

~ Chapter 4 ~

why is Business Intelligence useful for a Data Scientist ?

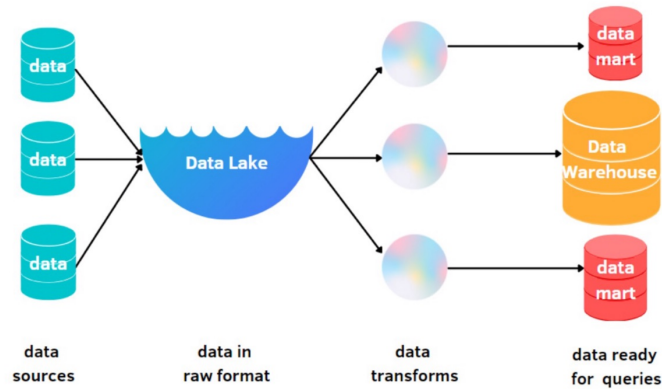


Illustration by author

structured → Data ที่ถูกเก็บไว้ใน Warehouse

Unstructured → Data ที่คอมพิวเตอร์ไม่สามารถเข้าใจได้ทันที เช่น รูป เสียง ข้อความ (คอมเมนต์)

ศัพท์ที่ควรรู้

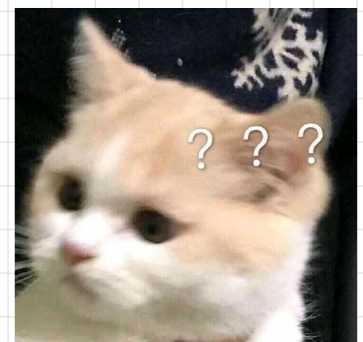
Table of Content:

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

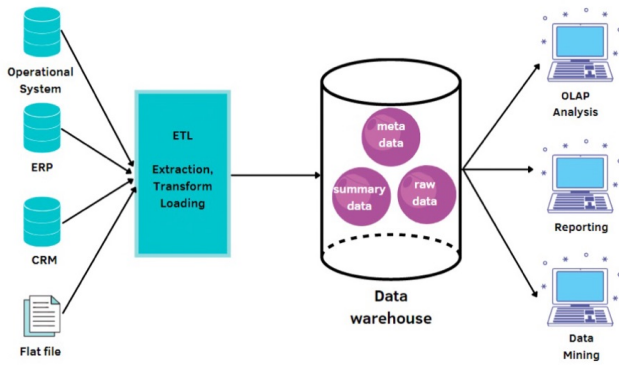
OLTP vs OLAP

OLTP → แอปพลิเคชันในการจัดการข้อมูล
เนื้อหาคือ (เพิ่ม, ลบ, อัปเดตข้อมูล ในฐาน
data base) เป็นกรณีการใช้งาน Data

OLAP → เกิดตอนที่ Data warehouse
เก็บข้อมูลไปวิเคราะห์ เพื่อดูแนวโน้มที่เราจัดการมาแล้ว



1. Data Warehouse



Data Warehouse architecture. Illustration by author.

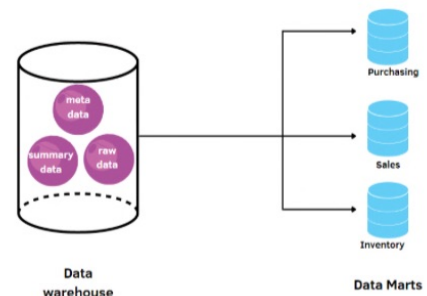
The data warehouses store three types of data :

🌸 **Metadata** → ข้อมูลที่อธิบายข้อมูลเกี่ยวกับตัวข้อมูลและรายละเอียดเพื่อให้เข้าใจตัวข้อมูลที่จัดเก็บ

🌸 **Summary data** → ข้อมูลที่รวบรวม / สรุปที่สร้างโดยผู้จัดการของตัวข้อมูล ซึ่งแบ่งประเภทกิจกรรมที่สัมพันธ์ข้อมูล

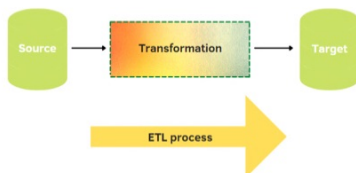
🌸 **Raw data** → ข้อมูลดิบเป็นข้อมูลที่ไม่ได้ผ่านการประมวลผล

2. Data Marts



Data warehouse vs Data Marts. Illustration by author.

