➤ Các nguồn khác.
- ➤ Dữ liệu có thể liên quan đến giao dịch, sản xuất, tiếp thị, nguồn nhân lực…
- ➤ Ngày nay, dữ liệu có thể là từ các mạng xã hội, web... (click, like…)
- ➤ Dữ liệu từ các cảm biến (sensor) được tích hợp trong máy móc phức tạp.

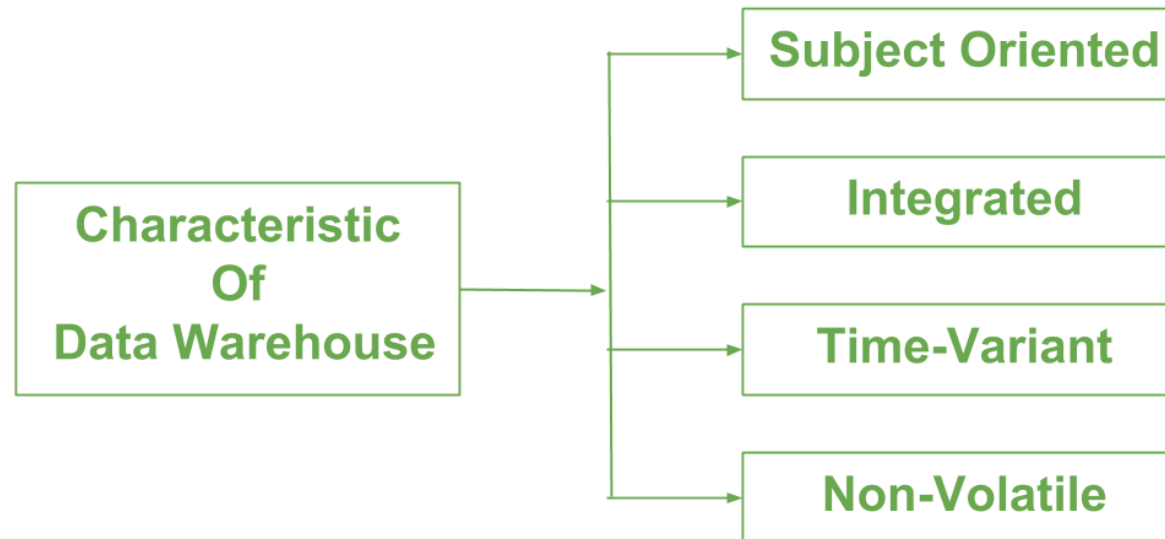# 1.2. MỤC ĐÍCH VÀ Ý NGHĨA CỦA KHO DỮ LIỆU

- **Kho dữ liệu khác biệt với hệ thống xử lý giao dịch trực tuyến (OLTP).**
  - Kho dữ liệu: tách biệt công việc phân tích với công việc giao dịch.
  - DW định hướng đọc hơn là ghi và chỉnh sửa.
    - Tăng hiệu suất phân tích
    - Tránh ảnh hưởng đến hệ thống giao dịch.
  - DW có thể được tối ưu hóa để hợp nhất dữ liệu từ nhiều nguồn và trở thành nguồn duy nhất của tổ chức.
  - Nguồn dữ liệu nhất quán cho tất cả người dùng
  - Ngăn ngừa được các tương tranh dữ liệu, nâng cao hiệu quả khi ra quyết định.

# 1.2. MỤC ĐÍCH VÀ Ý NGHĨA CỦA KHO DỮ LIỆU

# 1.3. ĐẶC TÍNH CỦA DW

- Hướng chủ đề (Subject-Oriented)

- Tích hợp (Integrated)

- Dữ liệu gắn thời gian và có tính lịch sử (Time-Variant)

- Dữ liệu không biến động (Non-volatile)

# 1.3. ĐẶC TÍNH CỦA DW

■ Hướng chủ đề (Subject-Oriented)

➤ Dữ liệu trong DW được xác định từ đầu là để phân tích về một hoặc một số chủ đề nhất định, ví dụ "doanh thu".

➤ DW không phải là nơi lưu trữ thông tin về mọi mặt hoạt động của công ty hoặc tổ chức.

■ Tích hợp (Integrated)

➤ Dữ liệu được tập hợp từ nhiều nguồn và lưu trữ nhất quán. Ví dụ, cùng một mặt hàng được quản lý với 2 tên khác nhau ở 2 hệ thống quản lý kho và hệ thống bán hàng.

➤ Khi tập hợp dữ liệu về DW, sẽ có một bước biến đổi (transform) để đưa về một tên duy nhất cho mặt hàng này.
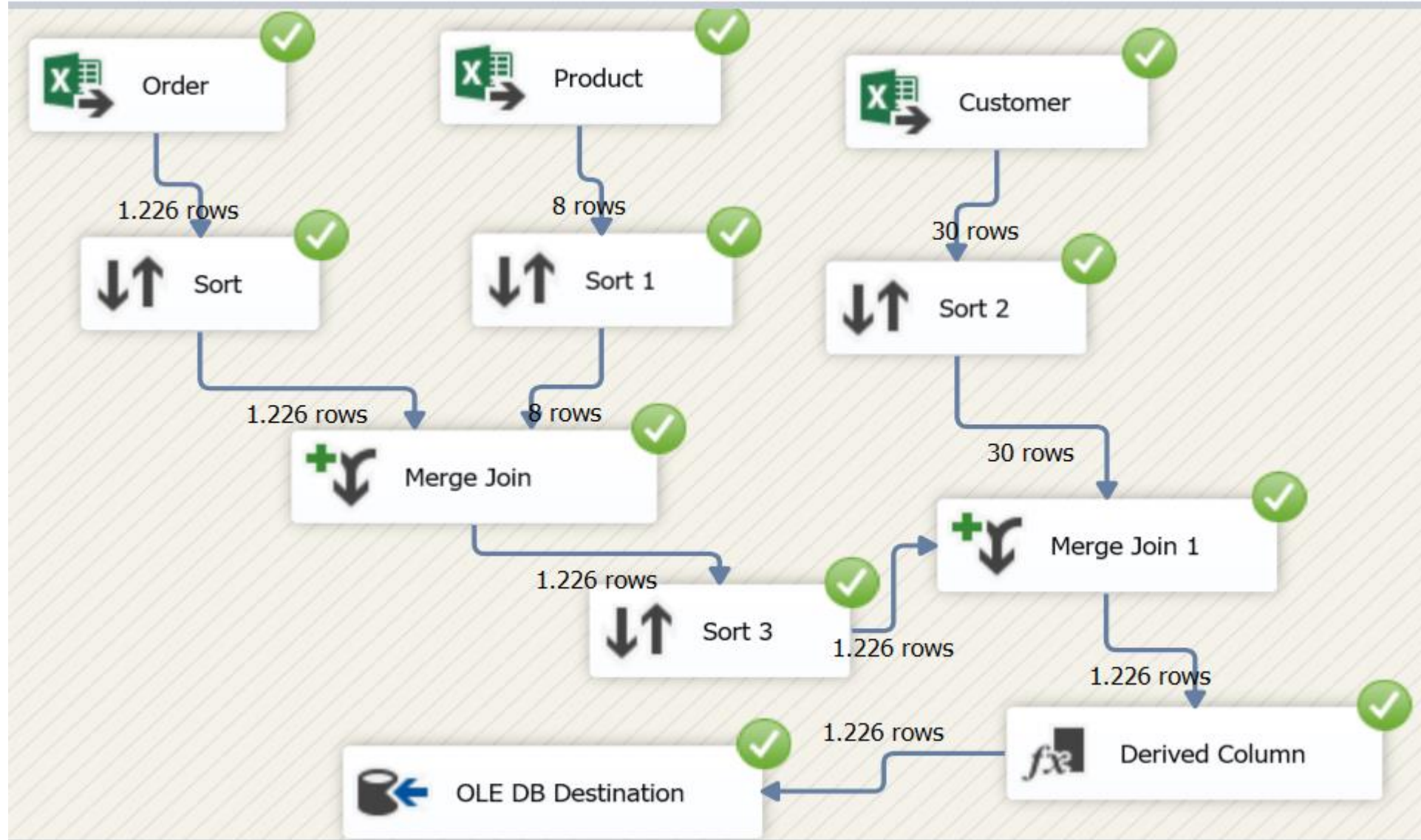
# 1.3. ĐẶC TÍNH CỦA DW

■ Biến đổi theo thời gian (Time-variant)

➤ Dữ liệu trong Data Warehouse luôn gắn với một thời điểm cụ thể trong một giới hạn thời gian nhất định. Ví dụ, người sử dụng dữ liệu có thể truy vấn lịch sử hàng tồn kho trong 3, 6 hoặc 12 tháng trước.

➤ Đây cũng là một trong những khác biệt căn bản của DW với các hệ thống OLTP, vốn chỉ lưu trữ trạng thái dữ liệu mới nhất (trong ví dụ trên là lượng hàng tồn kho hiện tại).

■ Ổn định (Non-volatile)

➤ Một khi đã được đưa vào Data Warehouse, dữ liệu sẽ không bị thay đổi hoặc xóa.

➤ Đặc điểm này cho phép người quản lý có được bức tranh toàn cảnh về toàn bộ lịch sử hoạt động.
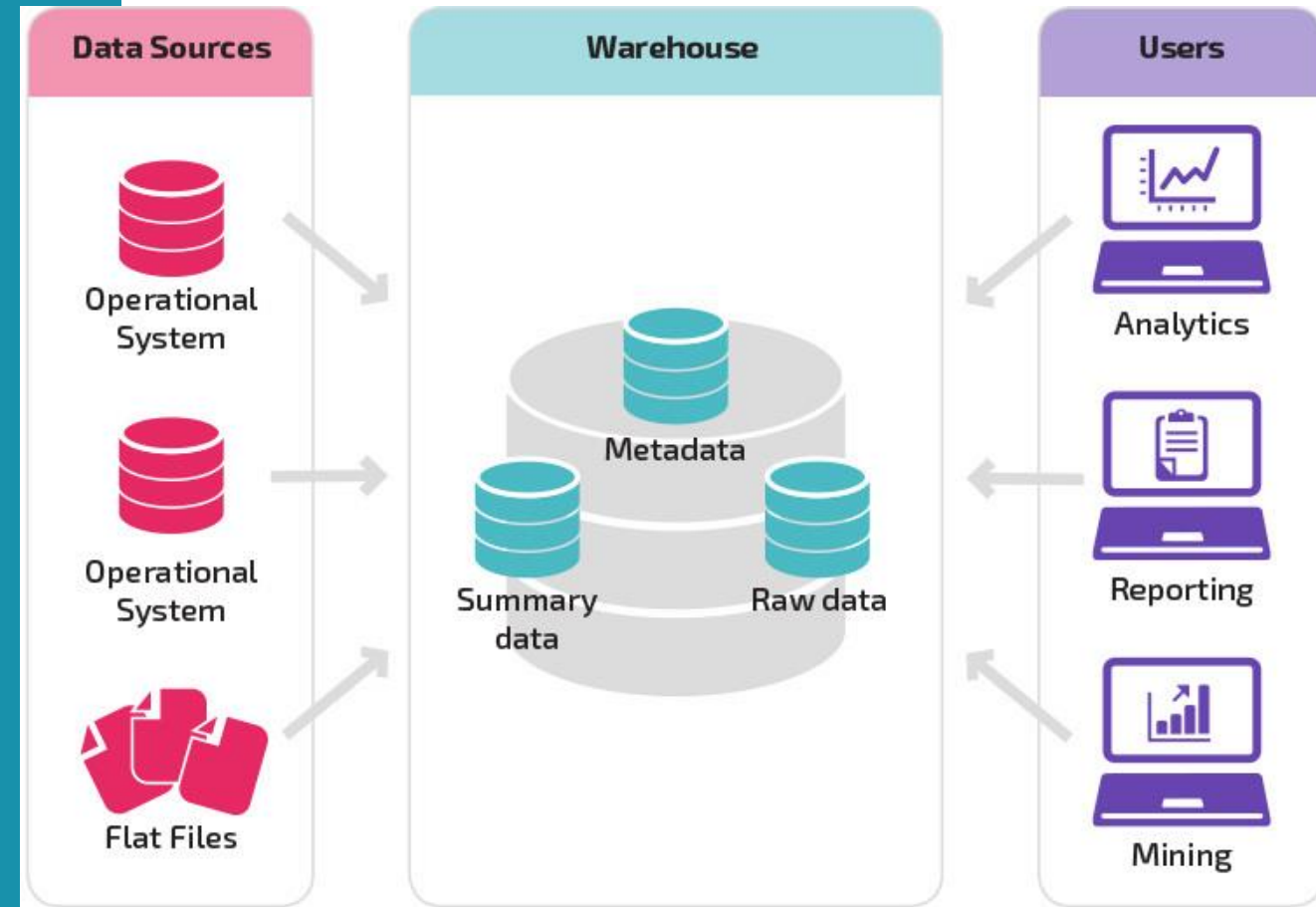
# DEMO ETL

# Questions

1. What type of data are there in the world?
2. What sources can create data?
3. How can connect, process and analyze these types of data effectively?
4. How can we store these types of data effectively?
5. What solutions and tools to support for storing and analyzing those data?
6. What is a data warehouse?
7. What are characteristics of a data warehouse?
8. What are the purposes of data warehouse?
9. Who is the father of the data warehouse?
10. How is the history of the data warehouse?

