

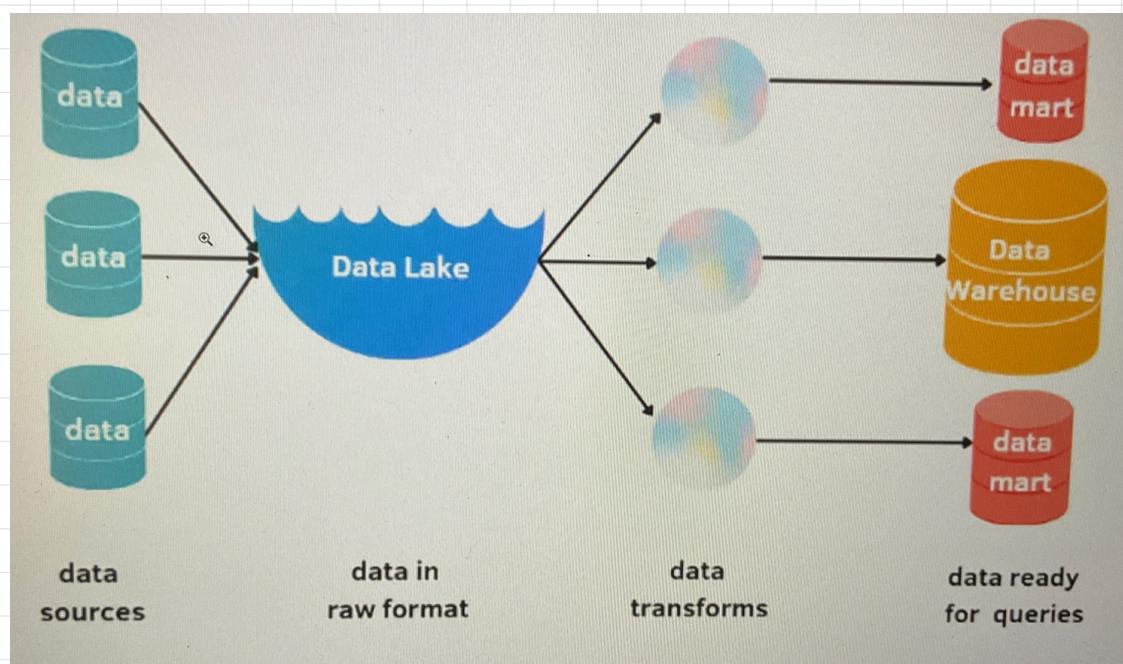


# **CS 412 Intro. to Data Mining**

## **Chapter 4. Data Warehousing and On-line Analytical Processing**

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

# Why is Business Intelligence useful for a Data Scientist?



ក្នុងវិមានអរកម្មអនុសាស្ត្រទេរច្បាប់ មែនត្រូវបាន  
រាយការណីជាន់ Data Lake សង្ឃ

គិតជាអក្សរខ្លួន និងគោរពទាំងមីនា

## Table of Content:

Structured  $\Rightarrow$  Data ຖែរការណ៍ Warehouse

Unstructured  $\Rightarrow$  Data ក្នុងអាជីវកម្ម និងការ  
រំលែករាយការណ៍សម្រាប់បង្កើតឡើង  
ដោយ សម្រាប់ គ្រប់គ្រងការងារ

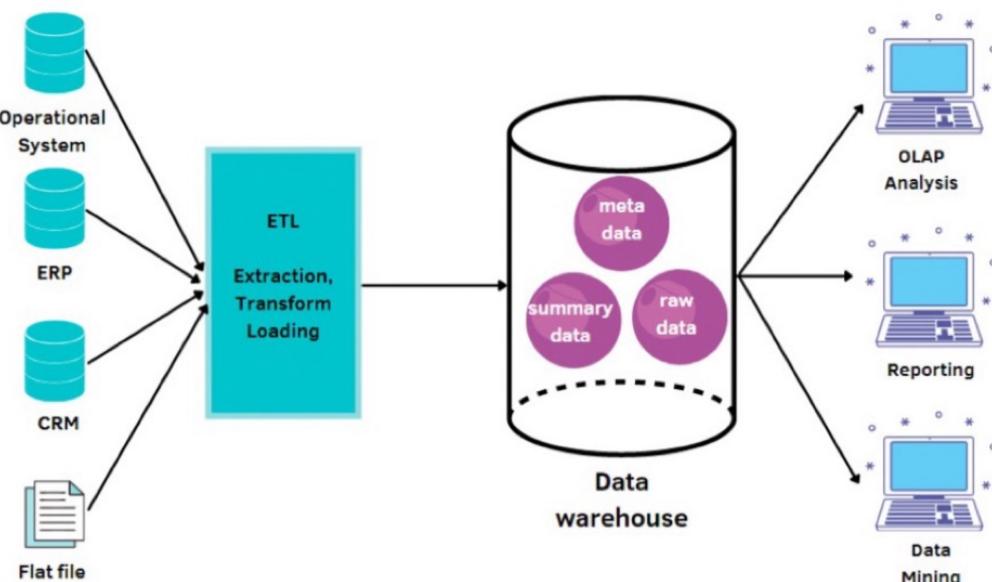
1. Data Warehouse
  2. Data Mart
  3. OLTP vs OLAP
  4. ETL
  5. Star vs Snowflake schemas
  6. Data Lake
  7. From ETL to ELT
  8. Batch vs Stream Processing