3. ETL



ETL process. Illustration by author

ETL : เป็นกระบวนการรวบรวมและประมวลผลที่ช่วยให้สามารถสำรวจ
คลังข้อมูลได้

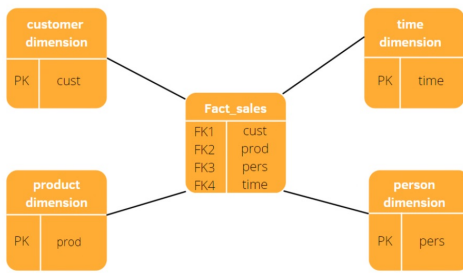
🥕 **Extract** : รับข้อมูลดิบจากแหล่งต่างๆ เช่น ไฟล์ CSV, JSON, XML

🥕 **Transform** : แปลงข้อมูลให้อยู่ในรูปแบบที่เป็นประโยชน์สำหรับวิเคราะห์
รวมกับที่ซ้ำ ตรวจสอบ การรวม และแปลงรายการที่ซ้ำกัน

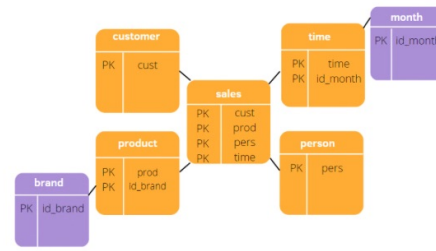
🥕 **Load** : ข้อมูลจะถูกบันทึกไว้ในปลายทางสุดท้ายนั้นคือคลังข้อมูล เกิดขึ้นได้
ทุกนาที ชั่วโมง วัน สัปดาห์ หรือในเวลาที่กำหนดทำใน การวิเคราะห์
ก็จะมีแผนกที่รับผิดชอบ



4. Star vs Snowflake Schemas



Star Schema. Illustration by author



Snowflake Schema. Illustration by author

Data warehouse มีทั้งข้อมูลโดยมี Schemas หลายมิติ Schemas มีประโยชน์ในการจัดเก็บข้อมูลจำนวนมากเพื่อวัตถุประสงค์ในการวิเคราะห์ มี Schemas ที่ใช้กันหลายรูปแบบ

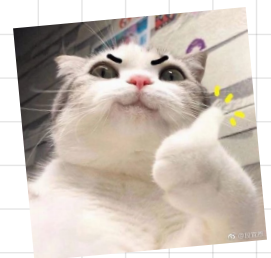
- **Star Schema** → มีโครงสร้างคล้ายดาวที่ตรงกลางมี fact table โดยมี Dimension tables อยู่รอบๆ
- **Snowflake Schema** → มีโครงสร้างแบบเกล็ดหิมะ จัดตั้งช่วยเพิ่มขนาดและใช้พื้นที่จัดเก็บน้อยลง

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.

- **Data Lake** มีความซับซ้อนน้อยกว่า Data warehouse เนื่องจากเก็บข้อมูลทั้งข้อมูลดิบ ที่ยังไม่มีการสร้าง และยังไม่มีการสร้าง อนุญาตให้เก็บข้อมูลโดยไม่มีการ Schemas อนุญาตให้ทั้งที่วิเคราะห์ข้อมูลในอดีต ปัจจุบัน และอนาคต เนื่องจากข้อมูลจะไม่ถูกลบ



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** หรือ Extract-transform-load เป็นกระบวนการดึงข้อมูลในหลายๆโปรแกรมแล้วมาแปลงและโหลดข้อมูลเข้าใน Data Warehouse การดึงข้อมูลเข้าใน Data Warehouse นั้นจะดึงข้อมูลจากหลายๆโปรแกรมแล้วมาแปลงข้อมูลก่อน ทำให้ผลลัพธ์ที่ออกมาจะมีความซับซ้อนมากขึ้น
- **ELT** หรือ Extract (สกัดข้อมูล) Load (ดึงข้อมูล) Transform (แปลงรูปแบบข้อมูล) คือการดึงข้อมูลจากแหล่งแล้ว รวบรวมเข้าใน Data Warehouse โดยในกระบวนการจะมี Transform ข้อมูลตามที่ต้องการ