# DATA WAREHOUSE

- Database vs. Data warehouse
- Operation database vs. Data warehouse
- Data warehouse architectures
- Data warehouse building process

# 1.4. DATABASE VS. DATA WAREHOUSE

| Database | Data Warehouse |
|---|---|
| Transaction Oriented | Subject Oriented |
| Detail Data | Summarized Data |
| For OLTP | For OLAP |
| Optimized for write operation | Optimized for read operation |
| Low performance for analytical queries | High performance for analytical queries |
| Current/ Real - time | Historical |
| Size data: MB - GB | Size data: GB-TB |
| Purpose for data retrieval, Updating and management | Purpose for data analysis and decision making |

# 1.4. DATABASE VS. DATA WAREHOUSE

| Database | Data Warehouse |
|---|---|
| Is designed to record | Is designed to analyze |
| Tables and joins of a database are complex as they are normalized. | Table and joins are simple in a data warehouse because they are denormalized. |
| Generally limited to a single application | Stores data from any number of applications |
| ER model | Star schema, snowflake |
| Capture data | Analyze data |

# 1.4. OPERATION DATABASE VS. DATA WAREHOUSE

**Operation database**

- Cơ sở dữ liệu tác nghiệp là nguồn cho kho dữ liệu.

- Bao gồm thông tin chi tiết được sử dụng để điều hành hoạt động hàng ngày của doanh nghiệp.

- Dữ liệu thường xuyên thay đổi.

- Quản lý dữ liệu động trong thời gian thực

- Cơ sở dữ liệu hoạt động được gọi là OLTP (xử lý giao dịch trực tuyến).

**Data warehouse**

- DW mục đích để phân tích dữ liệu và ra quyết định, còn được gọi là hệ thống Xử lý phân tích trực tuyến (OLAP).

- OLAP và OLTP đều là cơ sở dữ liệu quan hệ, nhưng mục tiêu khác nhau.

# 1.5. Data Warehouse Architectures

■ **Back-end tier**
  ➤ ETL process
  ➤ Data Staging Area (DSA): Data was integrated and transformed then load to DW
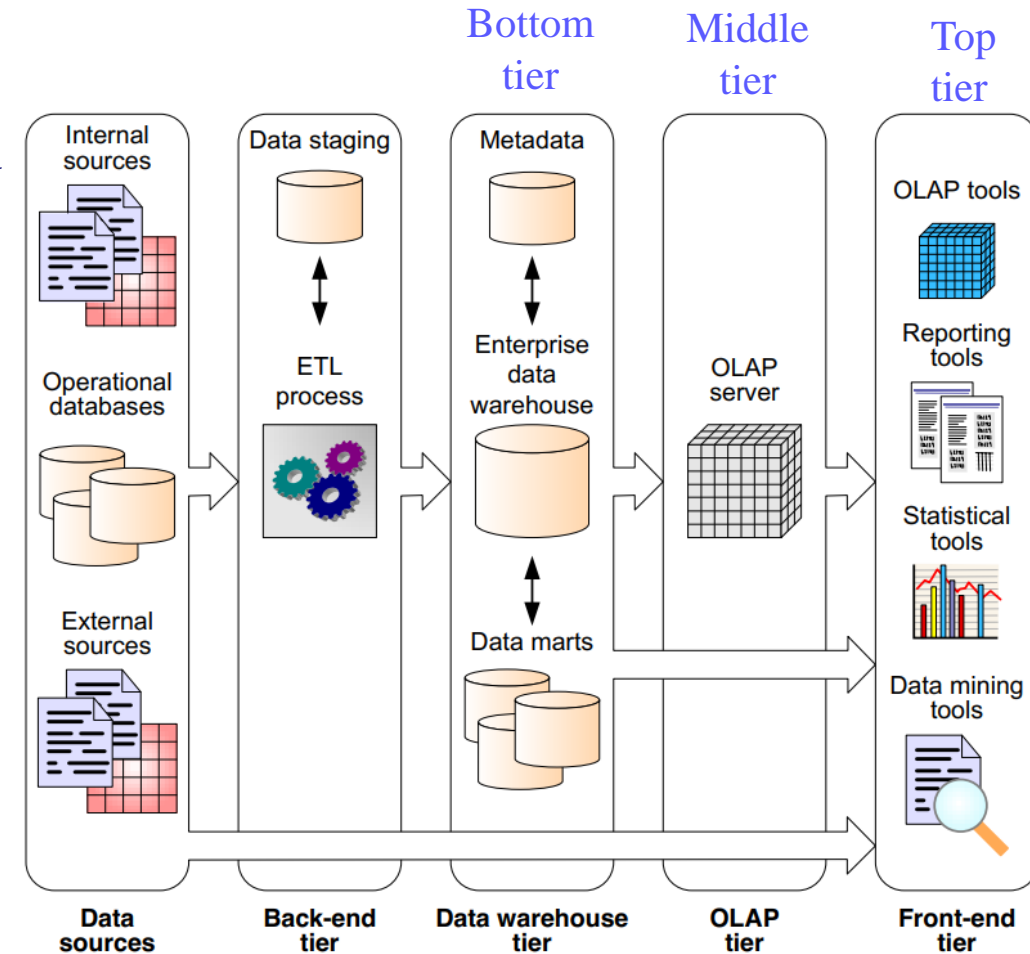
■ **Data warehouse tier**
  ➤ Enterprise DW and/or several data marts
  ➤ Metadata

■ **OLAP tier**
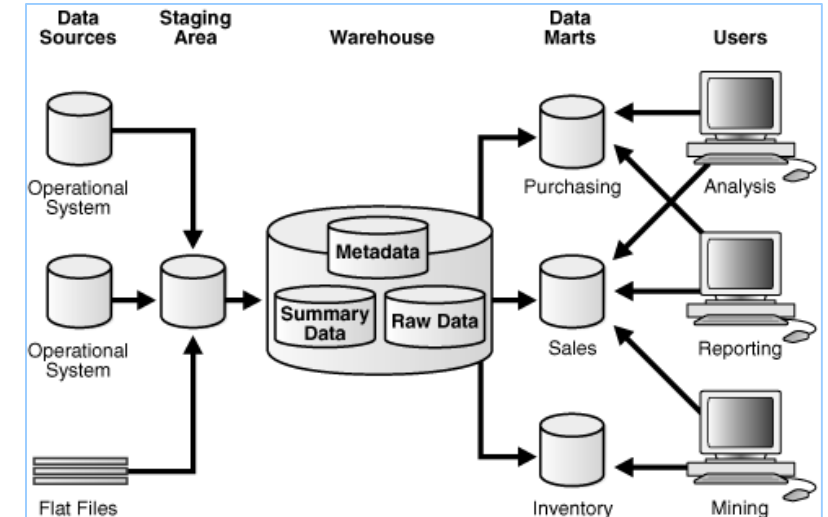  ➤ Provides a multidimensional view
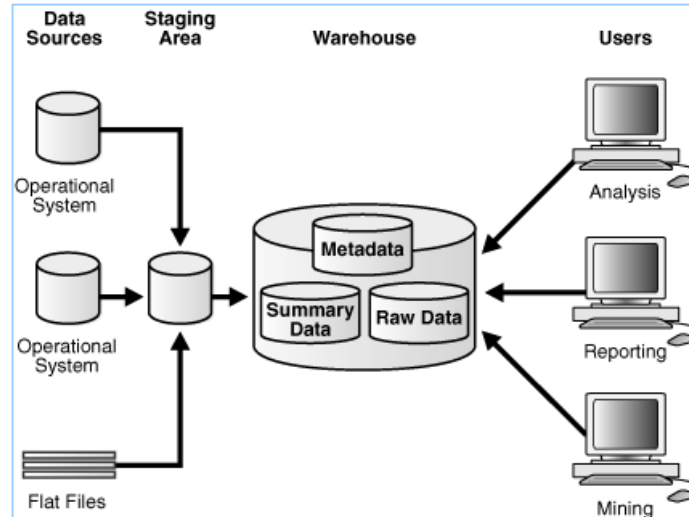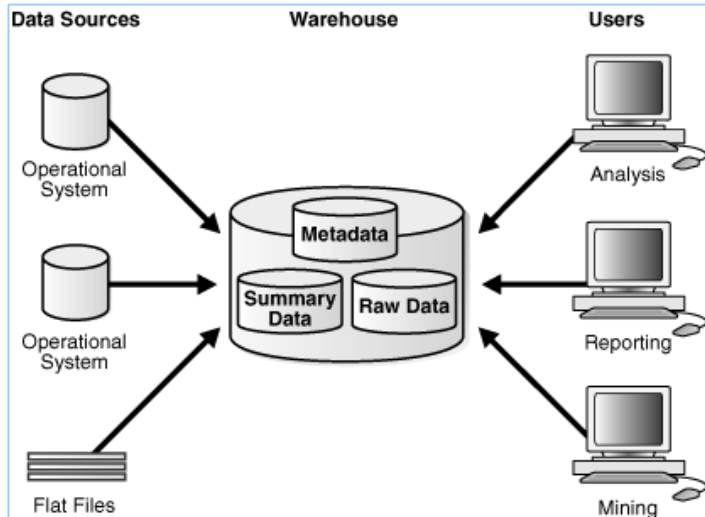
■ **Front-end tier**
  ➤ Used for data analysis and visualization
  ➤ Contains client tools: OLAP tools, reporting tools, statistical tools, and data-mining tools

# 1.5. Data Warehouse Architectures

- Three common architectures
  - Data Warehouse Architecture: Basic
  - Data Warehouse Architecture: with a Staging Area
  - Data Warehouse Architecture: with a Staging Area and Data Marts

# 1.5. Data Warehouse Architectures

**Bottom Up Vs Top Down Approach in Data Warehouse**

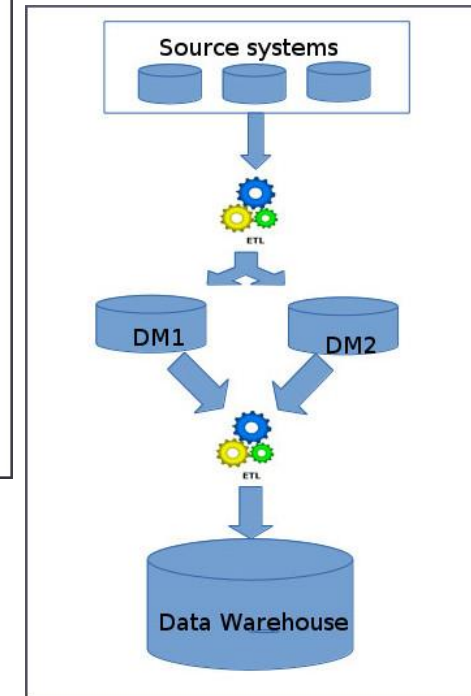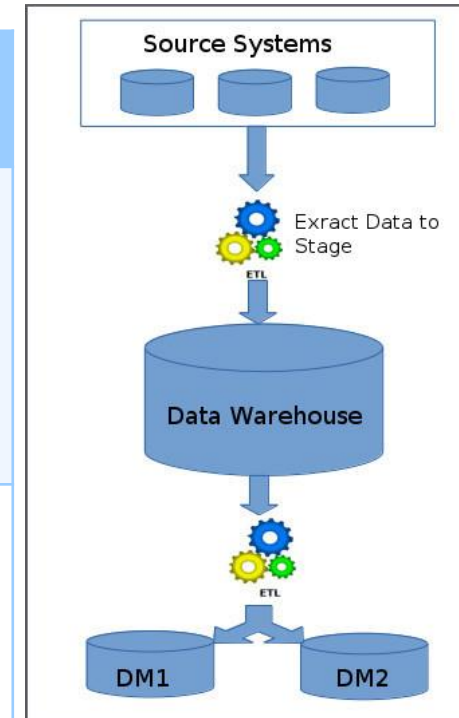| Top-Down Approach | Bottom-Up Approach |
|---|---|
| Provides a definite and consistent view of information as information from the data warehouse is used to create Data Marts | Reports can be generated easily as Data marts are created first and it is relatively easy to interact with data marts. |
| Strong model and hence preferred by big companies | Not as strong but data warehouse can be extended and the number of data marts can be created |
| Time, Cost and Maintenance is high | Time, Cost and Maintenance are low. |
| Bill Inmon | Kimball |

# 1.5. Data Warehouse Architectures

## Bottom Up Vs Top Down Approach in Data Warehouse

| Top-Down | Bottom-Up |
|---|---|
| **Advantages** | |
| - It is easier to maintain Top Down Design<br>- Provides consistent dimensional views of data across data marts, as all data marts are loaded from the DW.<br>- This approach is robust against business changes.<br>- Creating a new data mart from the data warehouse is very easy.<br>- Initial cost is high but subsequent project development cost is lower | - This model contains consistent data marts and these data marts can be delivered quickly.<br>- The data marts are created first to provide reporting capability<br>- It is easier to extend. Creating new data marts and then integrating with others.<br>- This Approach take less time. Initial set up is very quickly |
| **Disadvantage** | |
| - It represents a very large project and the cost of implementing the project is significant.<br>- It is time consuming and more time required for initial set up<br>- Highly skilled people required for set up | - Initial cost is low but each subsequent phase will cost same<br>- The positions of the DW and the data marts are reversed in the bottom-up approach design.<br>- It is difficult to maintain and often redundant and subject to revisions |

# 1.5. Data Warehouse Architectures

## Data marts

➤ Data Mart giúp tăng cường thời gian phản hồi của người dùng do giảm khối lượng dữ liệu.

➤ Chi phí triển khai Data Mart thấp hơn so với việc triển khai kho dữ liệu đầy đủ.

➤ Dependent Data Mart

- Data Mart phụ thuộc chứa những dữ liệu được lấy từ Data Warehouse

- Dữ liệu được trích lọc và tinh chế, tích hợp lại ở mức cao hơn để phụ vụ một chủ đề nhất định của Data Mart.

➤ Independent Data Mart

- Data Mart độc lập được xây dựng trước Data Warehouse

- Dữ liệu được lấy trực tiếp từ các nguồn khác nhau.

