

Metadata Repository

- Meta data is the data defining warehouse objects. It stores:
 - Description of the structure of the data warehouse
 - schema, view, dimensions, hierarchies, derived data defn, data mart locations and contents
 - Operational meta-data
 - data lineage (history of migrated data and transformation path), currency of data (active, archived, or purged), monitoring information (warehouse usage statistics, error reports, audit trails)
 - The algorithms used for summarization
 - The mapping from operational environment to the data warehouse
 - Data related to system performance
 - warehouse schema, view and derived data definitions
 - Business data
 - business terms and definitions, ownership of data, charging policies

14

Chapter 4: Data Warehousing and On-line Analytical Processing

- Data Warehouse: Basic Concepts
- Data Warehouse Modeling: Data Cube and OLAP
- Data Warehouse Design and Usage
- Data Warehouse Implementation
- Summary

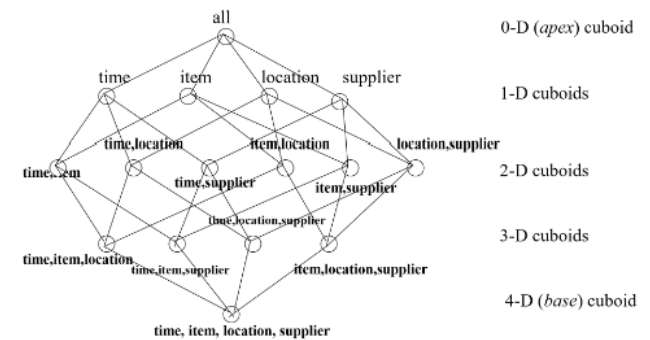
15

From Tables and Spreadsheets to Data Cubes

- A **data warehouse** is based on a multidimensional data model which views data in the form of a data cube
- A data cube, such as sales, allows data to be modeled and viewed in multiple dimensions
 - Dimension tables**, such as item (item_name, brand, type), or time(day, week, month, quarter, year) → *វិស័យ*
 - Fact table** contains **measures** (such as dollars_sold) and keys to each of the related dimension tables → *តម្លៃ*
- Data cube**: A lattice of cuboids
 - In data warehousing literature, an n-D base cube is called a **base cuboid**
 - The top most 0-D cuboid, which holds the highest-level of summarization, is called the **apex cuboid**
 - The lattice of cuboids forms a **data cube**.

16

Data Cube: A Lattice of Cuboids



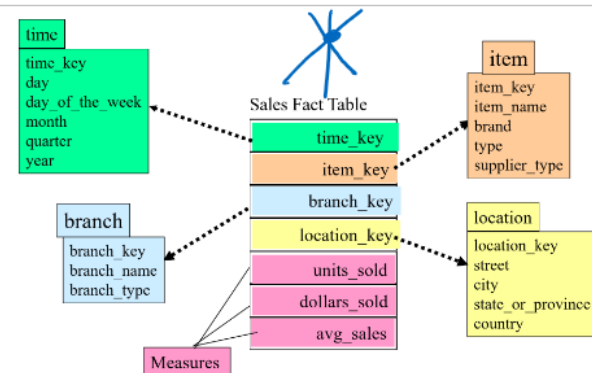
17

Conceptual Modeling of Data Warehouses

- Modeling data warehouses: dimensions & measures
 - Star schema**: A fact table in the middle connected to a set of dimension tables
 - Snowflake schema**: A refinement of star schema where some dimensional hierarchy is normalized into a set of smaller dimension tables, forming a shape similar to snowflake
 - Fact constellations**: Multiple fact tables share dimension tables, viewed as a collection of stars, therefore called **galaxy schema** or fact constellation