3. OLTP vs OLAP → integration in Data warehouse

→ ឧគ្គការនូវការទាំងរបស់ខ្លួន និងខ្លួន  
(ពេល, នូល, ចិត្តទៅលើវិវាទ Database)  
ដើម្បីការកិច្ចការណា Data

FIELD NAME	DATA TYPE	DESCRIPTION			
search_id	STRING	search id			
search_timestamp	TIMESTAMP	วันและเวลาดัชนีของการค้นหา			
user_agent	STRING	ประเภทของบราวเซอร์ที่ทำการค้นหา			
q	STRING	คำที่ผู้ใช้ใส่ในช่องค้นหา			
user_id	INTEGER	DUMMY user_id			
session_id	STRING	เลข session (จะเป็นเลขเดิมถ้าผู้ใช้ค้นเดิม ทำการค้นใหม่ภายใน 15 นาที)			
number_of_result	INTEGER	จำนวนผลการค้นหา			
lat	FLOAT	ละติจูด ของผู้ใช้			
long	FLOAT	ลองจิจูด ของผู้ใช้			
10M records	1 month		มิถุนายน, 2018		
14	Search type		June		
15	Keyword search when the field name "q" has any value		# of rows	%	
16	Nearby search when the field name "q" is empty or null		688441	6.88%	
17			9316789	93.12%	
18			sum	100005230	100.00%
19			size	1.79 GB	
20			compressed size	652.6 MB	
21					
22					
23					
24					
25					
26					
27					
28					

## 1. Data Warehouse



Data Warehouse architecture. Illustration by author.

The data warehouses store three types of data:

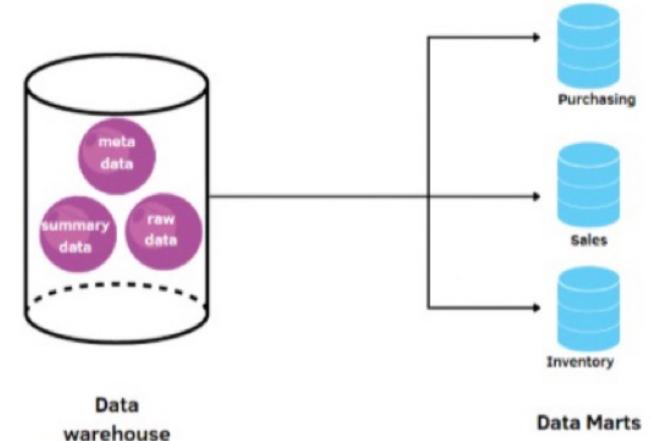
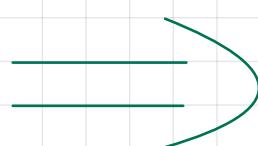
★ Meta data ⇒ ទិន្នន័យតាមចំណាំសម្រាប់បង្កើតការងារ និងការបង្ហាញ

★ Summary data  $\Rightarrow$  សិក្សានៅក្នុងរាជធានីភ្នំពេញ ត្រូវបានបង្កើតឡើងដើម្បី  
គ្រប់គ្រង និងអនុវត្តន៍ការងារសំគាល់សិក្សា