➤ Hybrid


Dependent Data Mart


Independent Data Mart


Hybrid Data Mart

# 1.6. Data warehouse building process

- Step 1: Determine Business Objectives and Identify Core Business Processes

- Step 2: Locate Data Sources

- Step 3: Design DW structure (Cube, Dimensions, Measures. Hierarchy…)

- Step 4: ETL (Extract, Transform, Load)

- Step 5: Implement the DW

- Step 6: Set Tracking Duration

# Questions

1. What is OLTP? What are the differences between Database and Data warehouse?
2. What are the differences between architecture approaches, the architectures of data warehouse?
3. How to build a data warehouse step by step?

## ETL (csv)

# DATA WAREHOUSE

- DW vs. Online Analytical Processing (OLAP)
- Star Schema
- Snowflake Schema
- Demo ETL to Star schema

# 1.7. DW & Online Analytical Processing (OLAP)

## Multidimensional Model

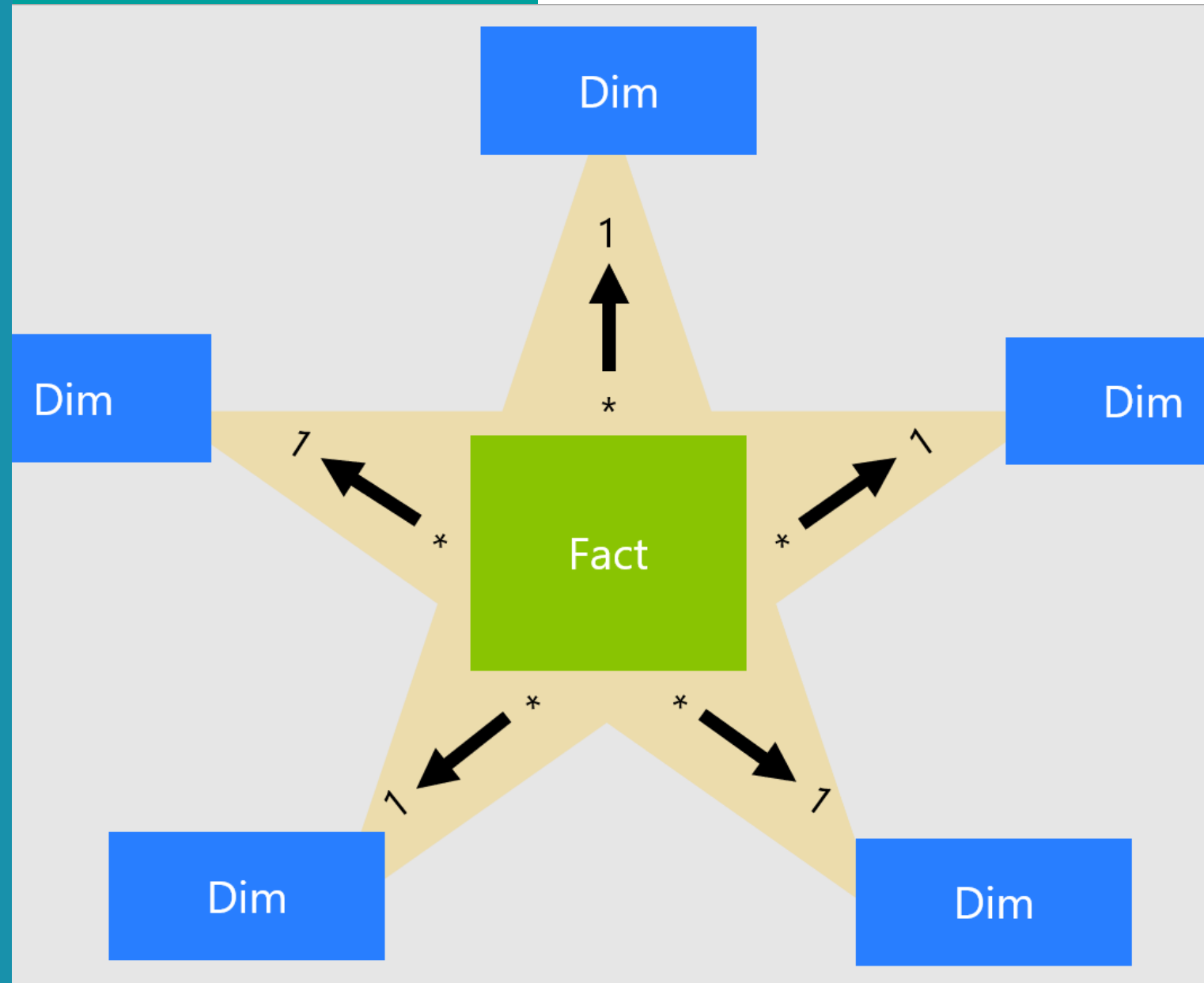➤ DWs and OLAP use a multidimensional view of data

➤ Represented as a data cube or an hypercube

- Dimensions: Perspectives for analyzing data

- Cells (facts): Contain measures, values that are to be analyzed

# 1.7. DW & Online Analytical Processing (OLAP)

**Hierarchies**

- Data granularity (độ mịn): Level of detail of measures
- Data analyzed at different granularities
- Hierarchies relate low-level (detailed) concepts to higher-level (general concepts)
    - Example: Store – City – Region/Province – Country
- Given two related levels in a hierarchy, lower level is called child, higher level is called parent
- Instances of these levels are called members

# 1.7. DW & Online Analytical Processing (OLAP)

**Product**

| All |
| --- |

| Category |
| --- |

| Product |
| --- |

**Time**

| All |
| --- |

| Year |
| --- |

| Semester |
| --- |

| Quarter |
| --- |

| Month |
| --- |

| Day |
| --- |

**Customer**

| All |
| --- |

| Continent |
| --- |

| Country |
| --- |

| State |
| --- |

| City |
| --- |

| Customer |
| --- |

**Hierarchies**

➤ Example

- Hierarchies of the Product,
- Time, and Customer dimension

# 1.7. DW & Online Analytical Processing (OLAP)

**Hierarchies**

Store dimension — All

Country level — France ... Italy

Region/Province level — Ile-de-France, Provence-Alpes-Côte d'Azur ... ..., Lazio, Lombardy

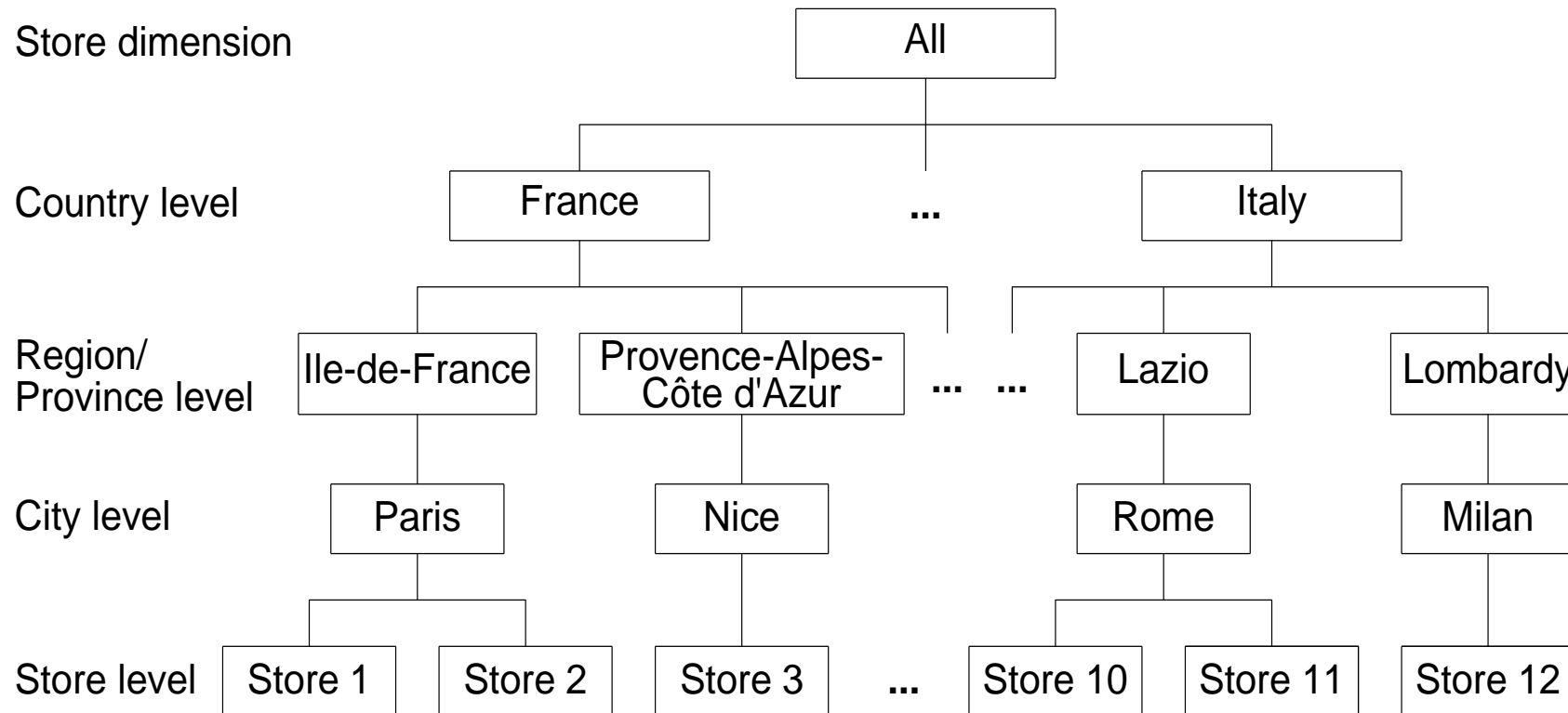City level — Paris, Nice, Rome, Milan

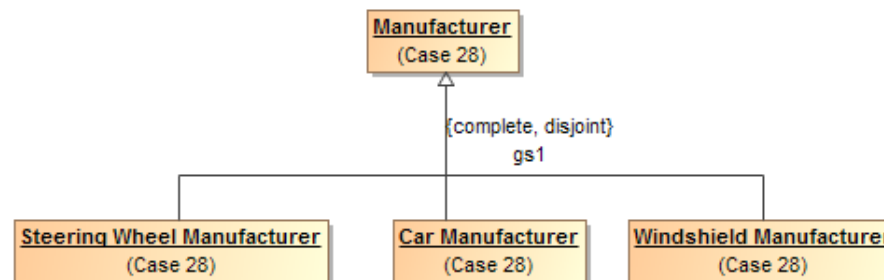Store level — Store 1, Store 2, Store 3 ... Store 10, Store 11, Store 12

# 1.7. DW & Online Analytical Processing (OLAP)

**Measure Aggregation and Summarizability**

➤ Measures are aggregated when using hierarchies for visualizing data at different abstraction levels

➤ Summarizability conditions ensure correct aggregation

■ Disjointness of instances: Grouping of instances in a level with respect to the parent in the next level must result in disjoint sets

■ Completeness: All instances are included in the hierarchy and each instance is related to one parent in the next level

■ Correct use of aggregation functions: Type of measures determine the kind of aggregation functions that can be applied

Manufacturer
(Case 28)

{complete, disjoint}
gs1

Steering Wheel Manufacturer
(Case 28)

Car Manufacturer
(Case 28)

Windshield Manufacturer
(Case 28)

# 1.7. DW & Online Analytical Processing (OLAP)

**Elements of Dimensional Data Model**

➤ A data structure technique optimized for data storage in a DW

➤ The purpose of dimensional modeling is to optimize the database for faster retrieval of data.

➤ The concept of Dimensional Modelling was developed by Ralph Kimball and consists of "fact" and "dimension" tables.

# 1.7. DW & Online Analytical Processing (OLAP)

**Elements of Dimensional Data Model**

➤ **Fact**

- Result from a business process or business event
    - Facts are usually numeric and additive
- Granularity/grain (độ mịn)
    - Identifies the fact level of detail
    - One row per sale, one row per service call, one row per claim, …
    - Atomic grain is most flexible

# 1.7. DW & Online Analytical Processing (OLAP)

**Elements of Dimensional Data Model**

- **Dimension**
    - Provides the context surrounding a business process event.
    - They give who, when, what, where of a fact.
        - Ex. Product, Customer, Date, Patient, Vendor,
    - Each dimension row is a unique occurrence

- **Dimension attributes**
    - The Attributes are the various characteristics of the dimension in dimensional data modeling.
    - Ex. Location dimension: State, Country, Zip code etc
    - Hierarchical relationships