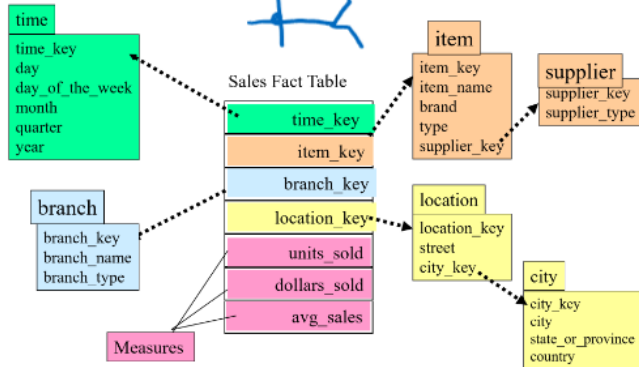
18

Star Schema: An Example



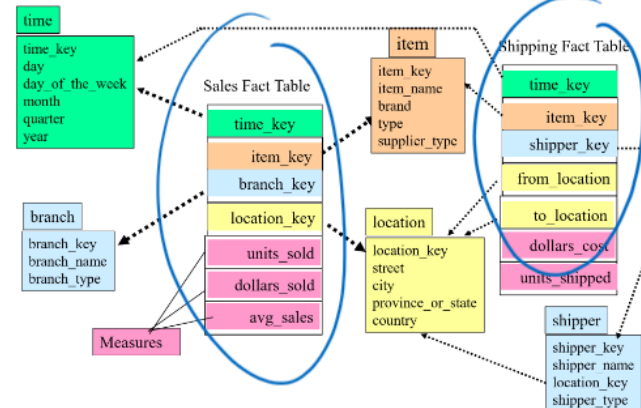
19

Snowflake Schema: An Example



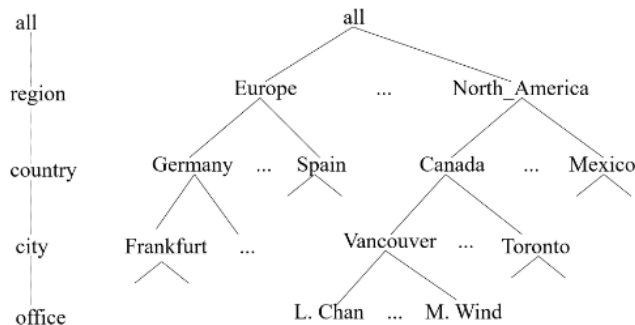
20

Fact Constellation: An Example



21

A Concept Hierarchy for a Dimension (location)



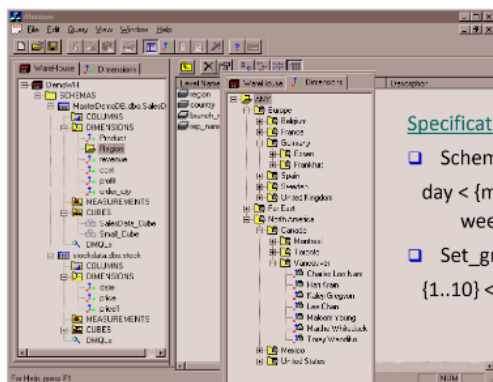
22

Data Cube Measures: Three Categories

- Distributive:** if the result derived by applying the function to n aggregate values is the same as that derived by applying the function on all the data without partitioning
 - E.g., count(), sum(), min(), max()
- Algebraic:** if it can be computed by an algebraic function with M arguments (where M is a bounded integer), each of which is obtained by applying a distributive aggregate function
 - avg(x) = sum(x) / count(x)
 - Is min_N() an algebraic measure? How about standard_deviation()?
- Holistic:** if there is no constant bound on the storage size needed to describe a subaggregate.
 - E.g., median(), mode(), rank()

