



# **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

# 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

จัดเก็บข้อมูลมาไว้ที่เดียว

หัวข้อเฉพาะ/เป้าหมาย/เพื่ออะไร    รวมจากหลายแหล่ง    Data ที่จัดเก็บในช่วงเวลาหนึ่ง, ไม่เปลี่ยนแปลง

Data ที่ไว้ใช้วิเคราะห์

# 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**



# 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

# 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”

# 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*

# 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



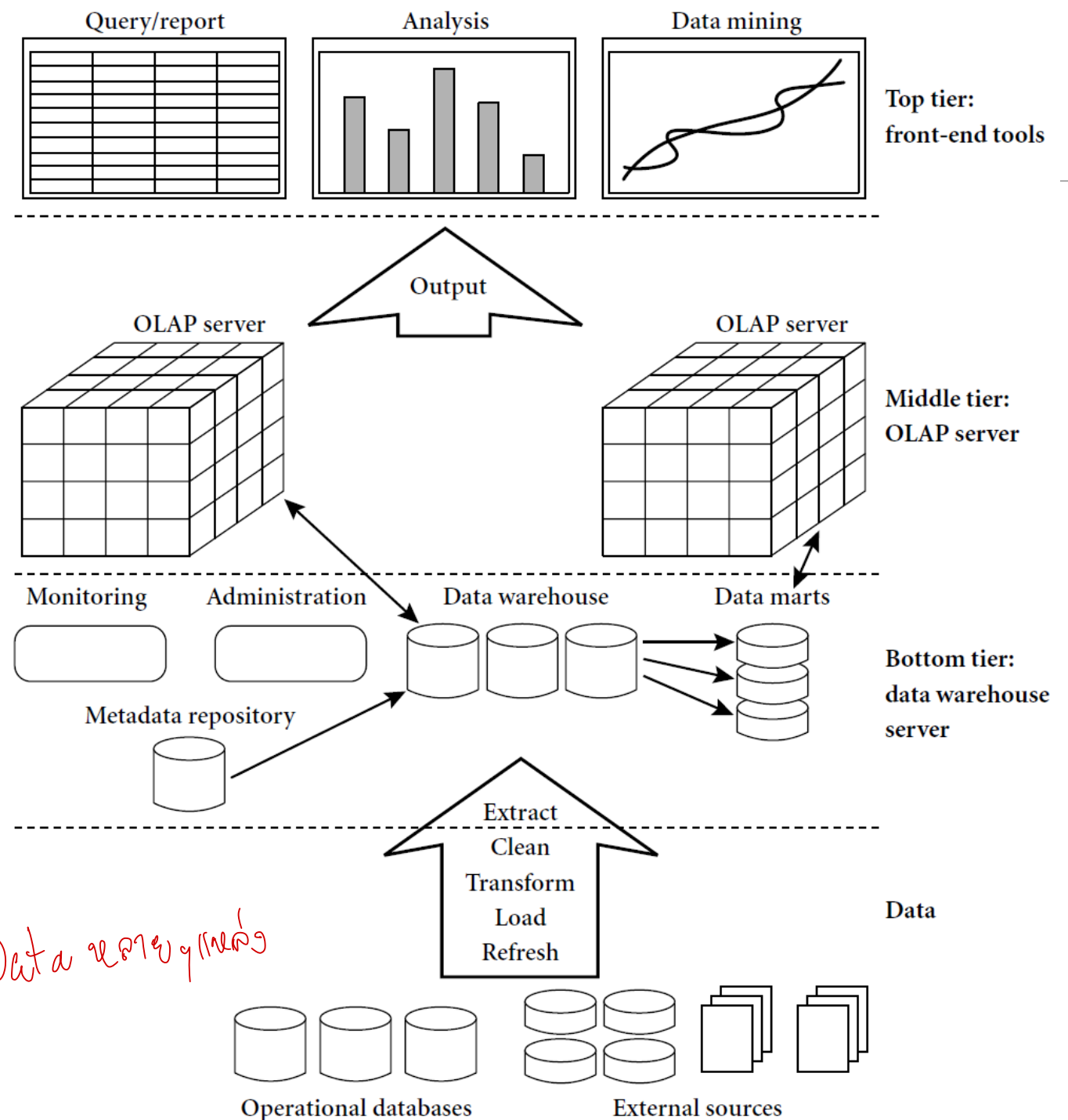
# Why a Separate Data Warehouse?

---

- ❑ 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

# Data Warehouse: A Multi-Tiered Architecture

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



# 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

# 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

# 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

# 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





# 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 เก็บใน data warehouse ไว้สองอัน
- ❑ **Dimension tables**, such as item (item\_name, brand, type), or time(day, week, month, quarter, year) *ข้อมูลสินค้า, สามารถจำแนกเป็นช่วงเวลา*
- ❑ **Fact table** contains **measures** (such as dollars\_sold) and keys to each of the related dimension tables เก็บตัวเลข
- ❑ **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**.

## 2D

2 ม, ๑๗๐๑  
ส่วนนี้เขาบอกได้ทั้งนี้?

Table 4.2: A 2-D view of sales data for *AllElectronics* according to the dimensions *time* and *item*, where the sales are from branches located in the city of Vancouver. The measure displayed is *dollars\_sold* (in thousands).

*location* = "Vancouver"

<i>time</i> (quarter)	<i>item</i> (type)			
	<i>home entertainment</i>	<i>computer</i>	<i>phone</i>	<i>security</i>
Q1	605	825	14	400
Q2	680	952	31	512
Q3	812	1023	30	501
Q4	927	1038	38	580

## 3D

3 ม, ๑๗๐๑  
1

Table 4.3: A 3-D view of sales data for *AllElectronics*, according to the dimensions *time*, *item*, and *location*. The measure displayed is *dollars\_sold* (in thousands).