# 1.7. DW & Online Analytical Processing (OLAP)

**Elements of Dimensional Data Model**

- **Dimension Table**
  - A dimension table contains dimensions of a fact.
  - They are joined to fact table via a foreign key.
  - Dimension tables are de-normalized tables.
  - Dimensions offers descriptive characteristics of the facts
  - No set limit set for given for number of dimensions
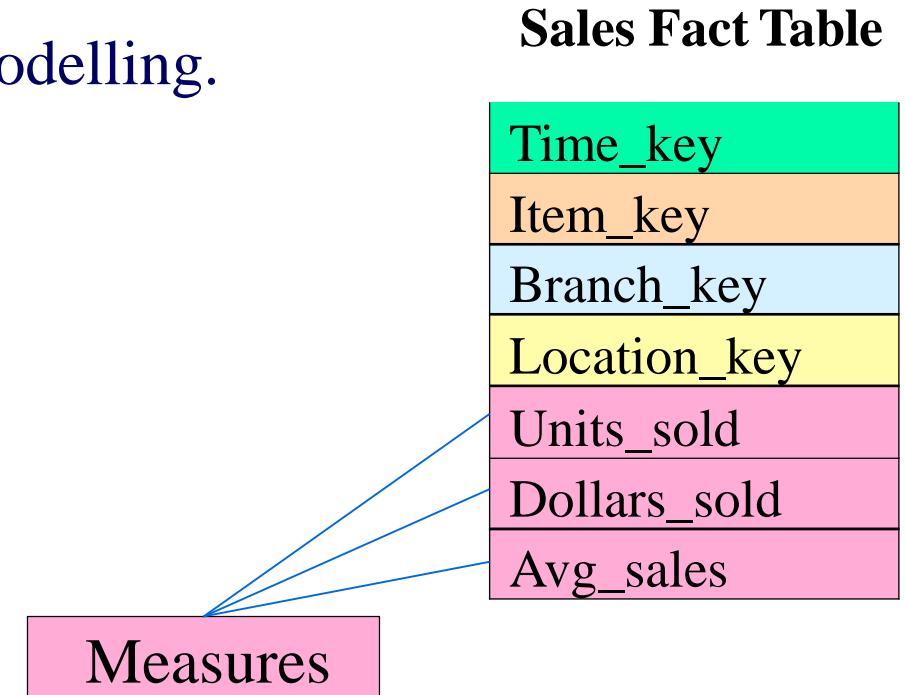  - The dimension can also contain hierarchical relationships

| Item |
|------|
| Item_key |
| Item_name |
| Brand |
| Type |
| Supplier_type |

# 1.7. DW & Online Analytical Processing (OLAP)

## Elements of Dimensional Data Model

### Fact Table

- A fact table is a primary table in dimension modelling.
- A Fact Table contains
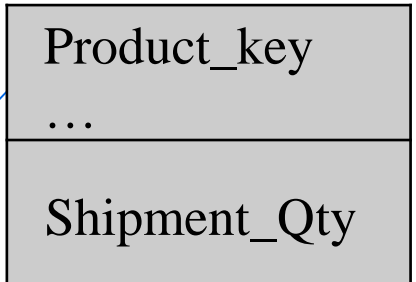  - Measurements/facts
  - Foreign key to dimension table

**Sales Fact Table**

| Time_key |
| --- |
| Item_key |
| Branch_key |
| Location_key |
| Units_sold |
| Dollars_sold |
| Avg_sales |

Measures
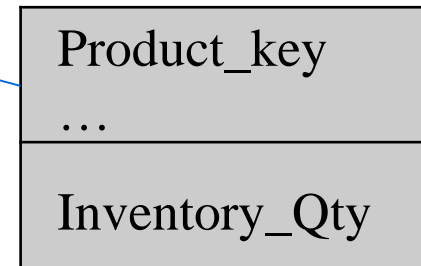
# 1.7. DW & Online Analytical Processing (OLAP)

- **Conformed Dimensions**

  ➤ Shared across business processes (fact tables) in the DW

  ➤ All fact tables use same standard dimensions

  ➤ Established via Bus Matrix, enforced in ETL

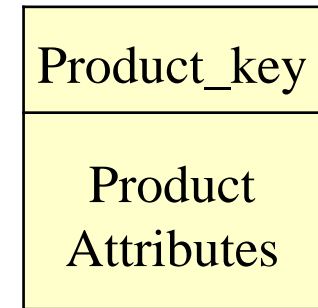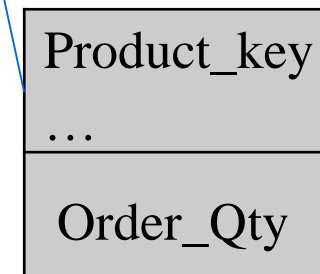**Shipment Facts**

| Product_key |
| --- |
| … |
| Shipment_Qty |

| Product_key |
| --- |
| Product Attributes |

Dimensions

**Inventory Facts**

| Product_key |
| --- |
| … |
| Inventory_Qty |

**Order Facts**

| Product_key |
| --- |
| … |
| Order_Qty |

**CONFORMED DIMENSIONS**

**ORDER**

Product Key
Date Key
Customer Key
Salesperson Key
Order Dollars
Cost Dollars
Margin Dollars
Sale Units

**PRODUCT**

Product Key
.....................

**CUSTOMER**

Customer Key
.....................

**SALESPERSON**

Salesperson Key
.....................

**DATE**

Date Key
.....................

**SHIPMENT**

Product Key
Date Key
Customer Key
Salesperson Key
Channel Key
Ship-to Key
Ship-from Key
Invoice Number
Order Number
Ship Date
Arrival Date

**CHANNEL**

Channel Key
.....................

**SHIP-TO**

Ship-to Key
.....................

**SHIP-FROM**

Ship-from Key
.....................
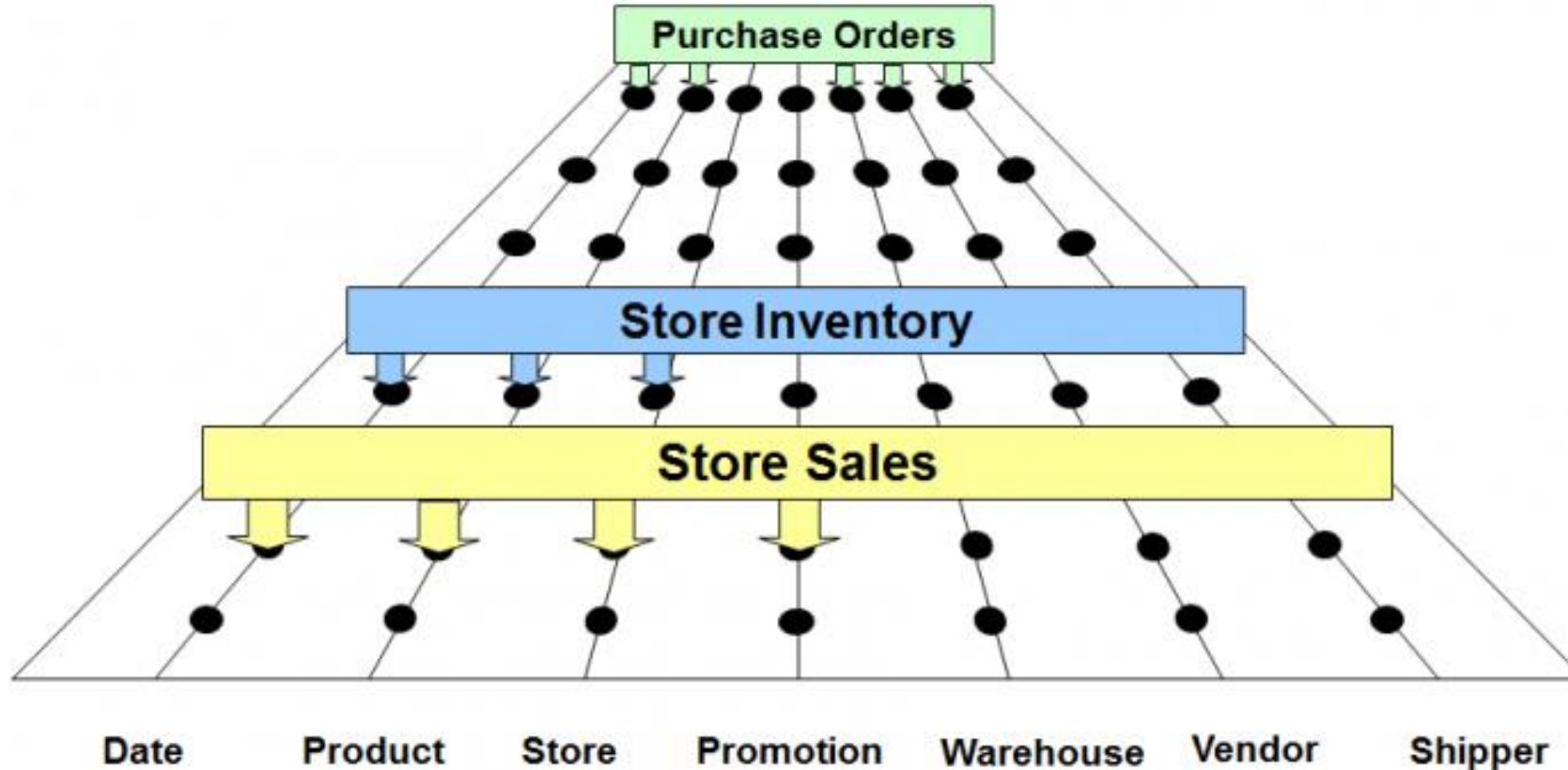
# 1.7. DW & Online Analytical Processing (OLAP)

**Enterprise DW bus Architecture**

# 1.7. DW & Online Analytical Processing (OLAP)

## DW Bus Matrix

➤ Rows = Business processes

➤ Columns = Conformed dimensions



| BUSINESS PROCESSES | COMMON DIMENSIONS | | | | | | |
|---|---|---|---|---|---|---|---|
| | Date | Product | Warehouse | Store | Promotion | Customer | Employee |
| Issue Purchase Orders | X | X | X | | | | |
| Receive Warehouse Deliveries | X | X | X | | | | X |
| Warehouse Inventory | X | X | X | | | | |
| Receive Store Deliveries | X | X | X | X | | | X |
| Store Inventory | X | X | | X | | | |
| Retail Sales | X | X | | X | X | X | X |
| Retail Sales Forecast | X | X | | X | | | |
| Retail Promotion Tracking | X | X | | X | X | | |
| Customer Returns | X | X | | X | X | X | X |
| Returns to Vendor | X | X | | X | | | X |
| Frequent Shopper Sign-Ups | X | | | X | | X | X |

# 1.7. DW & Online Analytical Processing (OLAP)

**Measure Classification: Additivity**

- **Additive** measures (**flow** or **rate** measures): Can be meaningfully summarized using addition along **all dimensions**
  - E.g., sales amount can be summarized when the hierarchies in Store, Time, and Product dimensions are traversed
- **Semi-additive** measures (**stock** or **level** measures): Can be meaningfully summarized using addition along **some (not all) dimensions**
  - E.g., inventory quantities, can be aggregated in the Store dimension, but cannot be aggregated in the Time dimension
- **Non-additive** measures (**value-per-unit** measures): Cannot be meaningfully summarized using addition along **any dimension**
  - E.g., item price, cost per unit, exchange rate

# 1.7. DW & Online Analytical Processing (OLAP)

**Measure Classification: Aggregation Complexity**

➤ Distributive (phân tán) measures: Defined by an aggregation function that can be computed in a distributed way

- Data is partitioned into n sets, aggregate function applied to each set, aggregated value is computed by applying a function to these n sub-aggregate values
- E.g., sum, min, max, count (distinct count is not)

➤ Algebraic (đại số) measures: Defined by an aggregation function that has can be expressed as a scalar function of distributive function

- E.g., average (can computed by sum/count)

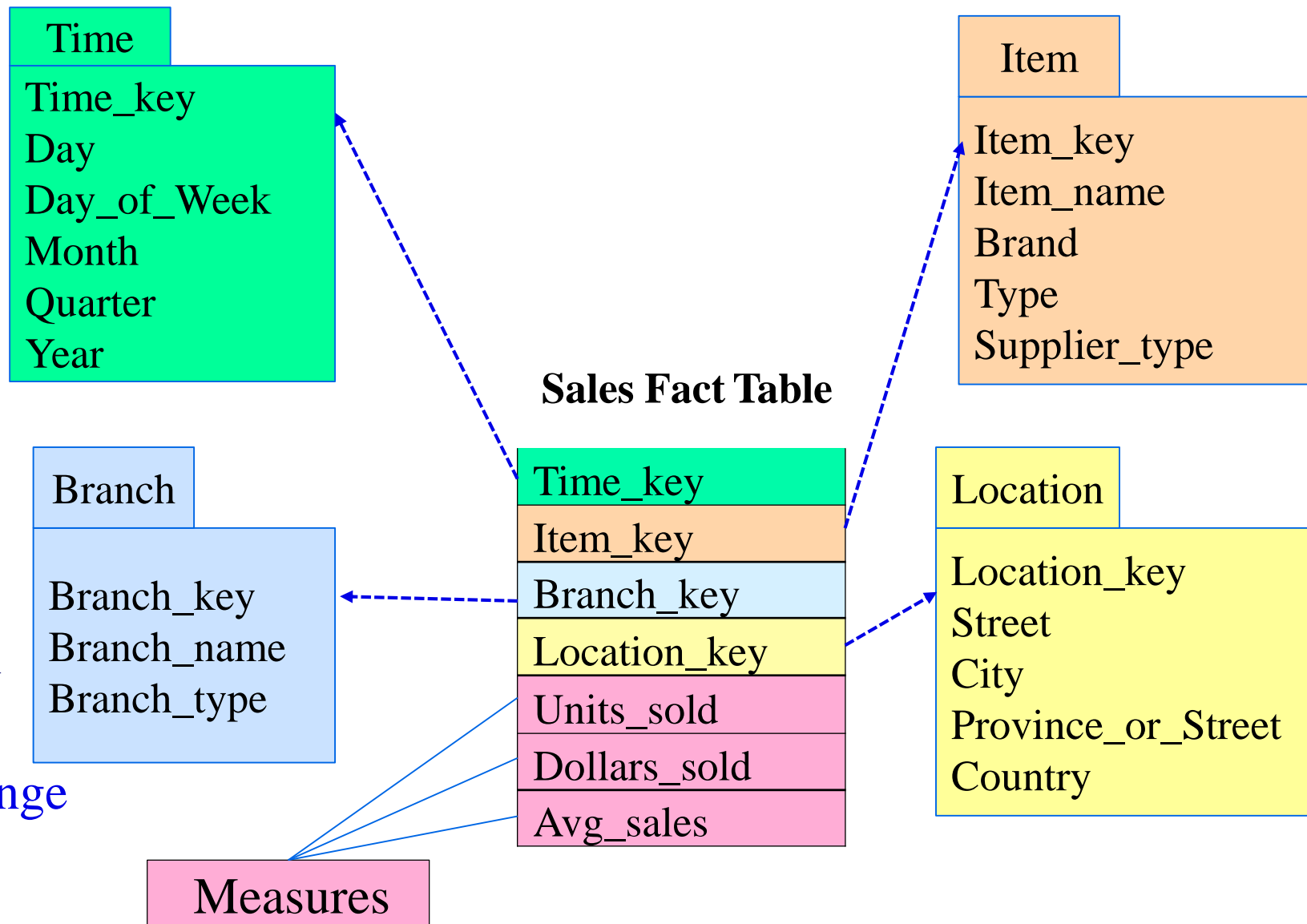➤ Holistic (tổng thể) measures : Cannot be computed from other sub-aggregate values