★ Raw data ⇒ ຖື່ມຂອບໃຈຂໍ້ມູນໄສຫຼວດກໍ່າລົງທຶນ  
ຕຽບກຳລັງຂອງຂໍ້ມູນທີ່ມີຄວາມຍຸດ

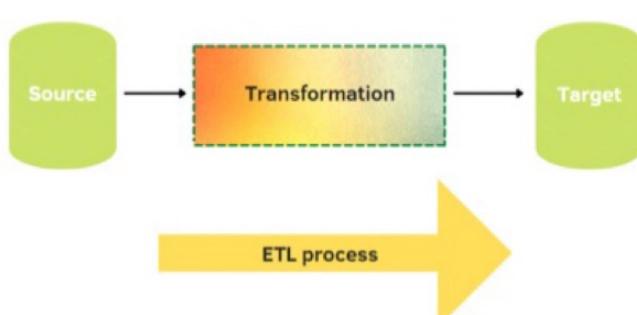
## 2. Data Marts

→ តើវាត្រូវបានចាប់សំរាប់ឡើងទៅលើ Sales, Inventory



## Data warehouse vs Data Marts. Illustration by author.

### 3. ETL



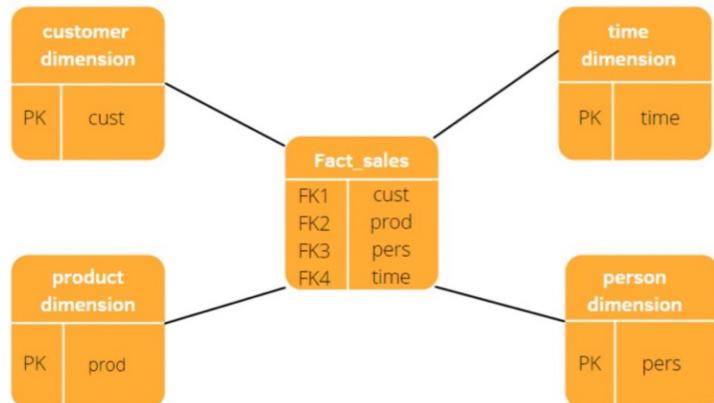
## ETL process. Illustration by author

→ ទីនេះការបង្កើតការងាររបស់ខ្លួន និងការរួមរាល់របស់ខ្លួន ដើម្បីស្ថានភាពនៃការអភិវឌ្ឍន៍

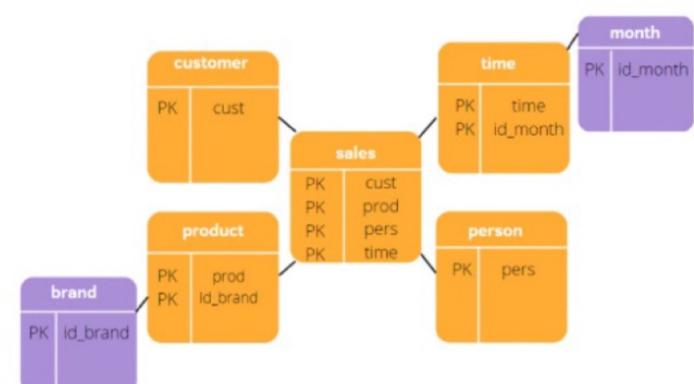
Extract : វិនិច្ឆ័យការពេក្តីរួម្ងាស់ទាមអត្ថបទ ផែន ឯក 'CSV, JSON, XML'

 **Transform:** ការផ្សេងៗបច្ចេកវិទ្យាបានរៀបចំឡើងដោយការសម្រេចក្នុងការរាយការណ៍ និងការរៀបចំការសំរាប់ នាយករដ្ឋមន្ត្រី ក្រោមការគាំទ្រការអភិវឌ្ឍន៍ការរាយការណ៍

## 4. Star vs Snowflake Schemas



Star Schema. Illustration by author



Snowflake Schema. Illustration by author

Data warehouse មានការចូលរួមនៃ Schemas ទាំងអស់ Schemas មិនបានមែនការចូលរួមនៃការគ្រប់គ្រងទៅការបង្កើតរាយការណ៍។ Star Schema នឹងធ្វើឡើងដោយសារតាមការចូលរួមនៃការបង្កើតរាយការណ៍ទៅការបង្កើតរាយការណ៍។

★ Star Schema ⇒ នឹងគ្រប់គ្រងគ្រប់គ្រងតាមការចូលរួមនៃ fact table និង Dimension tables ឬផ្សេងៗ

★ Snowflake Schema ⇒ នឹងគ្រប់គ្រងគ្រប់គ្រងតាមការចូលរួមនៃ fact table និង Dimension tables ឬផ្សេងៗ

## 5. Data Lake

DATA WAREHOUSE	DATA LAKE
structured, processed	structured / semi-structured/ unstructured, raw
schema-on-write	schema-on-read
expensive for large data volumes	designed for low-cost storage
less agile, fixed configuration	highly agile, possible updates
mature	maturing
business professionals	data scientists

Comparison of Data Warehouse and Data Lake. Illustration by author.