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

เป็น Data
ที่ไม่เคลื่อนไหว
เป็น Bath

→ เราต้องนำข้อมูลจากหลายๆ DW มาเพื่ออะไร? สรรวามาเพื่อมองคำถาม
สร้างเพื่อ support การตัดสินใจของผู้บริหาร

- ❑ Data warehousing:
 - ❑ The process of constructing and using data warehouses

From Tables and Spreadsheets to Data Cubes

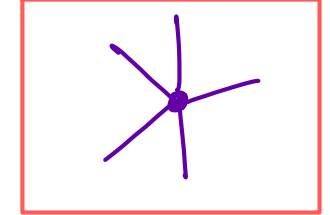
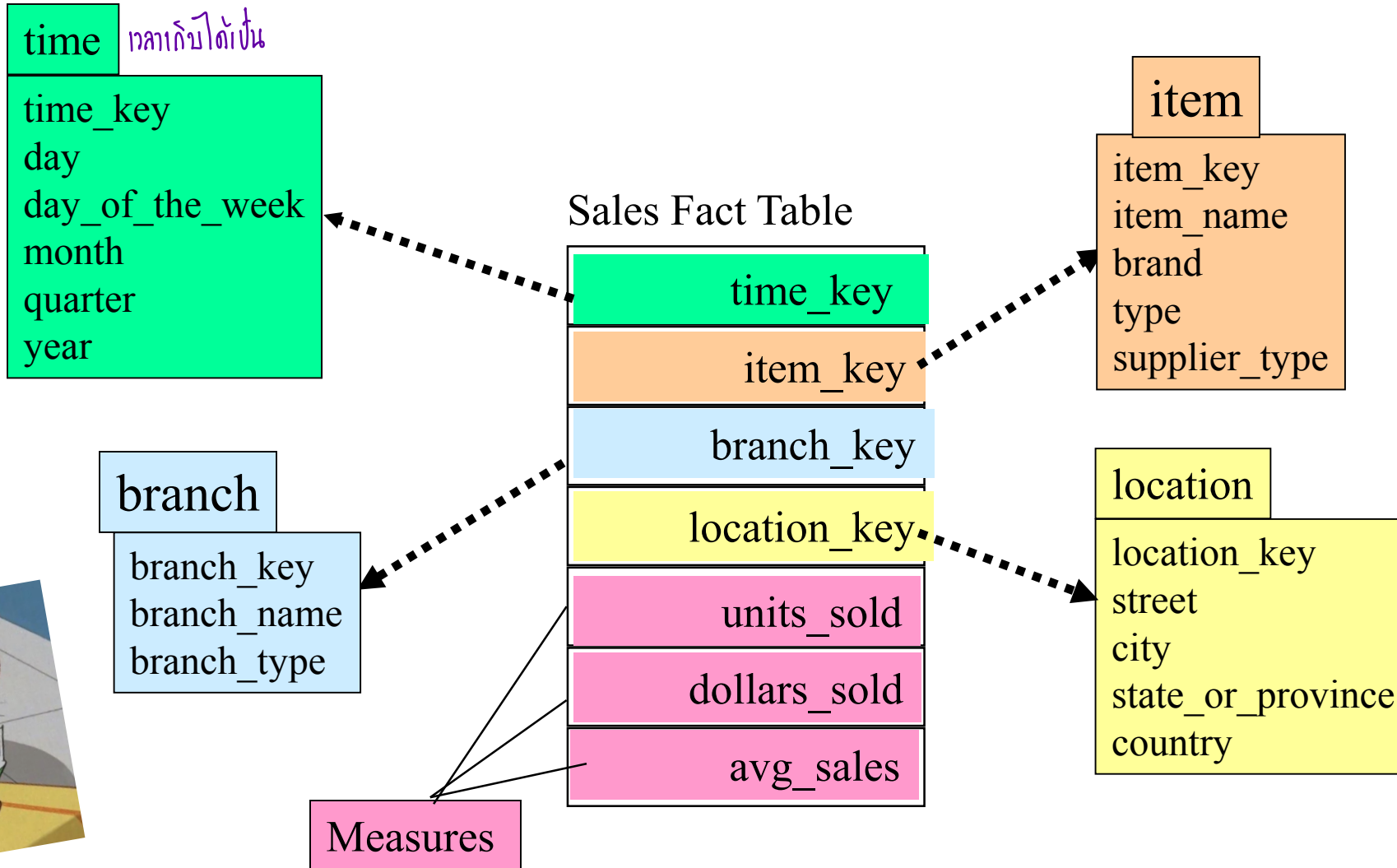
- ❑ A **data warehouse** is based on a multidimensional data model which views data in the form of a data cube
- ❑ A data cube, such as sales, allows data to be modeled and viewed in multiple dimensions
 - ❑ Dimension tables, such as item (item_name, brand, type), or time(day, week, month, quarter, year)
 - ❑ Fact table ^{ตารางที่เก็บค่าตัวเลข} 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**.