- E.g., median, mode, rank

# Star Schema

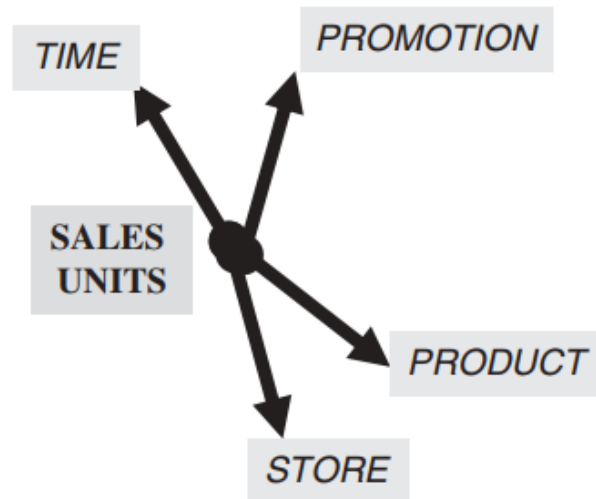- Fact table per business process / event, plus relevant dimensions
- Benefits
  - Easier to understand
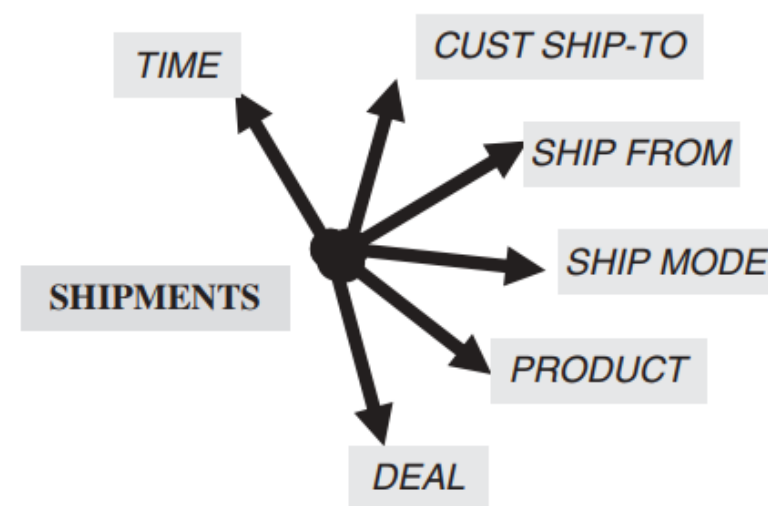  - Better performance from fewer joins
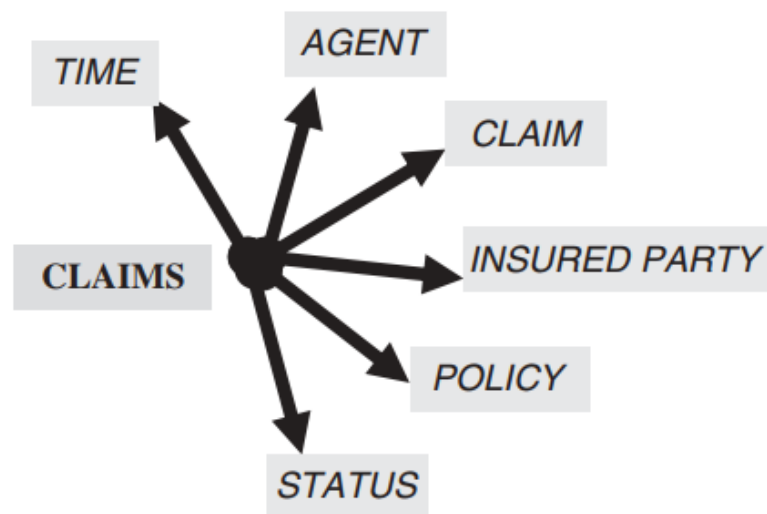  - Extensible to handle change

**Time**

Time_key
Day
Day_of_Week
Month
Quarter
Year

**Item**

Item_key
Item_name
Brand
Type
Supplier_type

**Branch**

Branch_key
Branch_name
Branch_type

**Sales Fact Table**

| Time_key |
| Item_key |
| Branch_key |
| Location_key |
| Units_sold |
| Dollars_sold |
| Avg_sales |

**Location**

Location_key
Street
City
Province_or_Street
Country

**Measures**

# Star Schema

- Example

### Supermarket Chain



TIME
PROMOTION
SALES UNITS
PRODUCT
STORE

### Manufacturing Company



TIME
CUST SHIP-TO
SHIP FROM
SHIP MODE
SHIPMENTS
PRODUCT
DEAL

### Insurance Business



TIME
AGENT
CLAIM
CLAIMS
INSURED PARTY
POLICY
STATUS

### Airlines Company



TIME
CUSTOMER
FLIGHT
FREQUENT FLYER FLIGHTS
FARE CLASS
AIRPORT
STATUS

# Snowflake Schema

**Time**
| |
|---|
| Time_key |
| Day |
| Day_of_Week |
| Month |
| Quarter |
| Year |

**Item**
| |
|---|
| Item_key |
| Item_name |
| Brand |
| Type |
| Supplier_key |

**Supplier**
| |
|---|
| Supplier_key |
| Supplier_type |

**Sales Fact Table**
| |
|---|
| Time_key |
| Item_key |
| Branch_key |
| Location_key |
| Units_sold |
| Dollars_sold |
| Avg_sales |

**Branch**
| |
|---|
| Branch_key |
| Branch_name |
| Branch_type |

**Location**
| |
|---|
| Location_key |
| Street |
| City_key |

**City**
| |
|---|
| City_key |
| City_name |
| State_or_province |
| Country |

**Measures**

# Assignment: ETL to Star schema

# Assignment: ETL to Star schema

# Practice (ELT to Star Schema)

# Questions

1. What is a multidimensional model, data cube, hierarchies?

2. What are elements of Dimensional Model, conformed dimensions, enterprise DW bus architecture, DW bus matrix?

3. What are facts, measures, measure classification?

4. What is Star schema, fact table, dimension table, Snowflake Schema?

# DATA WAREHOUSE

- Steps to Create dimensional modelling
- OLAP operations
- Slowly changing dimensions
- Rapidly changing dimensions
- Demo OLAP

Pivot

Drill Down

Slicing

Dicing

Roll UP

# Steps to Create Dimensional Modelling

■ The accuracy in creating your Dimensional modelling determines the success of your data warehouse implementation.

1. Identify Business Process
2. Identify Grain (độ mịn) (level of detail)
3. Identify Dimensions
4. Identify Facts
5. Build Schema

■ The model should describe

➤ Why, How much, When/Where/Who
➤ What of your business process

Select the Business Process • Why

Declare the Grain • How Much

Identify the Dimension • 3WS

Identify • What

Build the Schema

# Steps to Create Dimensional Modelling

- Identify Business Process
  - Identifying the actual business process a DW should cover.
    - Marketing, Sales, HR, etc.
    - Depends on the quality of data available for that process.
  - It is the most important step of the Data Modelling process
  - Type of business processes
    - Transaction
    - Accumulating Snapshot
    - Periodic Snapshot

# Steps to Create Dimensional Modelling

Identify Business Process

> **Transactions processes**

- The most basic fact grain
- One row is a transaction
- Ex. Sales, Return...

**CUSTOMER**
- CUSTOMER_ID
- CUSTOMER_NUMBER
- CUSTOMER_NAME
- [MORE.....]

**BRANCH**
- BRANCH_ID
- BRANCH_NAME
- BRANCH_CODE
- BRANCH_DESCRIPTION
- BRANCH_ADDRESS
- [MORE....]

**TRANSACTION_DATE**
- [TRANSACTION_DATE_ID ]
- TRANSACTION_DATE
- [TRANSACTION CALENDAR DAY]
- [MORE...]

**DAILY ACCOUNT**
- TRANSACTION_DATE_ID
- TRANSACTION_ID
- ACCOUNT_ID
- CUSTOMER_ID
- BRANCH_ID
- TRANSACTION_AMOUNT
- TRANSACTION_FEES

**ACCOUNT**
- ACCOUNT_ID
- ACCOUNT_NUMBER
- ACCOUNT_OPEN_DATE
- ACCOUNT_TYPE_CODE
- [MORE...]

**TRANSACTION_TYPE**
- TRANSACTION_TYPE_ID
- TRANSACTION_TYPE
- [MORE...]

(Deposit or Withdrawal) Transaction Type

**6 Rows inserted (1 row for each Transaction)**
Row 1: Withdrawal: $400,  Date: 2nd August 2,2005,    Time: 4:00AM
Row 2: Deposit:       $300, Date: 4th August 4,2005,     Time: 3:00AM
Row 3: Withdrawal: $600,  Date: 5nd August 5,2005,    Time: 2:00PM
Row 4: Withdrawal: $900,  Date: 6th August 6,2005,     Time: 9:00PM
Row 5: Deposit:       $900, Date: 18th August 18,2005, Time: 7:00AM
Row 6 :Deposit:       $800, Date: 23rd August 23,2005, Time: 1:00AM

# Steps to Create Dimensional Modelling

■ Identify Business Process

➤ **Accumulating Snapshot**

■ Capture a business process workflow

■ Fact row is initially inserted, then updated as milestones (mốc) occur

■ Ex. Order fulfillment, Job application tracking…

# Steps to Create Dimensional Modelling

Identify Business Process

## ➤ Periodic Snapshot

- At predetermined intervals snapshots of the same level of details are taken and stacked consecutively in the fact table

- Snapshots can be taken daily, weekly, monthly…

- Ex. Financial reports, Bank account values, GPA…



**ACCOUNT**
- 🔑 ACCOUNT_ID
- ACCOUNT_NUMBER
- ACCOUNT_OPEN_DATE
- ACCOUNT_TYPE_CODE
- [MORE…]

**CUSTOMER**
- 🔑 CUSTOMER_ID
- CUSTOMER_NUMBER
- CUSTOMER_NAME
- [MORE…..]

**MONTHLY_ACCOUNT**
- MONTH_END_DATE_ID
- ACCOUNT_ID
- CUSTOMER_ID
- [PREVIOUS BALANCE]
- [TOTAL DEPOSITS]
- [TOTAL WITHDRAWL]
- [AVAILABLE BALANCE]

**MONTH**
- 🔑 MONTH_END_DATE_ID
- [MONTH NAME]
- [CALENDAR MONTH ]
- [MORE…]

**1 Row for every Customer every Month**

# Steps to Create Dimensional Modelling

- Identify Grain
  - Describes the level of detail for the business problem/solution
  - The lowest level of information for any table in your data warehouse
  - Need to answer questions:
    - Do we need to store all the available products or just a few types of products?
    - Do we store the product sale information on a monthly, weekly, daily or hourly basis?
    - How do the above two choices affect the database size?

# Steps to Create Dimensional Modelling

- Identify Dimensions
  - Dimensions provide context for facts
  - We can easily identify dimensions because of "by", and/or "for" words. Like date, store, inventory, etc.
    - Ex. we want to find the sales for specific products in different locations on a daily basis.
      - Dimensions: Product, Location and Time
      - Attributes:  for Product: Product key (Foreign Key), Name, Type, Specifications
      - Hierarchies: For Location: Country, State, City, Street Address, Name

# Steps to Create Dimensional Modelling

- Identify Facts

  ➤ Facts are quantifiable numerical values associated with the business process

  ➤ This step is co-associated with the business users (they get access to data).

  ➤ Most of the fact table rows are numerical values like price or cost per unit, etc.

  ➤ Ex. we want to find the sales for specific products in different locations on a daily basis.

    - The fact is Sum of Sales by product by location by time.

# Steps to Create Dimensional Modelling

■ Build Schema

➤ Star Schema

➤ Snowflake Schema

# Rules for Dimensional Modelling

■ Load atomic (nguyên tử) data into dimensional structures.

■ Build dimensional models around business processes.

■ Every fact table has an associated date dimension table.

■ All facts in a single fact table are at the same grain or level of detail.

■ Dimension tables use a surrogate key (không liên hệ với DL, ex. auto number)

■ Balance requirements and realities to deliver business solution to support their decision-making

# Slowly Changing Dimensions

■ Dimensional data changes infrequently but when it does, you need a strategy for addressing the change.

➤ Ex. What happens when a customer has a new address, an employee has a name change?