<i>location</i> = "Chicago"					<i>location</i> = "New York"					<i>location</i> = "Toronto"					<i>location</i> = "Vancouver"				
<i>item</i>					<i>item</i>					<i>item</i>					<i>item</i>				
<i>time</i>	<i>home ent.</i>	<i>comp.</i>	<i>phone</i>	<i>sec.</i>	<i>home ent.</i>	<i>comp.</i>	<i>phone</i>	<i>sec.</i>		<i>home ent.</i>	<i>comp.</i>	<i>phone</i>	<i>sec.</i>		<i>home ent.</i>	<i>comp.</i>	<i>phone</i>	<i>sec.</i>	
Q1	854	882	89	623	1087	968	38	872		818	746	43	591		605	825	14	400	
Q2	943	890	64	698	1130	1024	41	925		894	769	52	682		680	952	31	512	
Q3	1032	924	59	789	1034	1048	45	1002		940	795	58	728		812	1023	30	501	
Q4	1129	992	63	870	1142	1091	54	984		978	864	59	784		927	1038	38	580	

# 3D

**location (cities)**

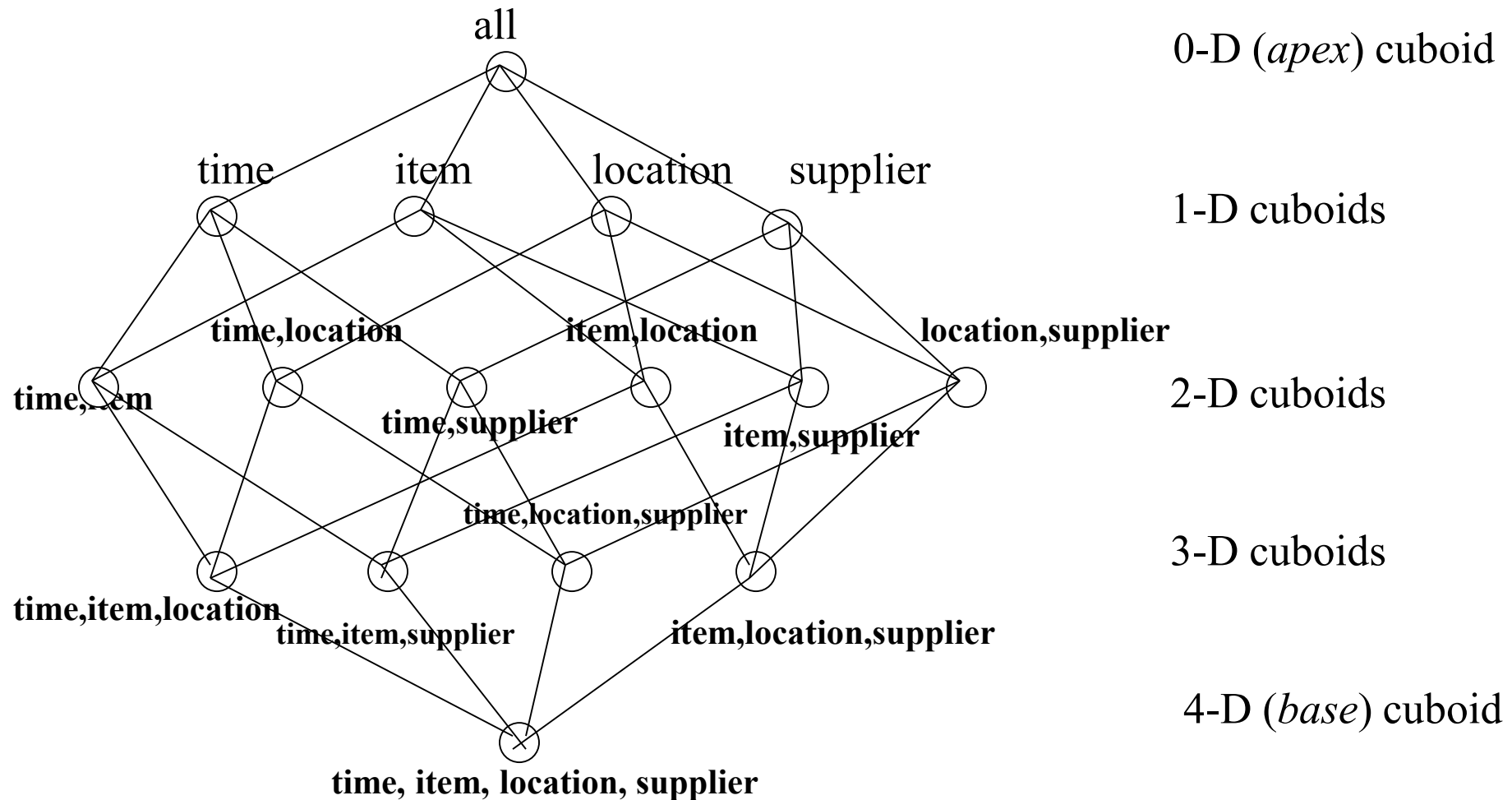
**time (quarters)**

**item (types)**

	Chicago	New York	Toronto	Vancouver
Q1	605	825	14	400
Q2	680	952	31	512
Q3	812	1023	30	501
Q4	927	1038	38	580