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

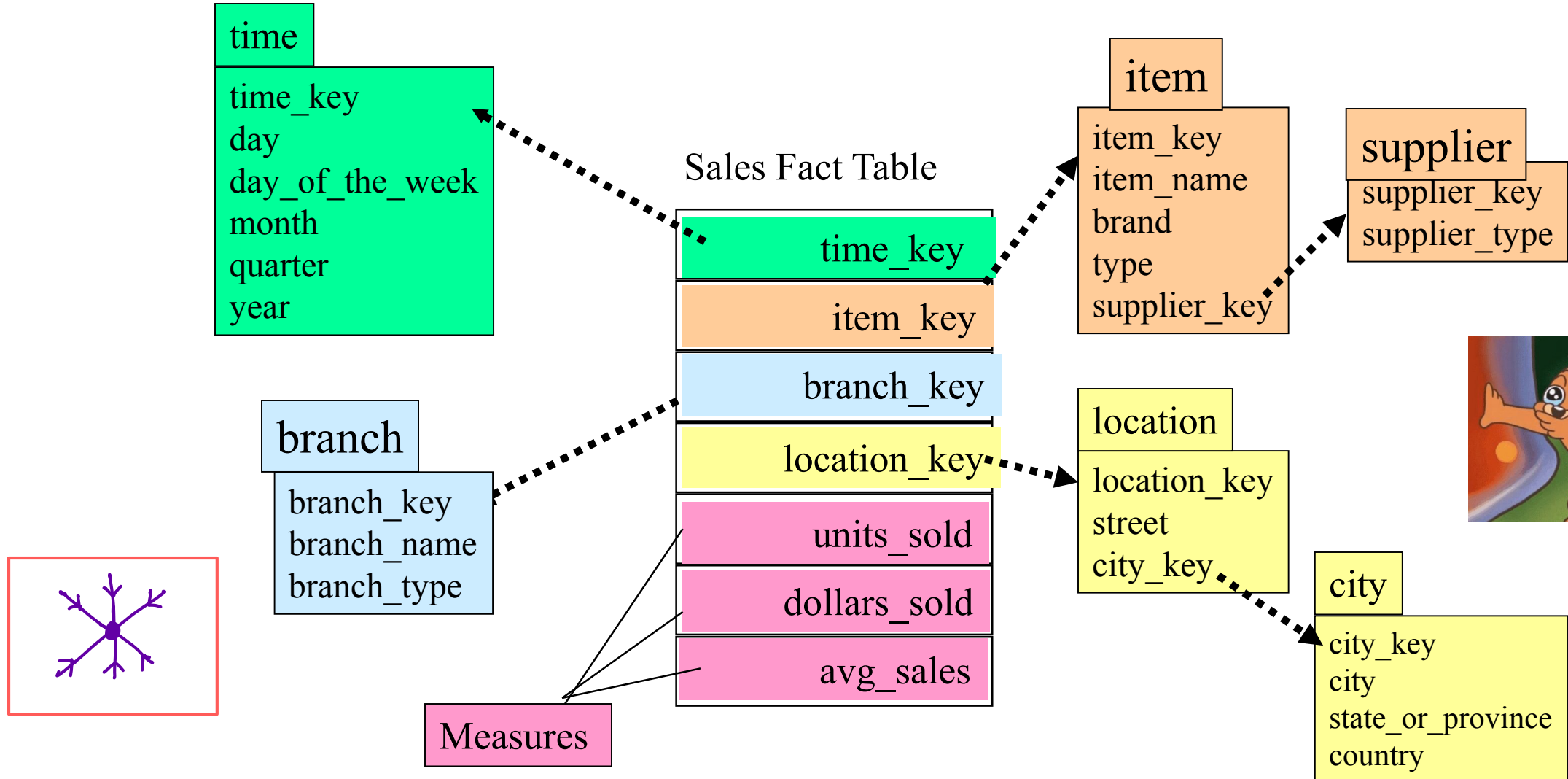
↗ มีจุดศูนย์กลาง ทุกอย่างมาเชื่อมจากศูนย์กลาง