■ **4 popular Strategies**

➤ Type 1: Overwrite the existing attribute

➤ Type 2: Add a new Dimension row

➤ Type 3: Add a new dimension attribute

➤ Type 6: (1+2+3)

# Slowly Changing Dimensions

**Type 1: Overwrite the existing attribute**

- Appropriate for
  - Correcting mistakes or error in data
  - Change where historical associations do not matter
  - The old value has no significance
- Ex. Employee name changes, Corrections…

| Key | ID | Name | Region |
|-----|-----|------|--------|
| 123 | VA-13 | ACME Products | Northeast |
| 234 | PA-07 | Ace Products & Services | Northeast |

| Key | ID | Name | Region |
|-----|-----|------|--------|
| 123 | VA-13 | ACME Products | Mid-Atlantic |
| 234 | PA-07 | Ace Products & Services | Northeast |

# Slowly Changing Dimensions

- **Type 2: Add a new Dimension row**
  - ➤ Most popular strategy, as it preserves history
  - ➤ Natural key is repeated (surrogate key (không liên hệ với DL, ex. auto number)
  - ➤ Old and new values are stored along with effective dates and indicator of which row is "current"

| Key | ID | Name | Region | ACTV RCRD | ACTV START | ACTV END |
|-----|-------|---------------|-----------|-----------|------------|----------|
| 123 | VA-13 | ACME Products | Northeast | 1 | 20140328 | 99999999 |
| 234 | PA-07 | Ace Products | Northeast | 1 | 20140508 | 99999999 |

| Key | ID | Name | Region | ACTV RCRD | ACTV START | ACTV END |
|-----|-------|---------------|--------------|-----------|------------|----------|
| 123 | VA-13 | ACME Products | Northeast | 0 | 20140328 | 20160728 |
| 234 | PA-07 | Ace Products | Northeast | 1 | 20140508 | 99999999 |
| 784 | VA-13 | ACME Products | Mid-Atlantic | 1 | 20160729 | 99999999 |

# Slowly Changing Dimensions

**Type 3: Add a new dimension attribute**

➤ Infrequently used, preserves history

➤ Useful for soft changes where users might want to choose between the old and new attribute, or need to access both values for a time

➤ The new value is written to the existing column, the old value is stored in a new column

➤ This way queries do not have to be re-written to access the new attribute

| Key | ID | Name | Region | Previous Region |
|-----|-----|------|--------|-----------------|
| 123 | VA-13 | Ace Hardware | Northeast | |
| 234 | PA-07 | Ace Products | Northeast | |

| Key | ID | Name | Region | Previous Region |
|-----|-----|------|--------|-----------------|
| 123 | VA-13 | Ace Hardware | Mid-Atlantic | Northeast |
| 234 | PA-07 | Ace Products | Northeast | |

# Slowly Changing Dimensions

**Type 6: (1+2+3)**

➤ Type 6 is a very rarely used (hiếm khi sử dụng)

➤ Start with a Type 2, add columns for the records you wish to capture the current value as well as the historical value. This allows one to filter or group on the Type 2 value in effect when the measure occurred or the current attribute value

| Key | ID | Name | Current Region | Historical Region | ACTV RCRD | ACTV RCRD Start | ACTV RCRD End |
|-----|------|--------------|----------------|-------------------|-----------|-----------------|---------------|
| 123 | VA-13 | ACE Hardware | Northeast | Northeast | 1 | 20140328 | 99999999 |

| Key | ID | Name | Current Region | Historical Region | ACTV RCRD | ACTV RCRD Start | ACTV RCRD End |
|-----|------|--------------|----------------|-------------------|-----------|-----------------|---------------|
| 123 | VA-13 | ACE Hardware | Mid-Atlantic | Northeast | 0 | 20140328 | 20160728 |
| 784 | VA-13 | ACE Hardware | Mid-Atlantic | Mid-Atlantic | 1 | 20160729 | 99999999 |

| Key | ID | Name | Current Region | Historical Region | ACTV RCRD | ACTV RCRD Start | ACTV RCRD End |
|-----|------|--------------|----------------|-------------------|-----------|-----------------|---------------|
| 123 | VA-13 | ACE Hardware | Virginia | Northeast | 0 | 20140328 | 20160728 |
| 784 | VA-13 | ACE Hardware | Virginia | Mid-Atlantic | 0 | 20160729 | 20161231 |
| 934 | VA-13 | ACE Hardware | Virginia | Virginia | 1 | 20170101 | 99999999 |

# Rapidly Changing Dimensions

■ The attribute values of dimension change frequently causing the dimension grow rapidly

■ The rapid growth of this dimension will impact maintenance and performance as the dimension grows.

➢ **Solution: Mini-Dimension**

➤ Mini-dimensions contain the rapidly changing attributes of the original dimension

```
                DIM_CUSTOMER
P  *  CUST_KEY                    NUMBER
      CUST_NAME                   VARCHAR2
      CUST_CITY                   VARCHAR2
      CUST_STATE                  VARCHAR2
      CUST_AGE                    NUMBER
      CUST_INCOME                 NUMBER
      CUST_LIFETIME_PURCHASES     NUMBER
      CUST_RATING                 VARCHAR2
      CUST_ACCOUNT_STATUS         VARCHAR2
      CUST_CREDIT_SCORE           NUMBER
      CUST_GENDER                 VARCHAR2
      CUST_ACTV_RCRD_FL           NUMBER
      CUST_ACTV_RCRD_START_DT     DATE
      CUST_ACTV_RCRD_END_DT       DATE

⚬ DIM_CUSTOMER_PK (CUST_KEY)
```

Rapidly changing attributes: customer's age, income, the number of lifetime purchases, rating, account status, and credit score

```
                DIM_CUSTOMER
P  *  CUST_KEY                    NUMBER
      CUST_NAME                   VARCHAR2
      CUST_CITY                   VARCHAR2
      CUST_STATE                  VARCHAR2
      CUST_ATTRIBUTE_KEY          NUMBER
      CUST_GENDER                 VARCHAR2
      CUST_ACTV_RCRD_FL           NUMBER
      CUST_ACTV_RCRD_START_DT     DATE
      CUST_ACTV_RCRD_END_DT       DATE

⚬ DIM_CUSTOMER_PK (CUST_KEY)
```
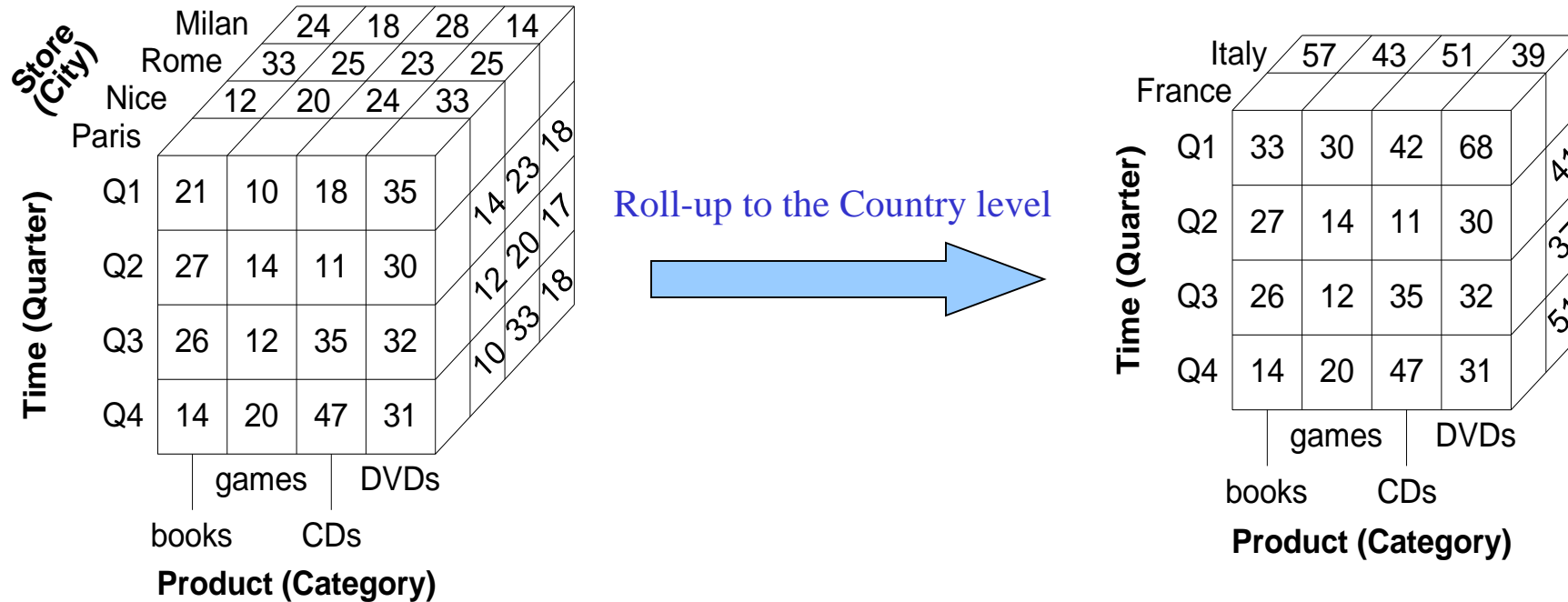
```
            DIM_CUSTOMER_ATTRIBUTE
P  *  CUST_ATTR_KEY                   NUMBER
      CUST_ATTR_AGE                   VARCHAR2
      CUST_ATTR_INCOME                VARCHAR2
      CUST_ATTR_LIFETIME_PURCHASES    VARCHAR2
      CUST_ATTR_RATING                VARCHAR2
      CUST_ATTR_ACCOUNT_STATUS        VARCHAR2
      CUST_ATTR_CREDIT_SCORE          VARCHAR2
      CUST_ATTR_ACTV_RCRD_FL          NUMBER
      CUST_ATTR_ACTV_RCRD_START_DT    DATE
      CUST_ATTR_ACTV_RCRD_END_DT      DATE

⚬ DIM_CUSTOMER_ATTRIBUTE_PK (CUST_ATTR_KEY)
```
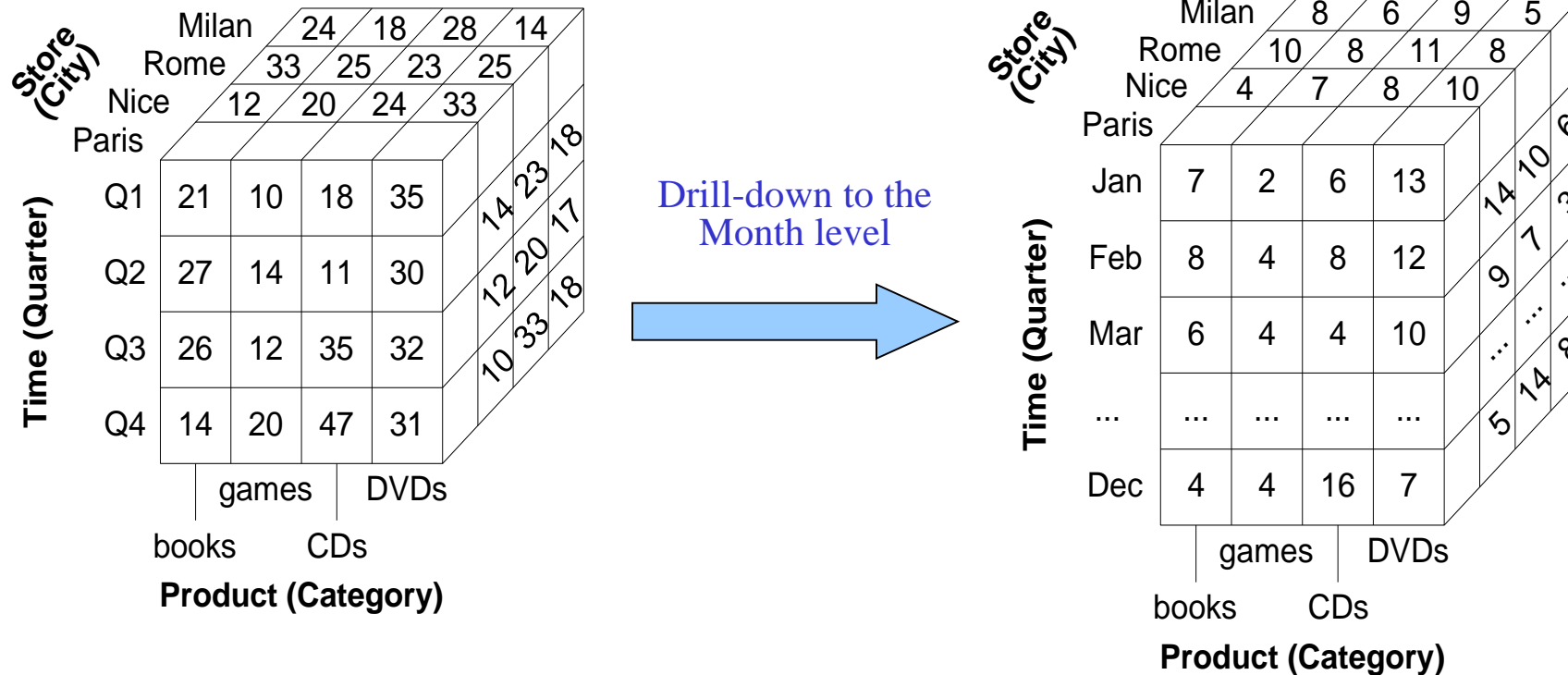
# OLAP Operations: Roll up