★ វិភាគនៃការចូលរួមនៃយុទ្ធសាស្ត្រ Data warehouse និងវឌ្ឍនភាពក្នុងការចូលរួមនៃយុទ្ធសាស្ត្រ។ Data warehouse នឹងបានគ្រប់គ្រងគ្រប់គ្រងតាមការចូលរួមនៃការបង្កើតរាយការណ៍។ Data lake នឹងបានគ្រប់គ្រងគ្រប់គ្រងតាមការចូលរួមនៃការបង្កើតរាយការណ៍។ Star Schema នឹងបានគ្រប់គ្រងតាមការចូលរួមនៃការបង្កើតរាយការណ៍។ Snowflake Schema នឹងបានគ្រប់គ្រងតាមការចូលរួមនៃការបង្កើតរាយការណ៍។

## 6. From ETL to ELT

ETL	ELT
data is transformed and then transferred to Data Warehouse DB	Data remains in the DB of Data Warehouse
At early stages, easier to implement	To implement ELT process deep knowledge of tools and expert skills are needed.
Supports relational and structured data.	Supports structured, unstructured data sources.
Does not support Data Lake	Allows use of Data Lake
High costs for small and medium businesses.	Low entry costs using online Software as a Service Platforms.
Complexity increase with the additional amount of data in the dataset.	Power of the target platform can process significant amount of data quickly.
The process is used for over two decades.	Relatively new concept and complex to implement.

Comparison of ETL and ELT. Illustration by author

★ ETL នឹង Extract-transform-load

នឹងការបញ្ចូនពីរឿងទិន្នន័យ និងការចូលរួមនៃយុទ្ធសាស្ត្រ និងការចូលរួមនៃការបង្កើតរាយការណ៍។

★ ELT នឹង Extract (ឈ្លឹនពីរឿង) load (ដំឡើង) transform (រៀបចំពីរឿង)

គឺការបញ្ចូនពីរឿងទិន្នន័យ និងការចូលរួមនៃយុទ្ធសាស្ត្រ និងការចូលរួមនៃការបង្កើតរាយការណ៍។

# 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



Model → Data Warehouse →

1. Data Cube → measures Dimension  
of 3 types → Star Schema,  
Snowflake Schema,  
Fact Constellation
2. OLAP (on-line Analytical Processing)

# What is a Data Warehouse?

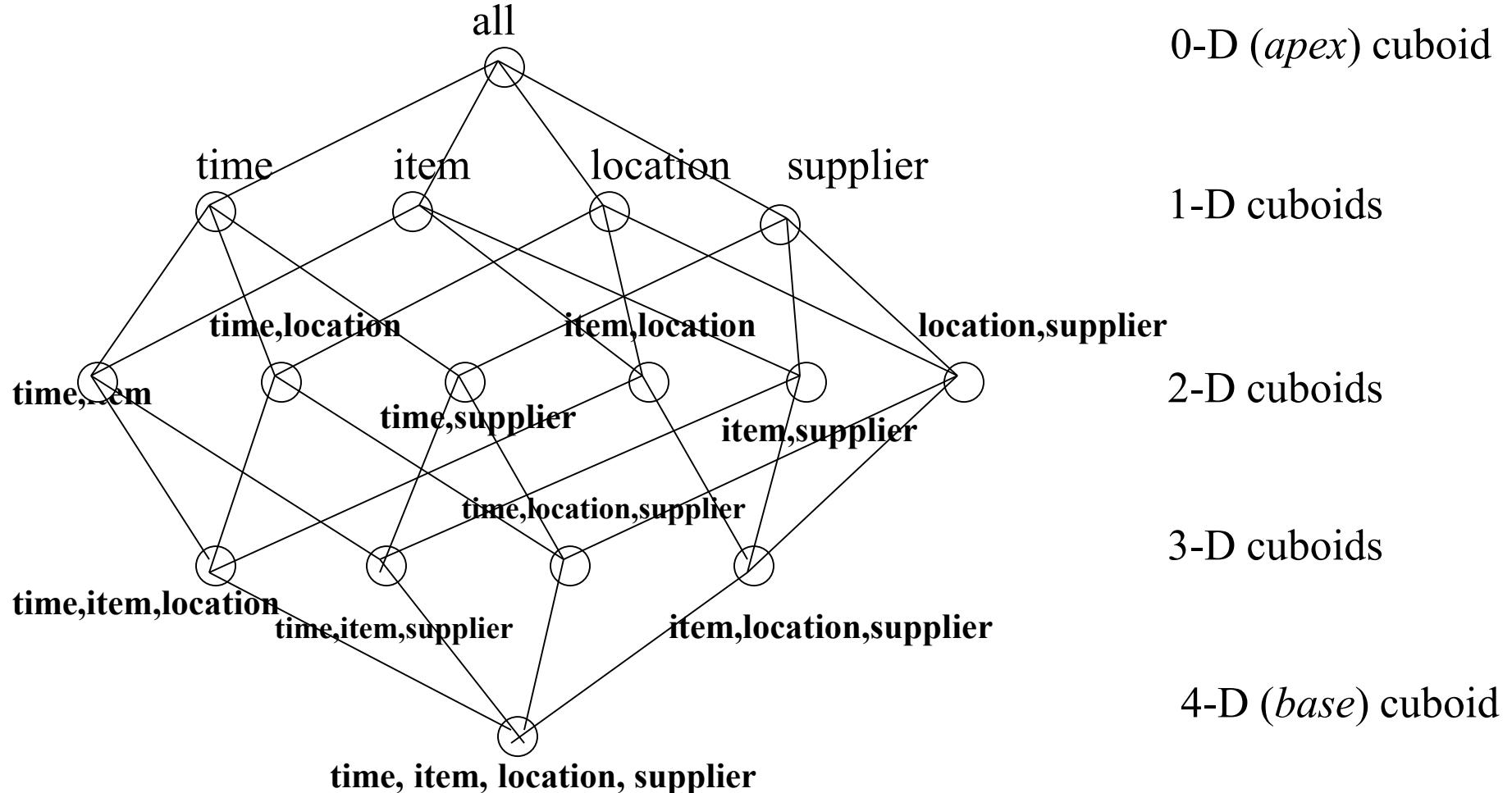
- ❑ Defined in many different ways, but not rigorously
  - ❑ A decision support database that is maintained **separately** from the organization's operational database
  - ❑ Support **information processing** by providing a solid platform of consolidated, historical data for analysis
  - ❑ “A data warehouse 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