23

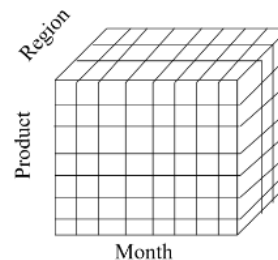
View of Warehouses and Hierarchies



24

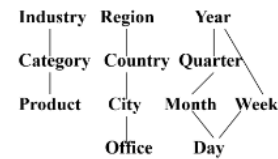
Multidimensional Data

- Sales volume as a function of product, month, and region



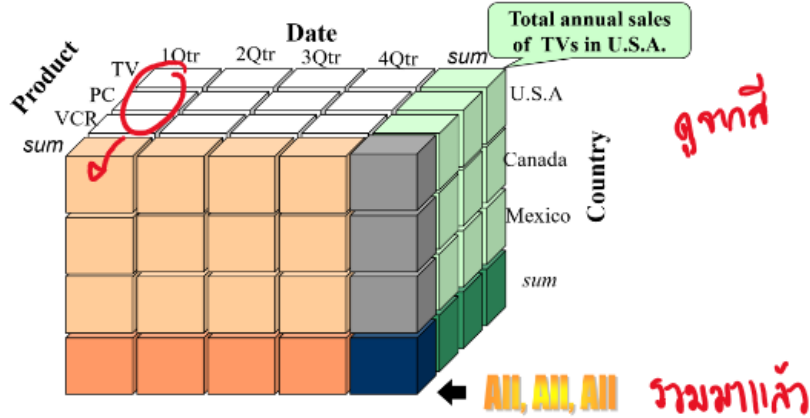
Dimensions: Product, Location, Time

Hierarchical summarization paths



25

A Sample Data Cube



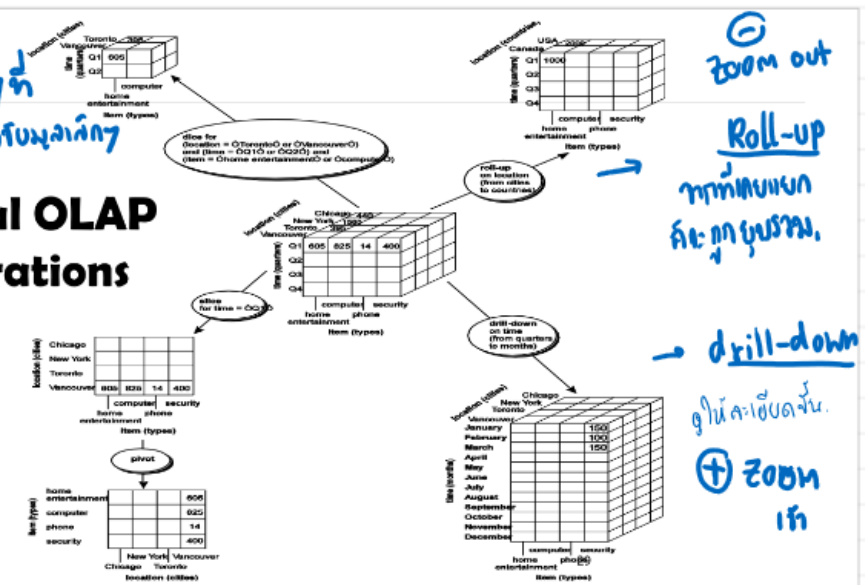
26

dice
ข้อมูลที่ใหญ่ที่สุด
ดูสถิติออกมาในมุมมองใหญ่

Typical OLAP Operations

Slice

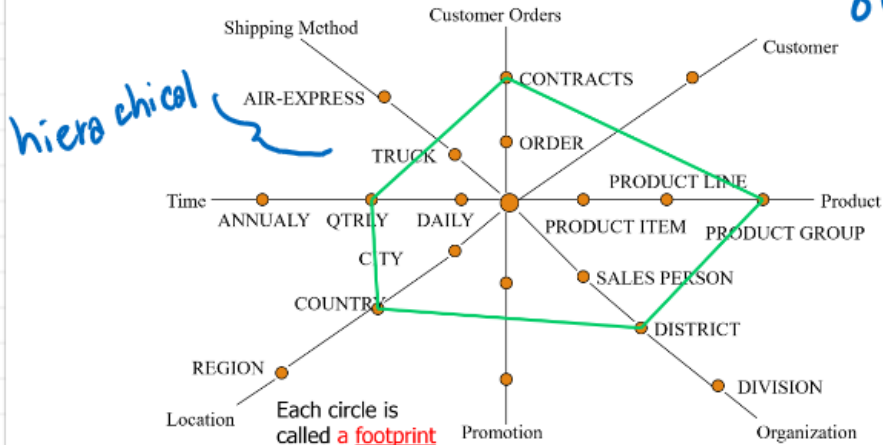
เลือก
ค่าเฉพาะ
มา 1 ตัว



29

A Star-Net Query Model

8 dimension



30



CS 412 Intro. to Data Mining

Chapter 4. Data Warehousing and On-line Analytical Processing

Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017



Chapter 4: Data Warehousing and On-line Analytical Processing

□ Data Warehouse: Basic Concepts



□ Data Warehouse Modeling: Data Cube and OLAP

□ Data Warehouse Design and Usage

□ Data Warehouse Implementation

□ Summary

3

What is a Data Warehouse?

- Defined in many different ways, but not rigorously
 - A decision support database that is maintained separately from the organization's operational database
 - Support information processing by providing a solid platform of consolidated, historical data for analysis
 - "A data warehouse is a subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process." —W. H. Inmon
- Data warehousing:
 - The process of constructing and using data warehouses

ข้อมูล
ประวัติ
ข้อมูล
คลัง

ข้อมูลสนับสนุน การตัดสินใจทางธุรกิจ

4

Data Warehouse—Subject-Oriented

- Organized around major subjects, such as customer, product, sales
- Focusing on the modeling and analysis of data for decision makers, not on daily operations or transaction processing
- Provide a simple and concise view around particular subject issues by excluding data that are not useful in the decision support process

5

Data Warehouse—Integrated

- Constructed by integrating multiple, heterogeneous data sources
 - relational databases, flat files, on-line transaction records
- Data cleaning and data integration techniques are applied.
 - Ensure consistency in naming conventions, encoding structures, attribute measures, etc. among different data sources
 - Ex. Hotel price: differences on currency, tax, breakfast covered, and parking
 - When data is moved to the warehouse, it is converted