Star Schema: An Example

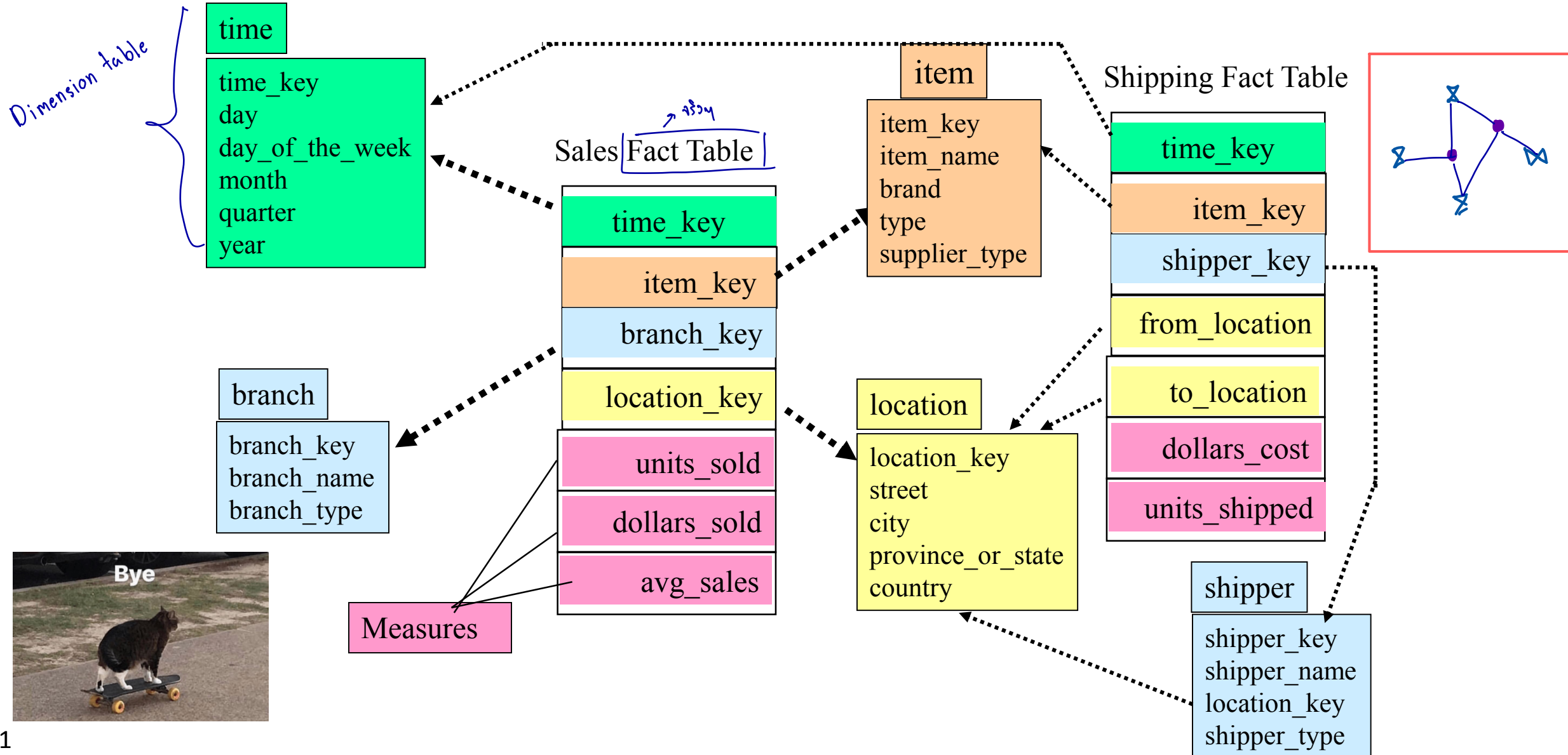


👉 ฝึกฝน: คล้ายๆ เกิดนี้ละ !!