computer security  
home phone  
entertainment

# Data Cube: A Lattice of Cuboids



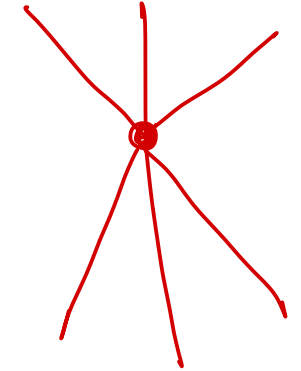
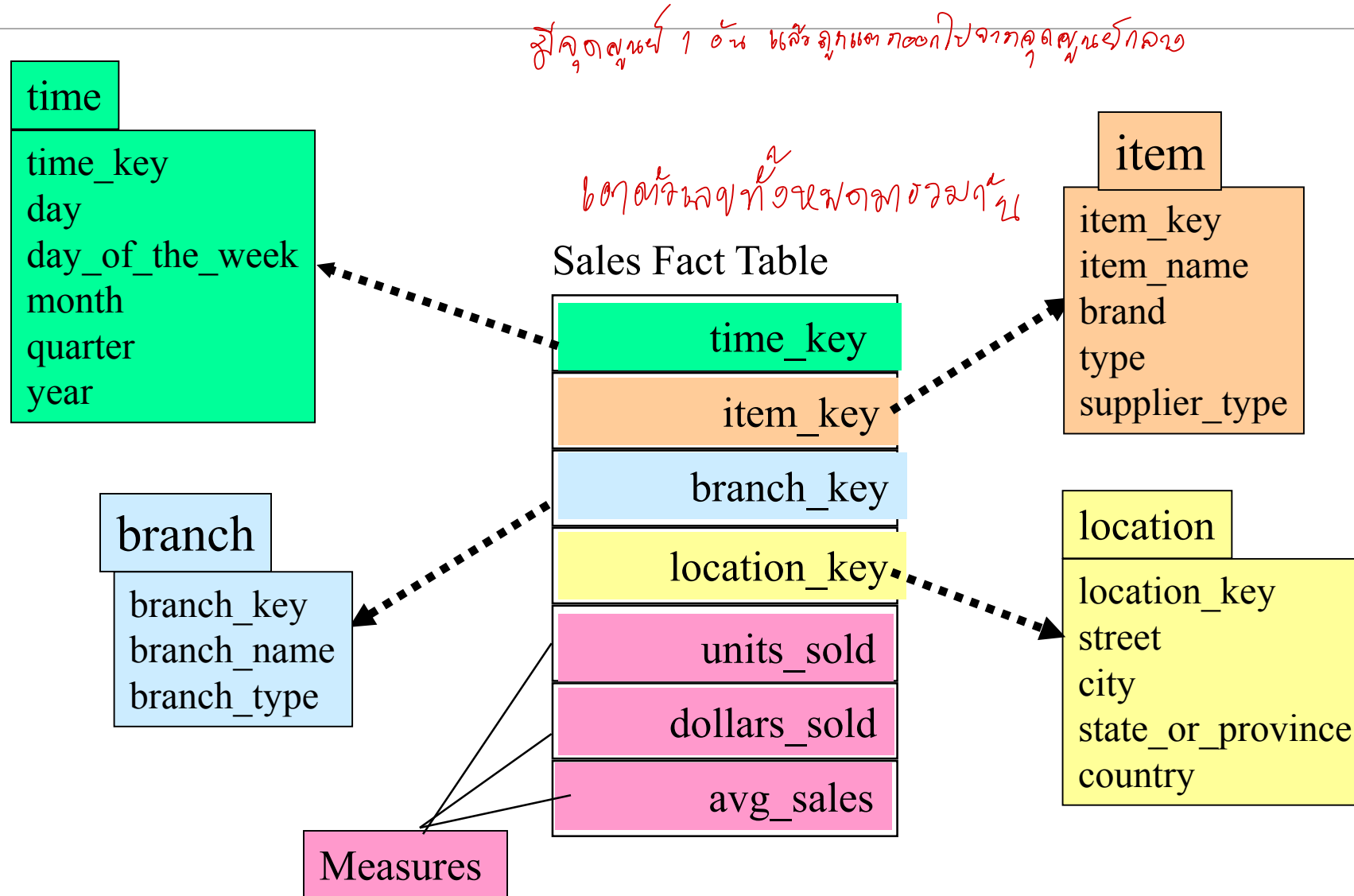