→ ပုဂ္ဂန် Data ကိုလဲ အဆင့်မြင်စွာ ပြန် Bath

# 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)
  - Fact table contains measures (such as dollars\_sold) and keys to each of the related dimension tables
- **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**.

# Data Cube: A Lattice of Cuboids



## 2D

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)			
	home	entertainment	computer	phone
Q1	605		825	14
Q2	680		952	31
Q3	812		1023	30
Q4	927		1038	38

⇒ 2 维

↳ 1 维

- មួយអាមេរិកមេគា time (quarter)

- មួយអាមេរិកក្នុងនៅទី item (type)

## 3D

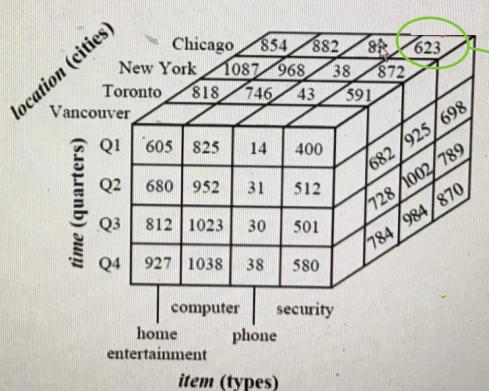
Table 4.3: A 3-D view of sales data for *AllElectronics*, according to the dimensions *time*, *item*, and *location*. The measure displayed is *dollars\_sold* (in thousands).

<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>	
home	home	home	home	
time	ent.	comp.	phone	
Q1	854	882	89	623
Q2	943	890	64	698
Q3	1032	924	59	789
Q4	1129	992	63	870
	1087	968	38	872
	894	769	52	682
	940	795	58	728
	978	864	59	784
	818	746	43	591
	892	952	31	512
	940	1023	30	501
	927	1038	38	580