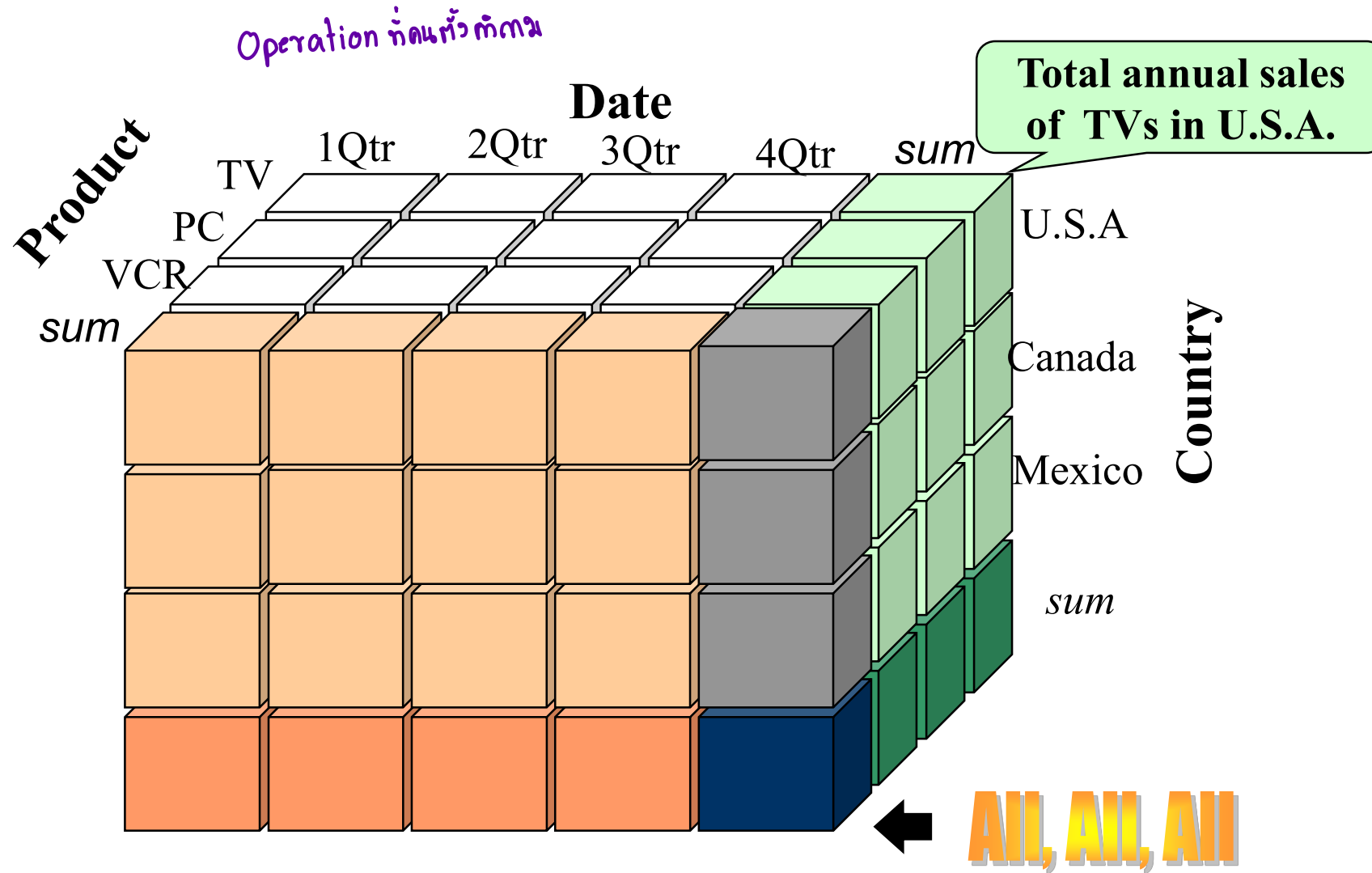
Snowflake Schema: An Example



Fact Constellation: An Example