# Conceptual Modeling of Data Warehouses

---

- ❑ 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

# Star Schema: An Example



นี่คือรูปแบบ 1 อัน สำหรับเก็บข้อมูลจากข้อมูลดิบ

เก็บข้อมูลที่เกี่ยวข้องมาไว้ที่นี่

Sales Fact Table

item

item\_key  
item\_name  
brand  
type  
supplier\_type

time

time\_key  
day  
day\_of\_the\_week  
month  
quarter  
year

branch

branch\_key  
branch\_name  
branch\_type

location

location\_key  
street  
city  
state\_or\_province  
country

time\_key

item\_key

branch\_key

location\_key

units\_sold

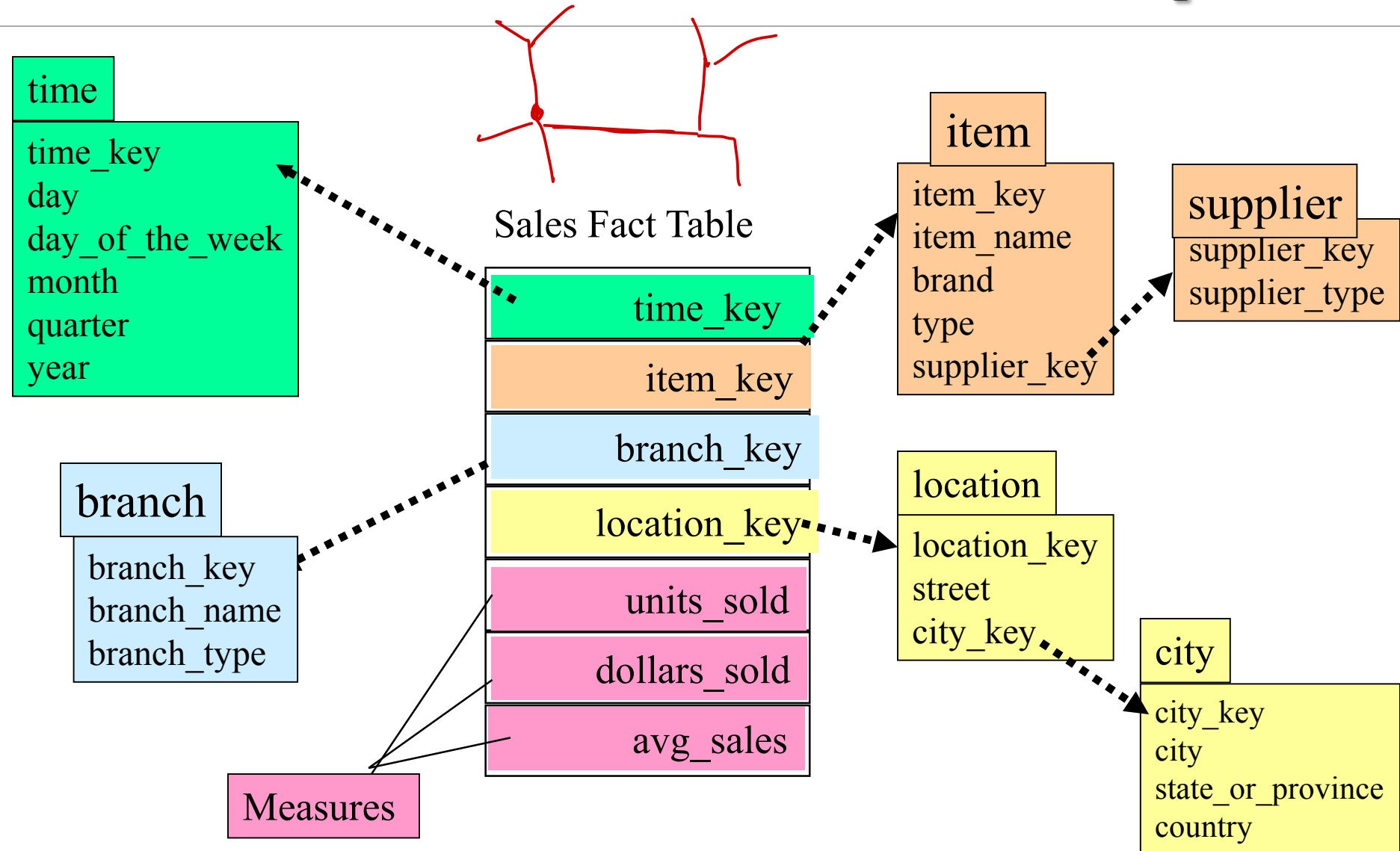
dollars\_sold

avg\_sales

Measures



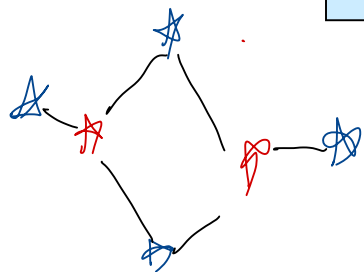
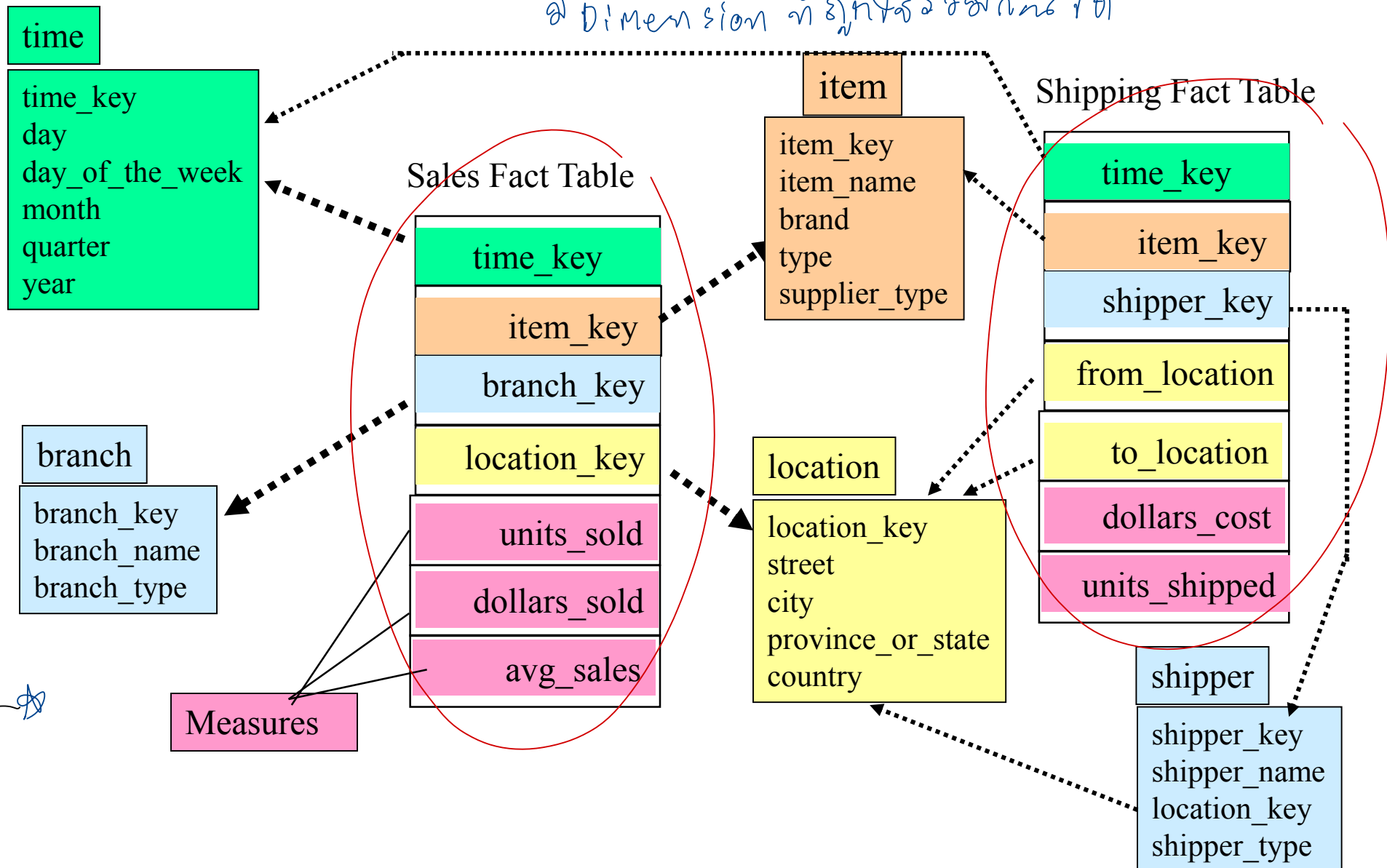
# Snowflake Schema: An Example