6

Data Warehouse—Time Variant

- The time horizon for the data warehouse is significantly longer than that of operational systems
 - Operational database: current value data
 - Data warehouse data: provide information from a historical perspective (e.g., past 5-10 years)
- Every key structure in the data warehouse
 - Contains an element of time, explicitly or implicitly
 - But the key of operational data may or may not contain "time element"

7

Data Warehouse—Nonvolatile

- ❑ Independence
 - ❑ A **physically separate store** of data transformed from the operational environment
- ❑ Static: Operational **update of data does not occur** in the data warehouse environment
 - ❑ Does not require transaction processing, recovery, and concurrency control mechanisms
 - ❑ Requires only two operations in data accessing:
 - ❑ **initial loading of data** and **access of data**

8

OLTP vs. OLAP

- ❑ OLTP: Online transactional processing
 - ❑ DBMS operations
 - ❑ Query and transactional processing
- ❑ OLAP: Online analytical processing
 - ❑ Data warehouse operations
 - ❑ Drilling, slicing, dicing, etc.

| | OLTP | OLAP |
|--------------------|--|--|
| users | clerk, IT professional | knowledge worker |
| function | day to day operations | decision support |
| DB design | application-oriented | subject-oriented |
| data | current, up-to-date detailed, flat relational isolated | historical, summarized, multidimensional integrated, consolidated |
| usage | repetitive | ad-hoc |
| access | read/write index/hash on prim. key | lots of scans |
| unit of work | short, simple transaction | complex query |
| # records accessed | tens | millions |
| #users | thousands | hundreds |
| DB size | 100MB-GB | 100GB-TB |
| metric | transaction throughput | query throughput, response |

9

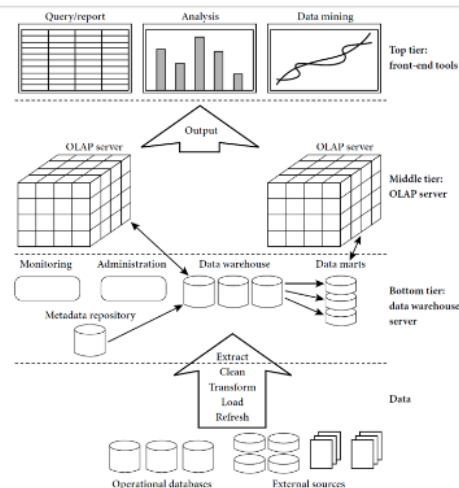
Why a Separate Data Warehouse?

- ❑ High performance for both systems
 - ❑ DBMS—tuned for OLTP: access methods, indexing, concurrency control, recovery
 - ❑ Warehouse—tuned for OLAP: complex OLAP queries, multidimensional view, consolidation
- ❑ Different functions and different data:
 - ❑ **missing data**: Decision support requires historical data which operational DBs do not typically maintain
 - ❑ **data consolidation**: DS requires consolidation (aggregation, summarization) of data from heterogeneous sources
 - ❑ **data quality**: different sources typically use inconsistent data representations, codes and formats which have to be reconciled
- ❑ Note: There are more and more systems which perform OLAP analysis directly on relational databases

10

Data Warehouse: A Multi-Tiered Architecture

- ❑ Top Tier: Front-End Tools
- ❑ Middle Tier: OLAP Server
- ❑ Bottom Tier: Data Warehouse Server
- ❑ Data



11

Three Data Warehouse Models

- ❑ **Enterprise warehouse**
 - ❑ Collects all of the information about subjects spanning the entire organization
- ❑ **Data Mart**
 - ❑ A subset of corporate-wide data that is of value to a specific groups of users
 - ❑ Its scope is confined to specific, selected groups, such as marketing data mart
 - ❑ Independent vs. dependent (directly from warehouse) data mart
- ❑ **Virtual warehouse**
 - ❑ A set of views over operational databases
 - ❑ Only some of the possible summary views may be materialized

12

Extraction, Transformation, and Loading (ETL)

- ❑ **Data extraction**
 - ❑ get data from multiple, heterogeneous, and external sources
- ❑ **Data cleaning**
 - ❑ detect errors in the data and rectify them when possible
- ❑ **Data transformation**
 - ❑ convert data from legacy or host format to warehouse format
- ❑ **Load**
 - ❑ sort, summarize, consolidate, compute views, check integrity, and build indices and partitions
- ❑ **Refresh**
 - ❑ propagate the updates from the data sources to the warehouse

13