Transforms detailed measures into summarized ones when one moves up in a hierarchy
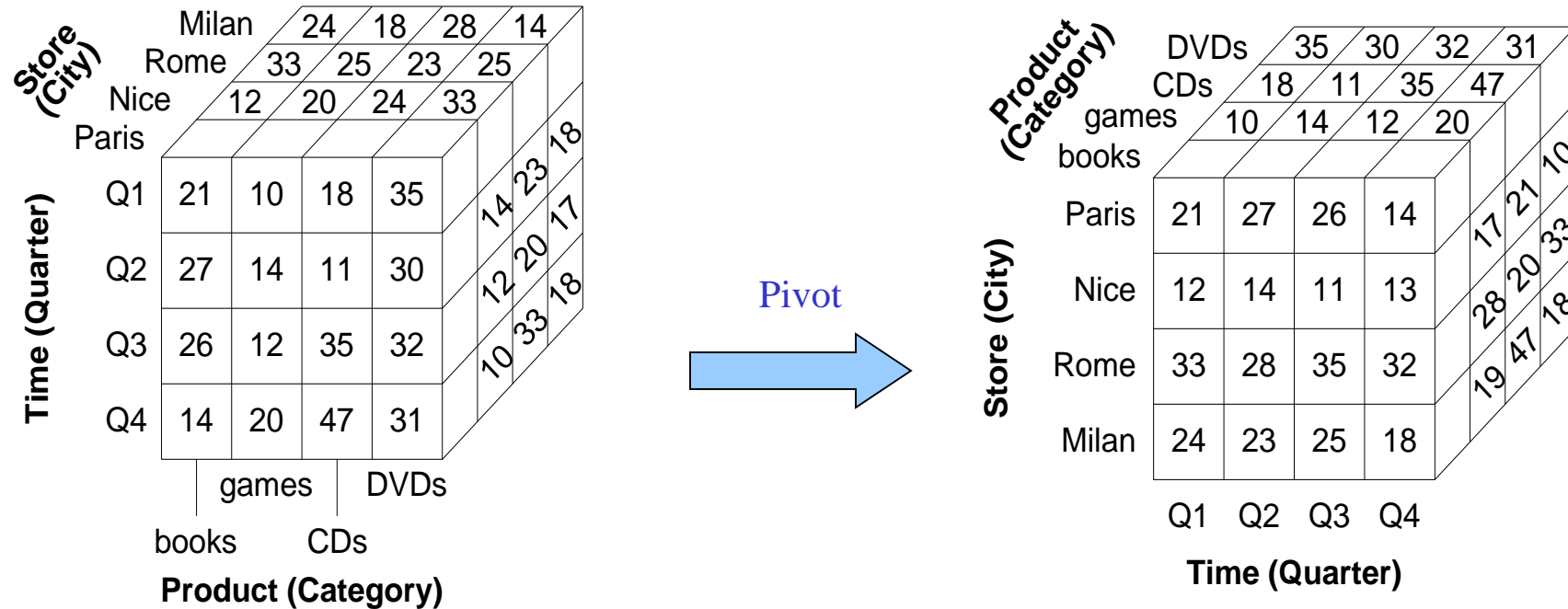


Roll-up to the Country level

# OLAP Operations: Drill down

■ Opposite to the roll-up operation, i.e., it moves from a more general level to a detailed level in a hierarchy
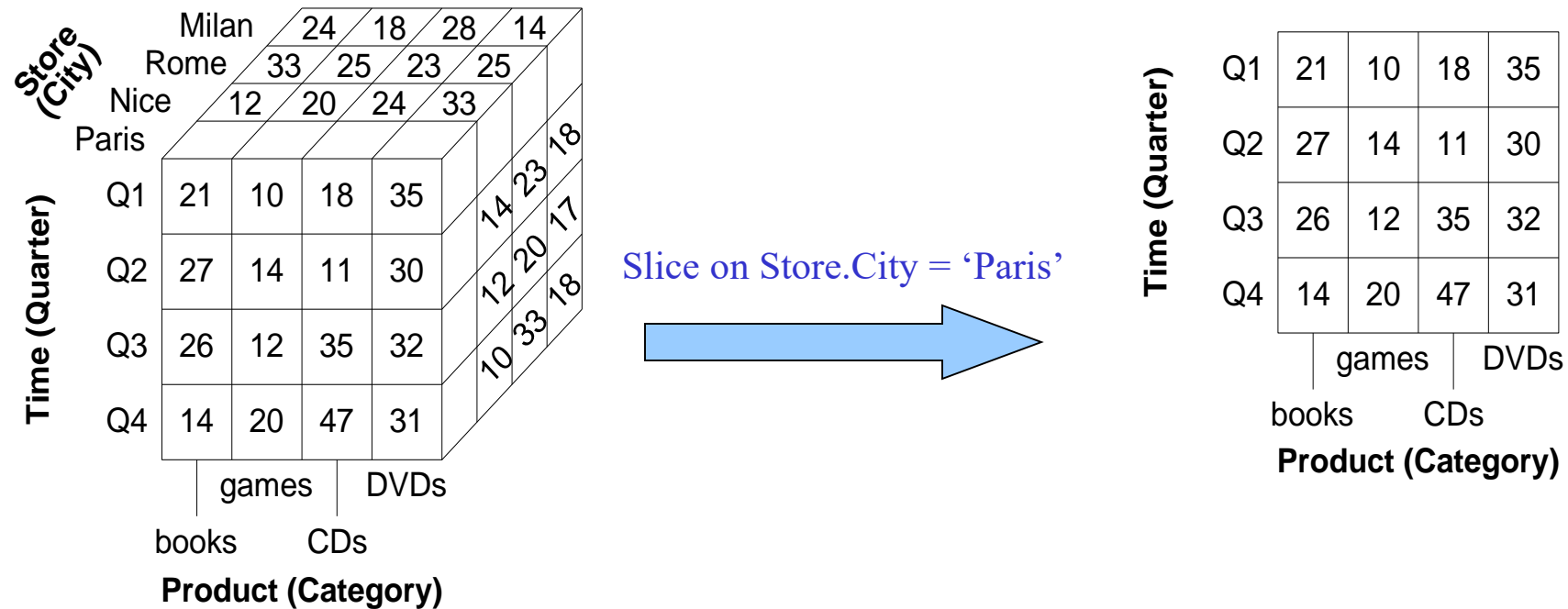


Drill-down to the Month level

# OLAP Operations: Pivot or Rotate

Rotates the axes of a cube to provide an alternative presentation of the data

# OLAP Operations: Slice

- Performs a selection on one dimension of a cube, resulting in a subcube



Slice on Store.City = 'Paris'

# OLAP Operations: Dice

- Defines a selection on two or more dimensions, thus again defining a subcube



Dice on Store.Country = 'France'
and Time.Quarter= 'Q1' or 'Q2'

# Questions

1. How to build dimensional modelling?

2. What is Slowly changing dimension?

3. What is Rapidly changing dimension?

4. What is OLAP (Online Analytical Processing)?

5. What are OLAP operations?

# DATA WAREHOUSE

- MDX Definition
- MDX Query Syntax
- Demo MDX

# MDX (MultiDimensional Expressions)

- **Syntax (basic)**

  - **SELECT <measures|dimensions> [on columns|rows]**
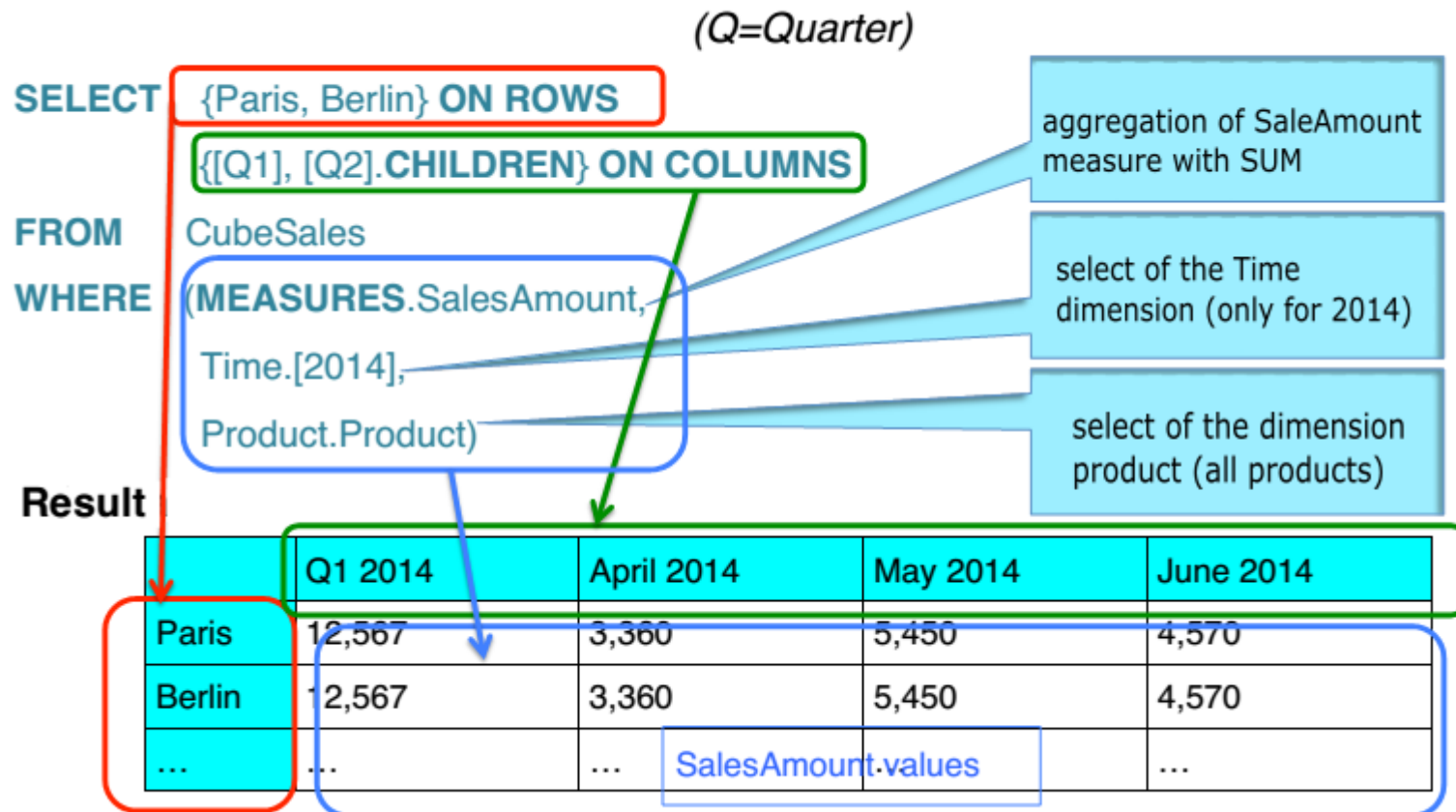
  - **FROM <Cube>**

  - **[WHERE <Slicing conditions>]**

- **How to use MDX?**

- **Create Datawarehouse (OLAP|Cube) using SSIS**

- **Using MDX to query Data in the cube**

# MDX (MultiDimensional Expressions)

# MDX (MultiDimensional Expressions)

**Demo**

1. Total Sales

2. Total Sales by Country

2a. Total Sales of Vietnam

2b. Total Sales of Vietnam

3. Total Sales by Year

4. Total Sales by Year (skip null value)

5. Total Sales by Year of Vietnam

6. Total Sales and Profit by Year (skip null value)

7. Total Sale by Year and Sub-Category

# DATA WAREHOUSE

- What is Meta Data?
- Why metadata is important
- A Critical Need in the Data Warehouse
- RAID Technology

# What is Meta Data?

■ Metadata is data about the data or documentation about the information which is required by the users.

■ Metadata includes the following:

➤ The location and descriptions of warehouse systems and components.

➤ Names, definitions, structures, and content of data-warehouse and end-users views.

➤ Identification of authoritative data sources.

➤ Integration and transformation rules used to populate data.

➤ Integration and transformation rules used to deliver information to end-user analytical tools.

➤ Subscription information for information delivery to analysis subscribers.

➤ Metrics used to analyze warehouses usage and performance.

➤ Security authorizations, access control list, etc.

# Why metadata is important

■ Metadata in a data warehouse contains the answers to questions about the data in the data warehouse.

➤ You keep the answers in a place called the metadata repository.

➤ Here is a sample list of definitions:

- Data about the data
- Table of contents for the data
- Catalog for the data
- Data warehouse atlas (bản đồ)
- Data warehouse roadmap (lộ trình)
- Data warehouse directory
- Glue that holds the data warehouse contents together

# Why metadata is important

Metadata element for the Customer entity

| Entity Name: | Customer |
|---|---|
| Alias Names: | Account, Client |

Definition:      A person or an organization that purchases goods or services from the company.

Remarks:      Customer entity includes regular, current, and past customers.

Source Systems:      Finished Goods Orders, Maintenance Contracts, Online Sales.

| | |
|---|---|
| Create Date: | January 15, 1999 |
| Last Update Date: | January 21, 2001 |
| Update Cycle: | Weekly |
| Last Full Refresh Date: | December 29, 2000 |
| Full Refresh Cycle: | Every six months |
| Data Quality Reviewed: | January 25, 2001 |
| Last Deduplication: | January 10, 2001 |
| Planned Archival: | Every six months |
| Responsible User: | Jane Brown |

**Figure 9-1**    Metadata element for *Customer* entity.

# A Critical Need in the Data Warehouse

■ For Using the Data Warehouse

➢ Users retrieve information from the data warehouse.

➢ Users themselves create ad hoc queries and run these against the data warehouse.

➢ They format their own reports.

➢ Before they can create and run their queries, users need to know about the data in the data warehouse

=> They need metadata.

# A Critical Need in the Data Warehouse

- For Building the Data Warehouse.
  - Metadata is absolutely essential for building your data warehouse in every activity and every task.
    - Know the source systems and their data structures.
    - Know the structures and the data content in the data warehouse.
    - Determine the mappings and the data transformations.

# A Critical Need in the Data Warehouse

For Administering the Data Warehouse

**Data Extraction/Transformation/Loading**

How to handle data changes?

How to include new sources?

Where to cleanse the data? How to change the data cleansing methods?

How to cleanse data after populating the warehouse?

How to switch to new data transformation techniques?

How to audit the application of ongoing changes?