การเชื่อมโยง ที่ชัดเจน กับ

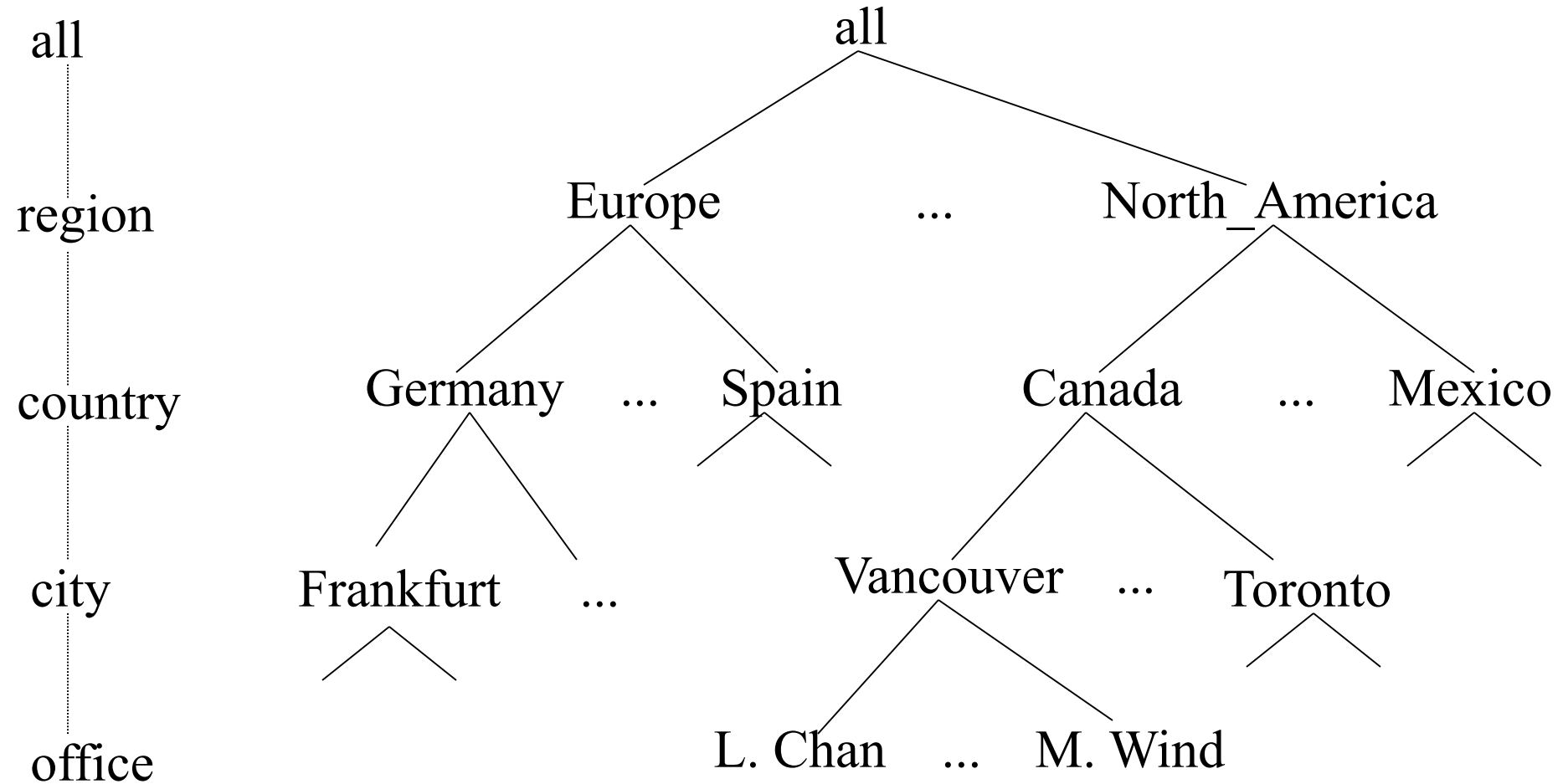
# Fact Constellation: An Example

คือ Dimension ที่ถูกใช้ร่วมกัน 7 ตัว



สามารถแบ่งแยกย่อยได้อีก

# A Concept Hierarchy for a Dimension (location)

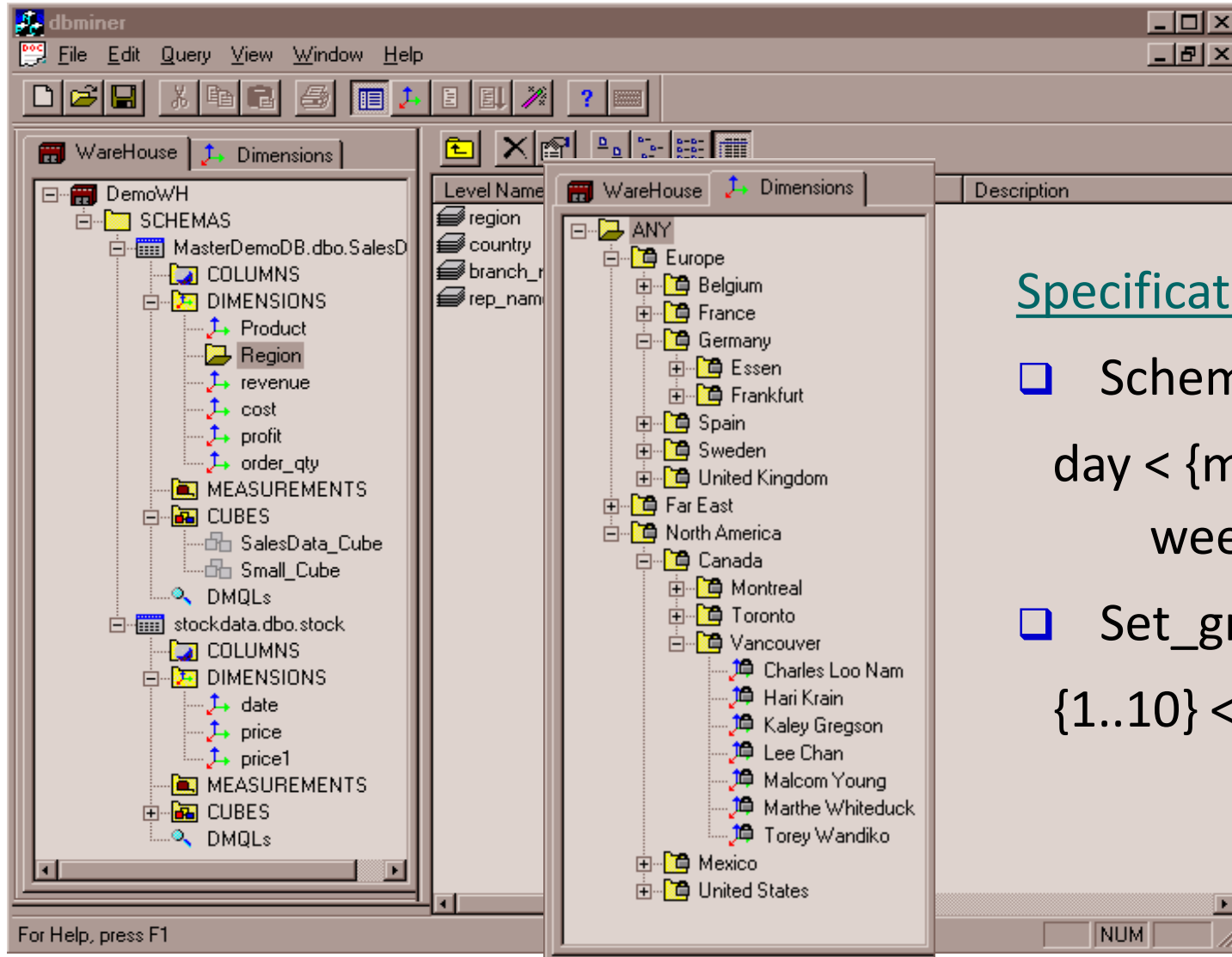


# Data Cube Measures: Three Categories

---

- ❑ 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
  - ❑  $\text{avg}(x) = \text{sum}(x) / \text{count}(x)$
  - ❑ 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()`

# View of Warehouses and Hierarchies



## Specification of hierarchies

- ☐ Schema hierarchy  
day < {month < quarter;  
week} < year
- ☐ Set\_grouping hierarchy  
{1..10} < inexpensive