⇒ 3 维

↳ ចាប់ពី 1 លើ 4 ដែលត្រូវបានបង្ហាញឡើង  
នៅក្នុងការបង្ហាញទី 3 នេះ

## 3D

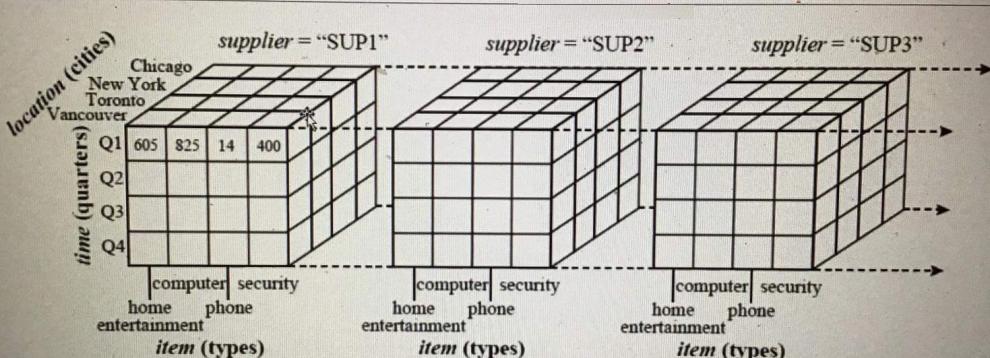


⇒ 3 维

Q1 too location

Chicago = 123

## 4D



⇒ 4 维

ចាប់ពី 1 លើ 4 ដែលត្រូវបានបង្ហាញឡើង  
នៅក្នុងការបង្ហាញទី 4 នេះ

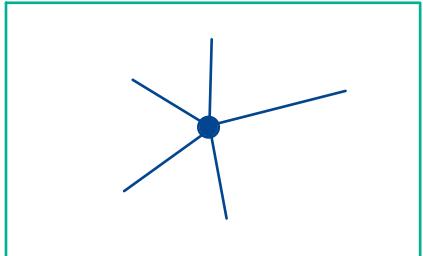
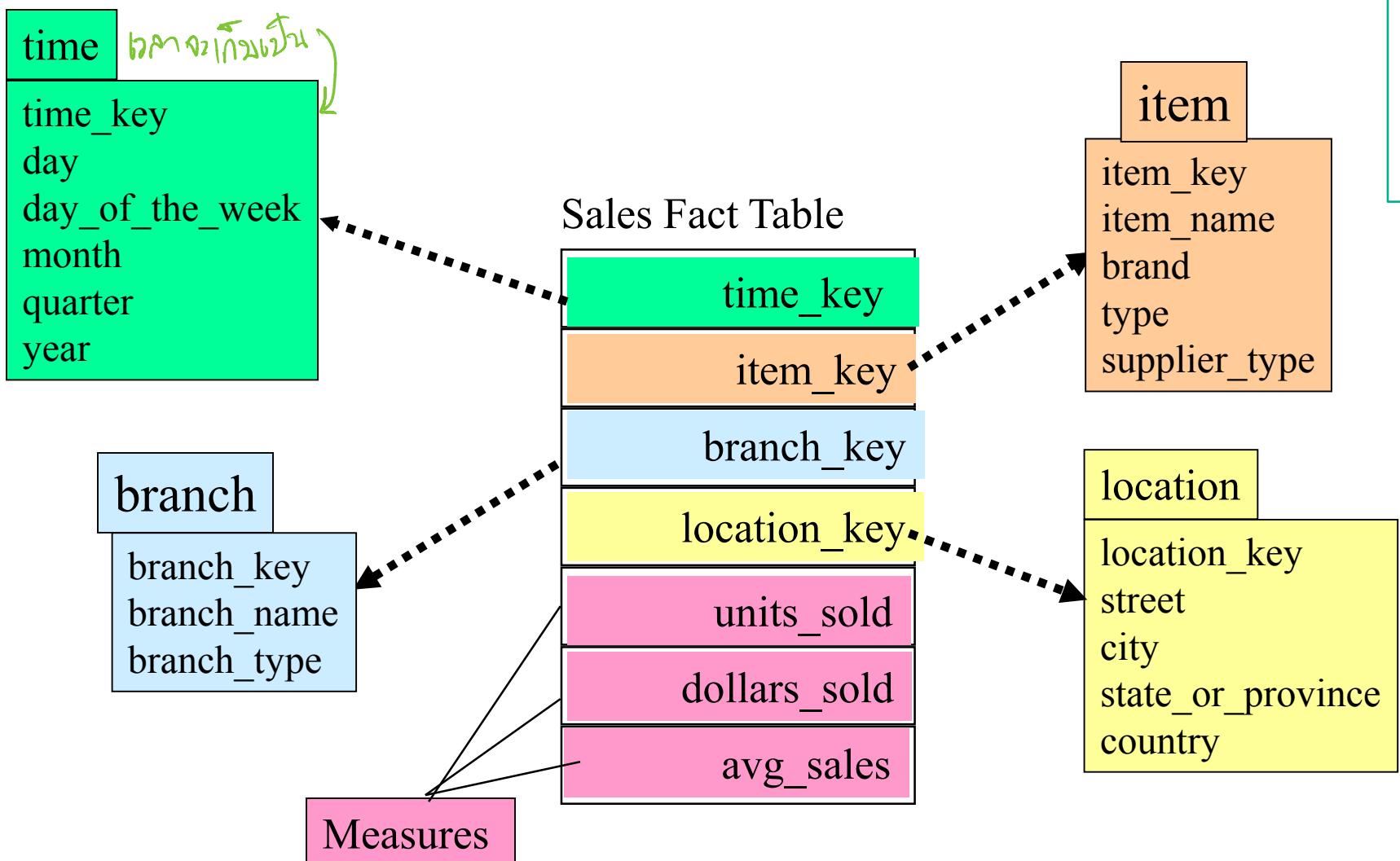
Figure 4.4: A 4-D data cube representation of sales data, according to the dimensions *time*, *item*, *location*, and *supplier*. The measure displayed is *dollars\_sold* (in thousands). For improved readability, only some of the cube values are shown.