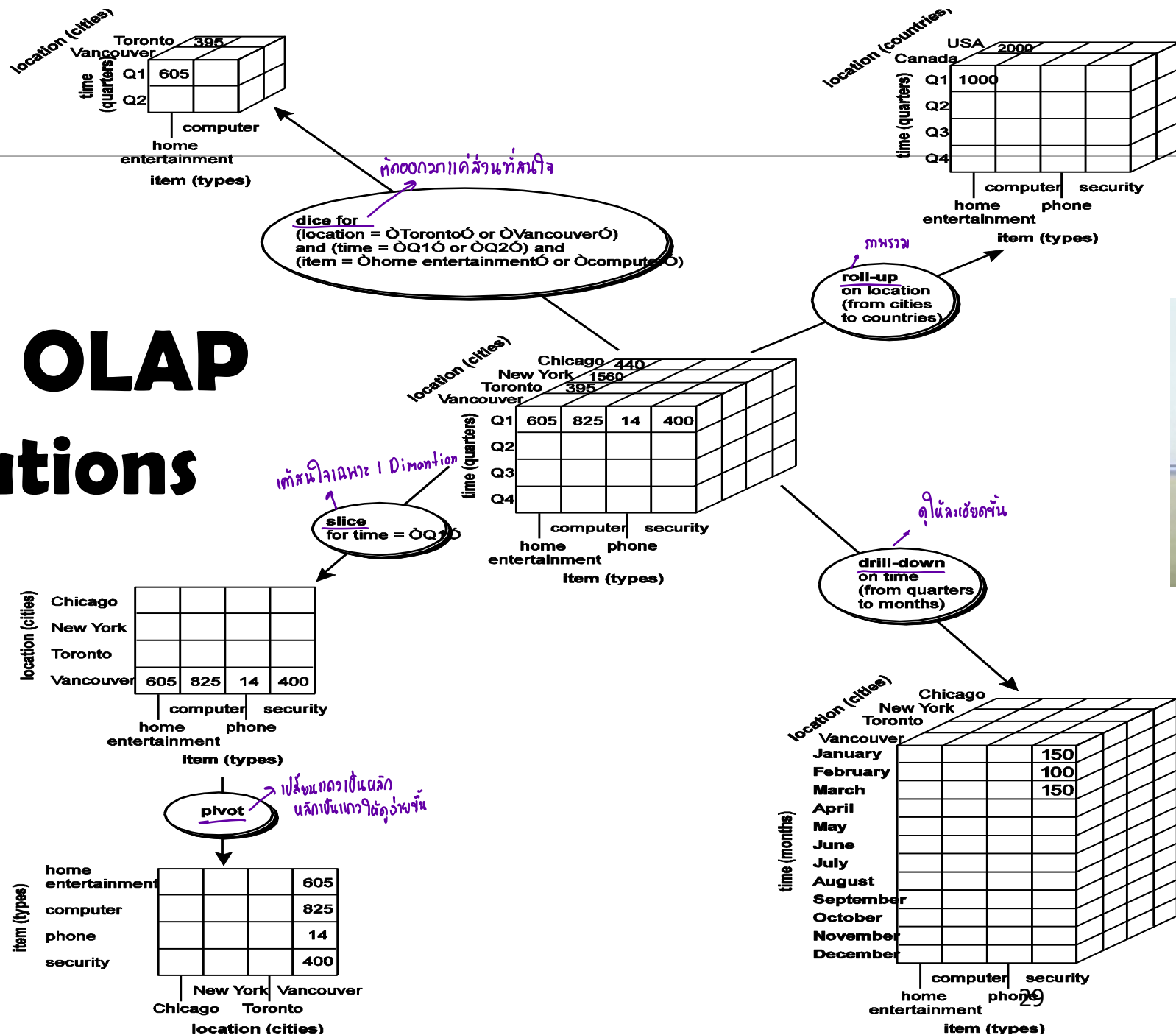
A Sample Data Cube



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



Roll-up



the end...