# Multidimensional Data

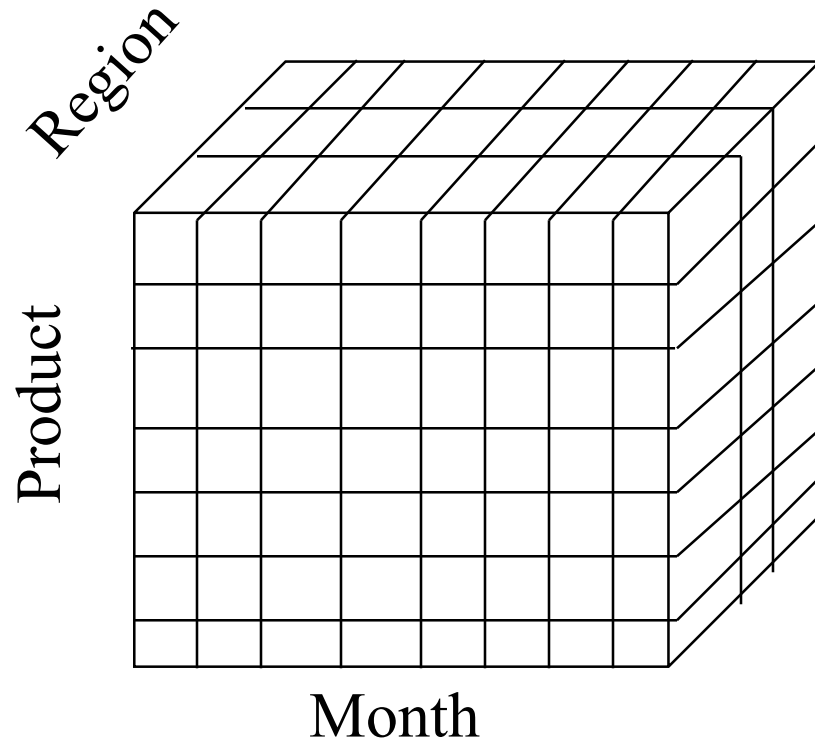
ပျော်လှယ်မှုများကို ဖော်ပြသော

ဒီမင်ရှင်

ရေဒီယ

အချိန်

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

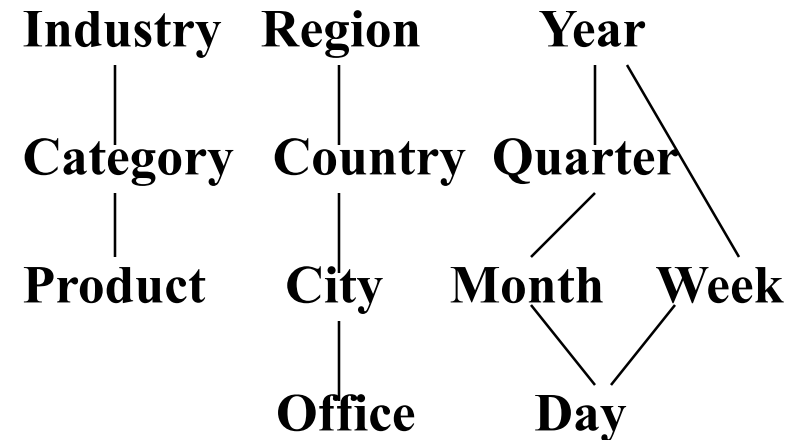


ဒီမင်ရှင်များကို  
အချိန်အတိုင်း

Dimensions: *Product, Location, Time*

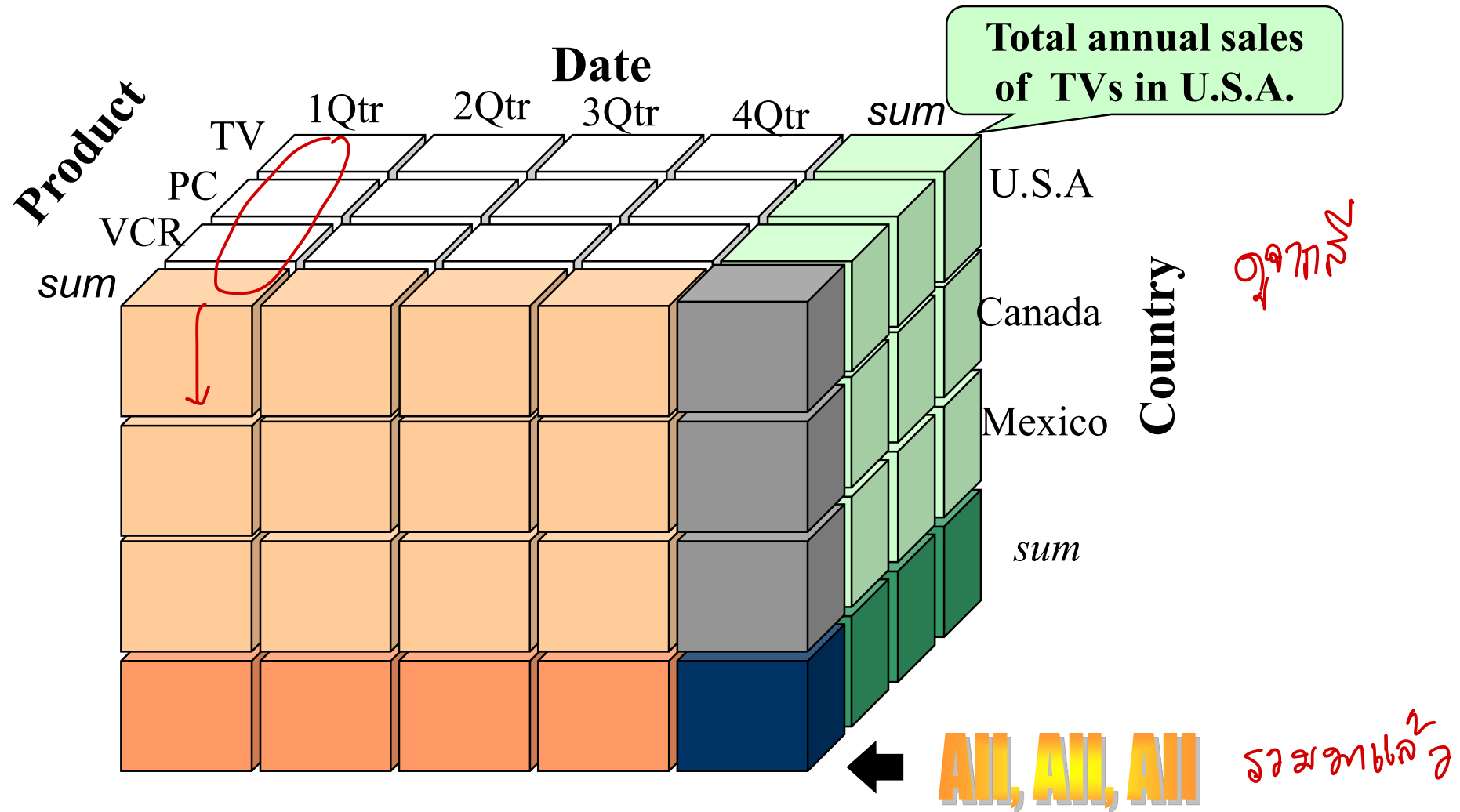
Hierarchical summarization paths

ပျော်လှယ်မှုများကို ဖော်ပြ



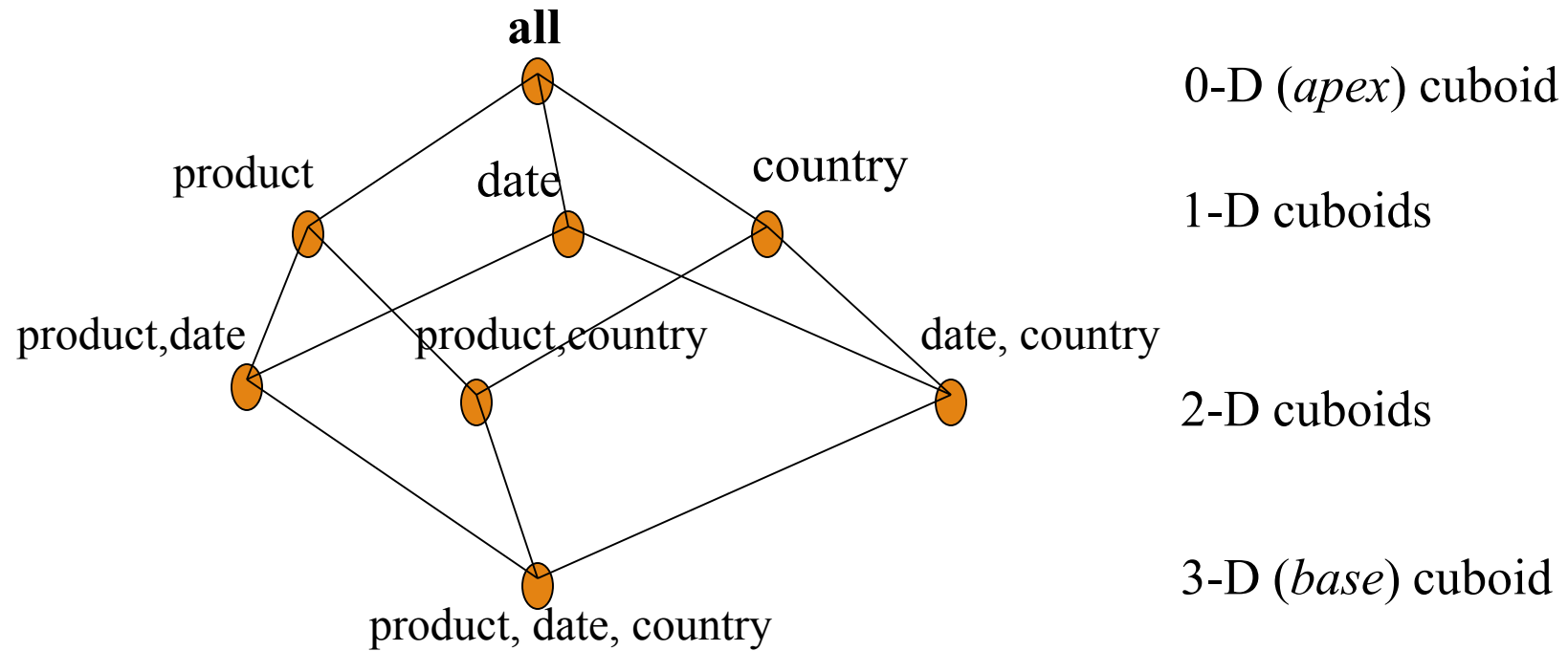