# Conceptual Modeling of Data Warehouses

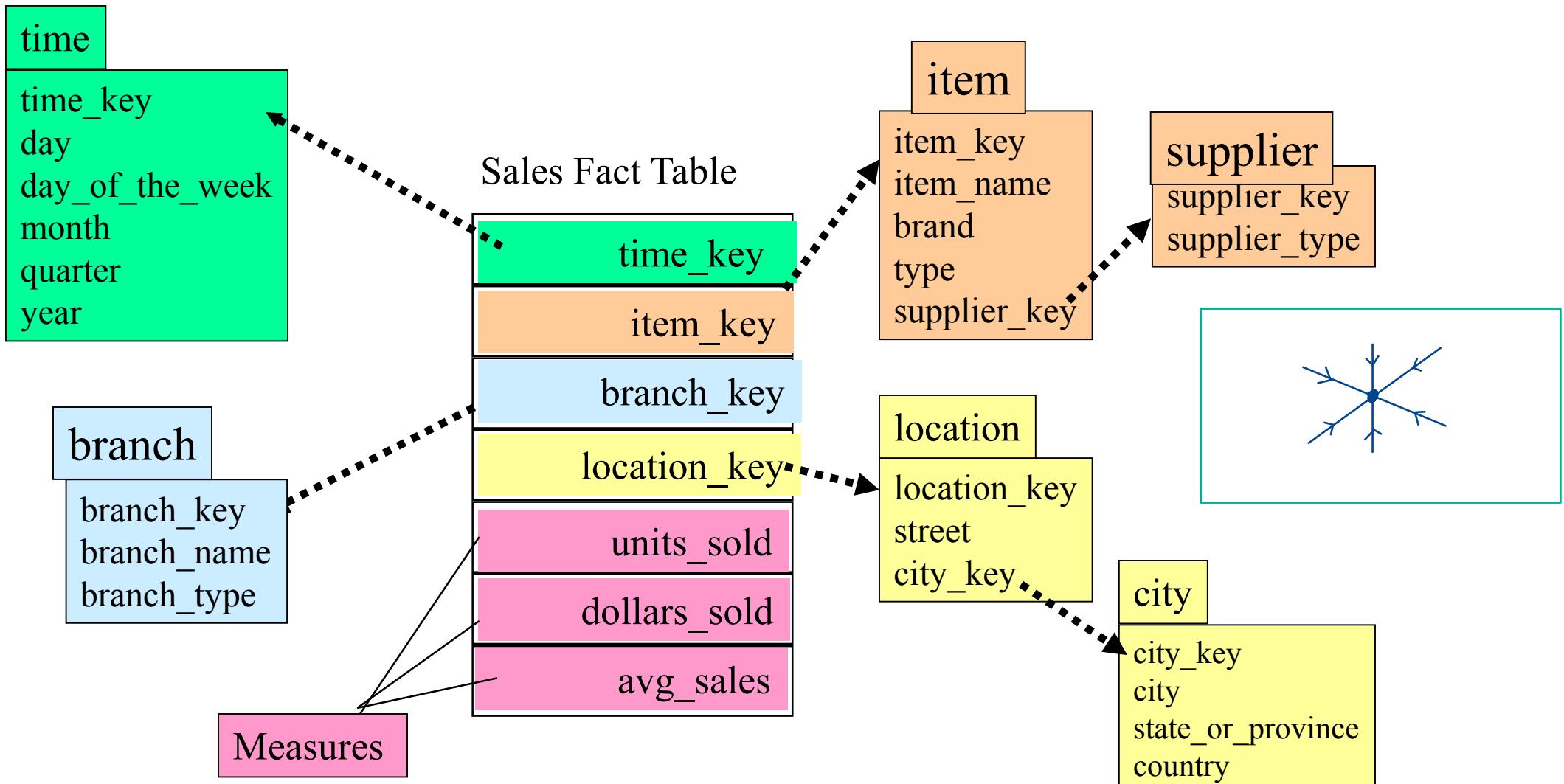
---

- 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



# Snowflake Schema: An Example



# Fact Constellation: An Example

*Dimension Table*

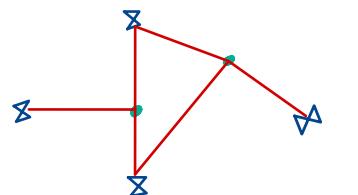
time
time_key
day
day_of_the_week
month
quarter
year