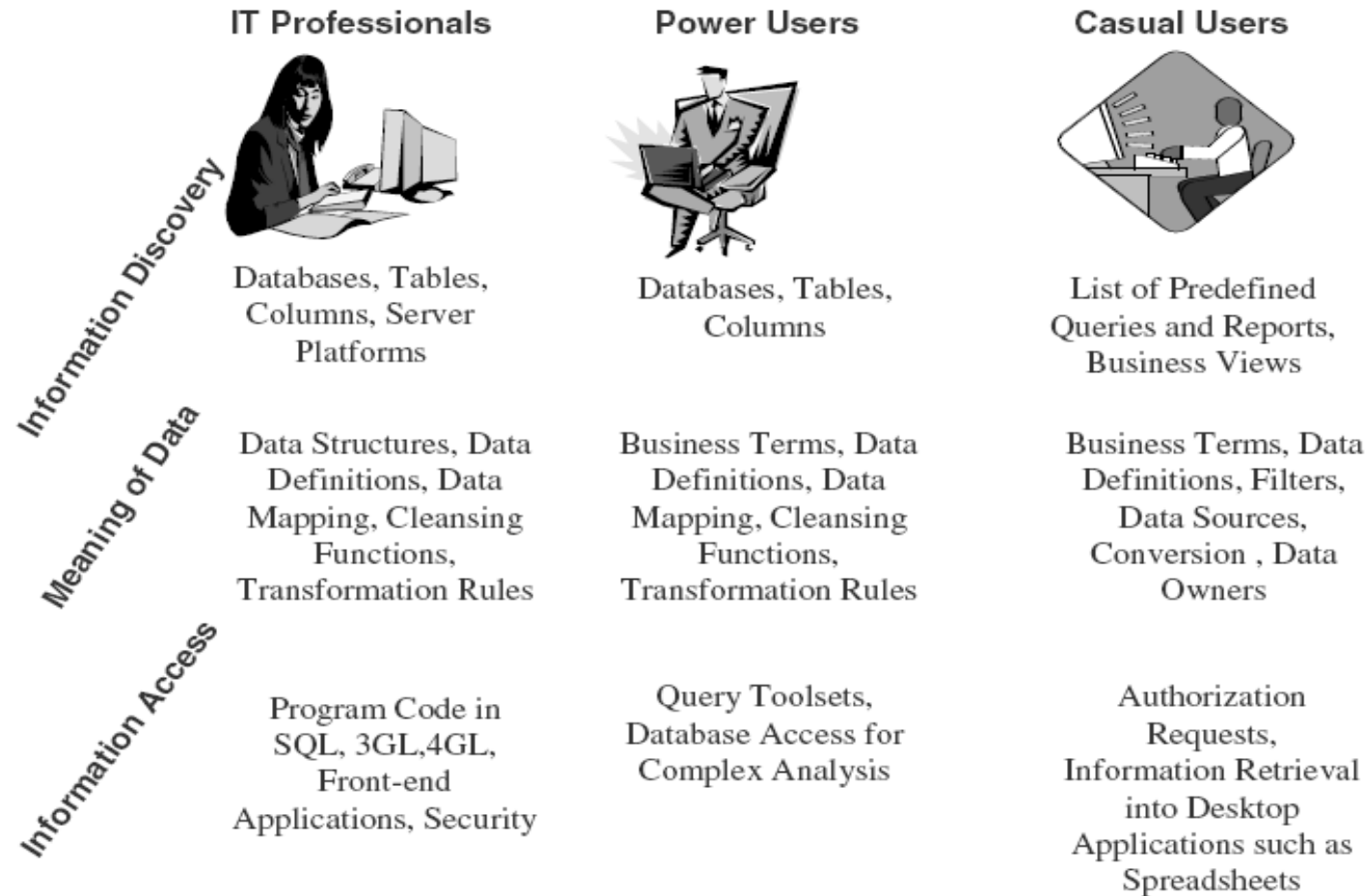**Data from External Sources**

How to add new external data sources?

How to drop some external data sources?

When mergers and acquisitions happen, how to bring in new data to the warehouse?

How to verify all external data on ongoing basis?

**Data Warehouse**

How to add new summary tables?

How to control runaway queries?

How to expand storage?

When to schedule platform upgrades?

How to add new information delivery tools for the users?

How to continue ongoing training?

How to maintain and enhance user support function?

How to monitor and improve ad hoc query performance?

When to schedule backups?

How to perform disaster recovery drills?

How to keep data definitions up-to-date?

How to maintain the security system?

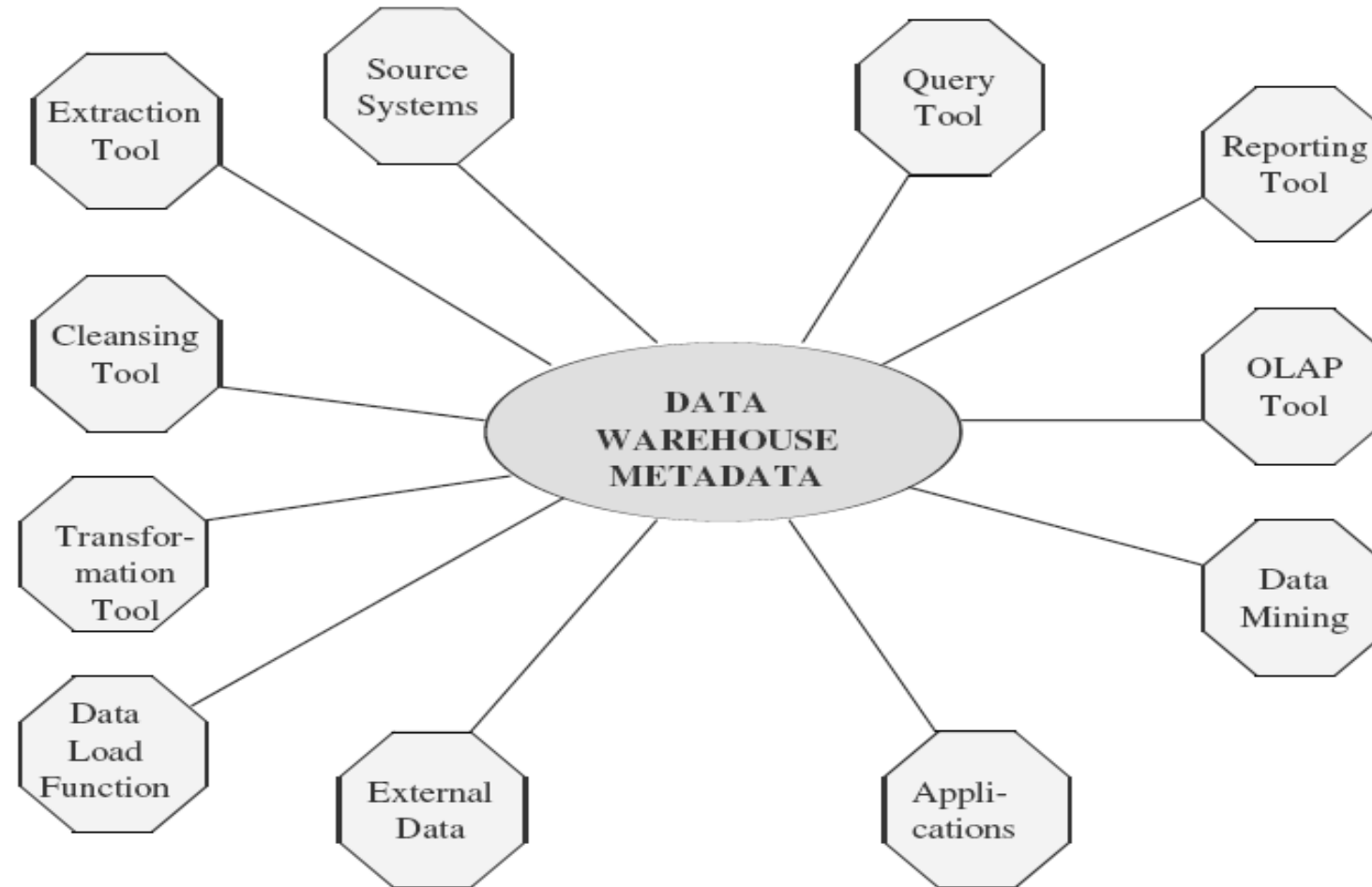How to monitor system load distribution?

# A Critical Need in the Data Warehouse

## Who Needs Metadata

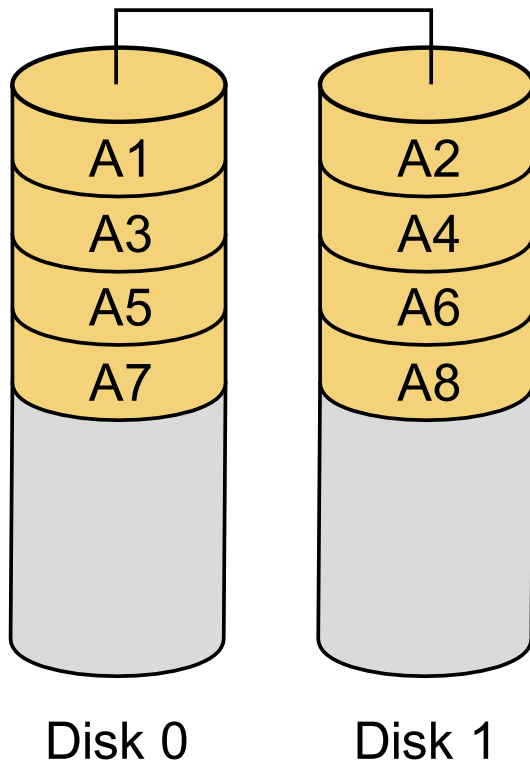| | IT Professionals | Power Users | Casual Users |
|---|---|---|---|
| **Information Discovery** | Databases, Tables, Columns, Server Platforms | Databases, Tables, Columns | List of Predefined Queries and Reports, Business Views |
| **Meaning of Data** | Data Structures, Data Definitions, Data Mapping, Cleansing Functions, Transformation Rules | Business Terms, Data Definitions, Data Mapping, Cleansing Functions, Transformation Rules | Business Terms, Data Definitions, Filters, Data Sources, Conversion, Data Owners |
| **Information Access** | Program Code in SQL, 3GL,4GL, Front-end Applications, Security | Query Toolsets, Database Access for Complex Analysis | Authorization Requests, Information Retrieval into Desktop Applications such as Spreadsheets |

# A Critical Need in the Data Warehouse

Metadata is Like a Nerve Center

# RAID Technology

## Some types of RAID



RAID 10 (RAID 1+0)
Mirror + Stripe

RAID 0

Disk 0    Disk 1

RAID 1

Disk 0    Disk 1

RAID 5

Disk 0    Disk 1    Disk 2    Disk 3
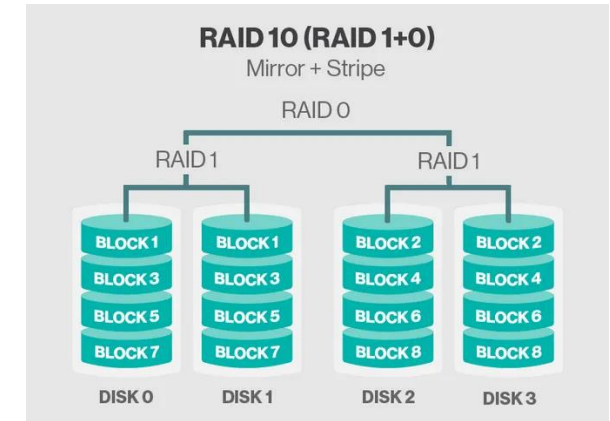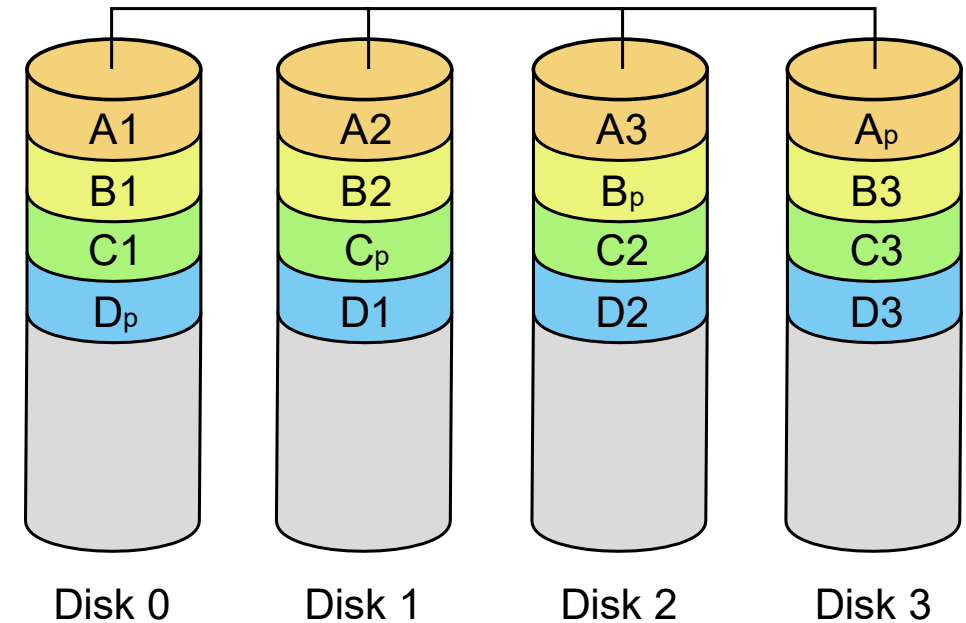
# Questions

1. What is MDX (MultiDimensional eXpressions) language? How do we use MDX to retrieve data in a data warehouse?

2. What is metadata? What are metadata types? (methods for classification of metadata; detail in Metadata types by functional areas)

3. What is the RAID technology? Explain more detail RAID 0, 1, 5.