# A Sample Data Cube



# Cuboids Corresponding to the Cube

---

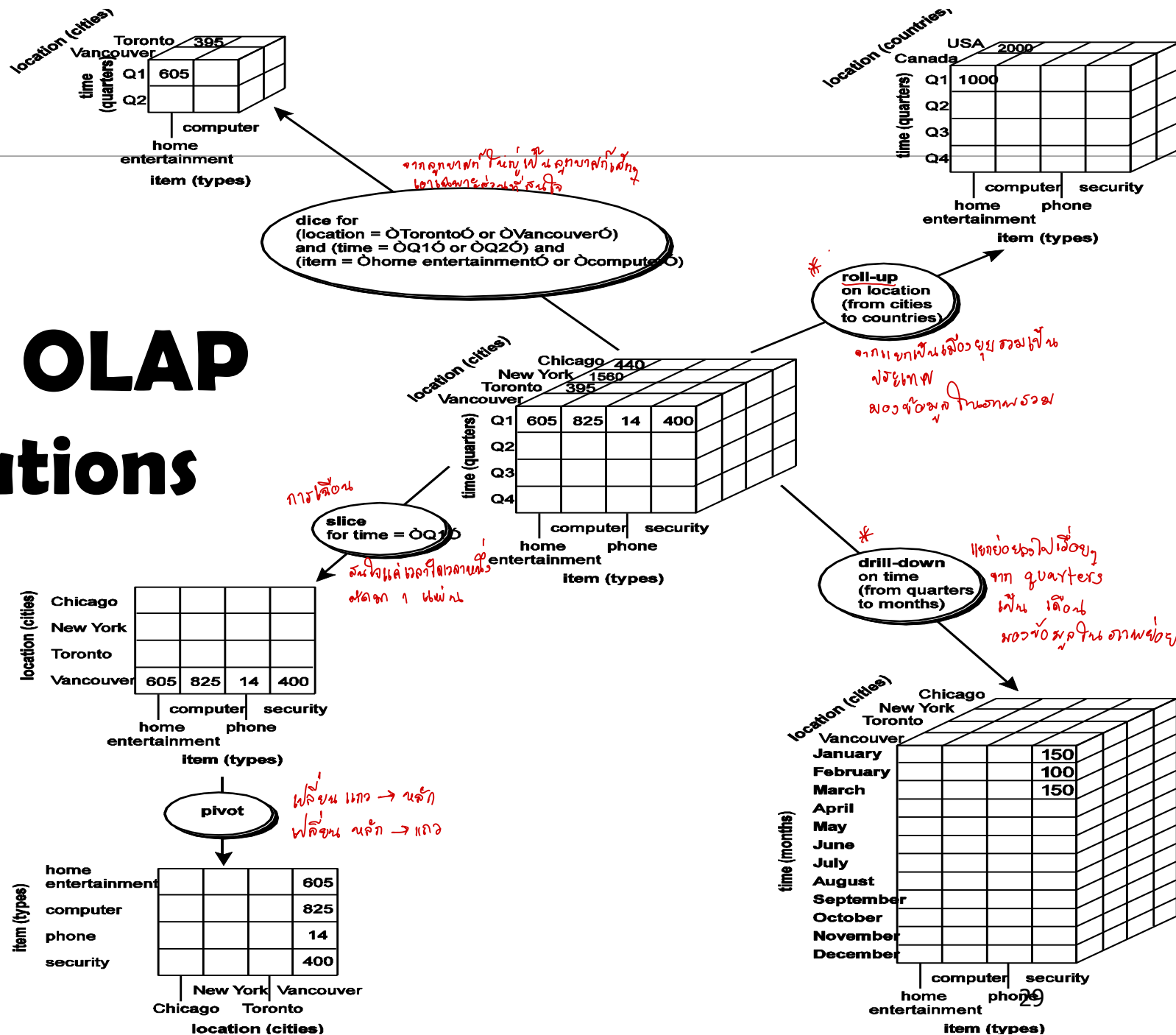


# Typical OLAP Operations

---

- ❑ **Roll up (drill-up):** summarize data
  - ❑ *by climbing up hierarchy or by dimension reduction*
- ❑ **Drill down (roll down):** reverse of roll-up
  - ❑ *from higher level summary to lower level summary or detailed data, or introducing new dimensions*
- ❑ **Slice and dice:** *project and select*
- ❑ **Pivot (rotate):**
  - ❑ *reorient the cube, visualization, 3D to series of 2D planes*
- ❑ **Other operations**
  - ❑ **Drill across:** *involving (across) more than one fact table*
  - ❑ **Drill through:** *through the bottom level of the cube to its back-end relational tables (using SQL)*

# Typical OLAP Operations



# A Star-Net Query Model

ข้อมูล data ในหน่วยต่างๆ  
ที่เชื่อมโยง

9 dimension

hierarchical