Sales Fact Table

time_key
item_key
branch_key
location_key
units_sold
dollars_sold
avg_sales

branch
branch_key
branch_name
branch_type

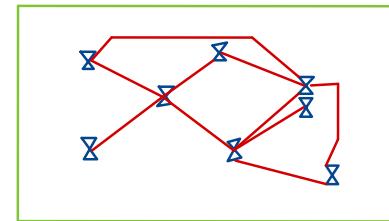
Measures



item
item_key
item_name
brand
type
supplier_type

Shipping Fact Table

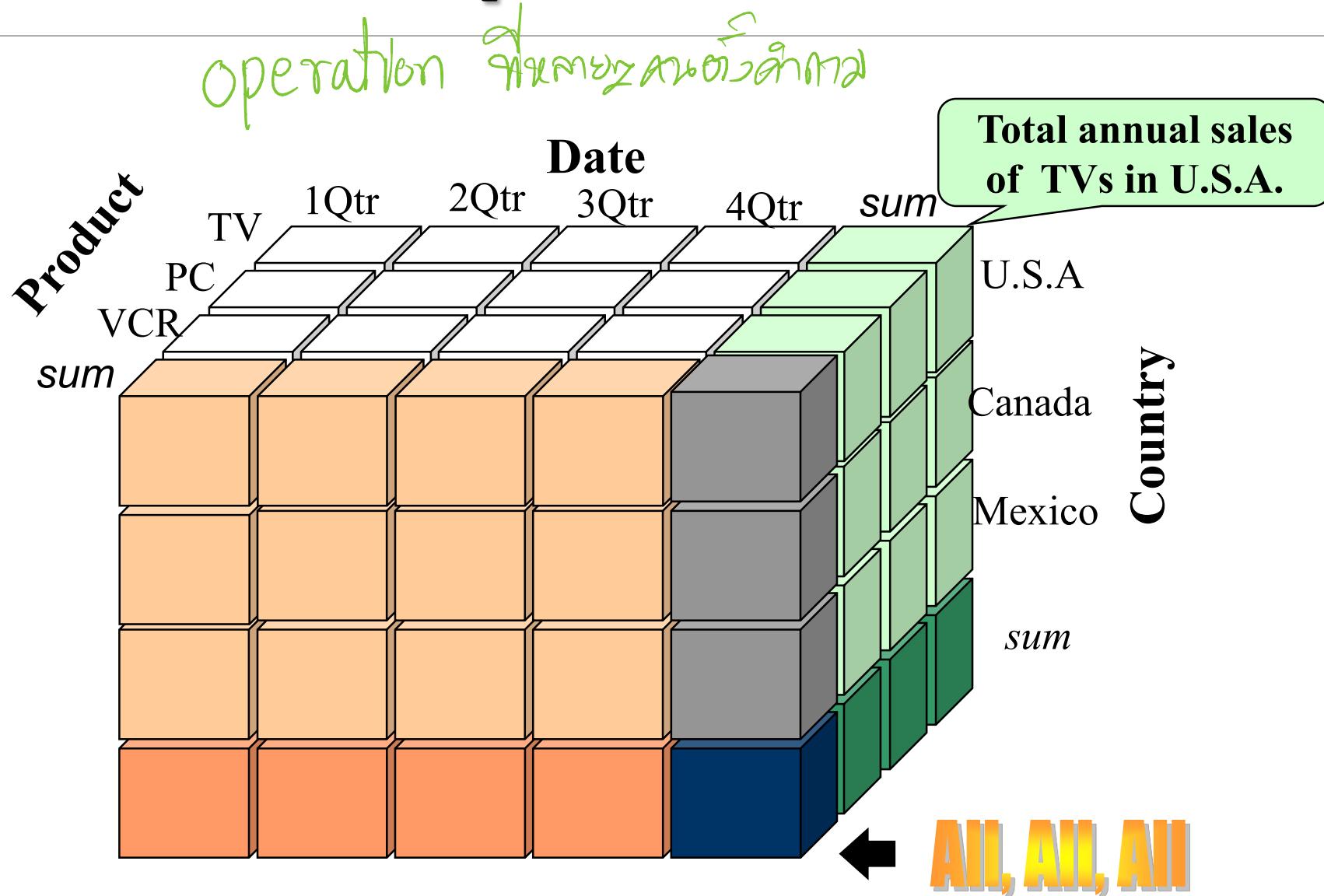
time_key
item_key
shipper_key
from_location
to_location
dollars_cost
units_shipped



location
location_key
street
city
province_or_state
country

shipper
shipper_key
shipper_name
location_key
shipper_type

# A Sample Data 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

