

# Ch3 Data Warehouse

---

**Dr. Bernard Chen Ph.D.**

University of Central Arkansas

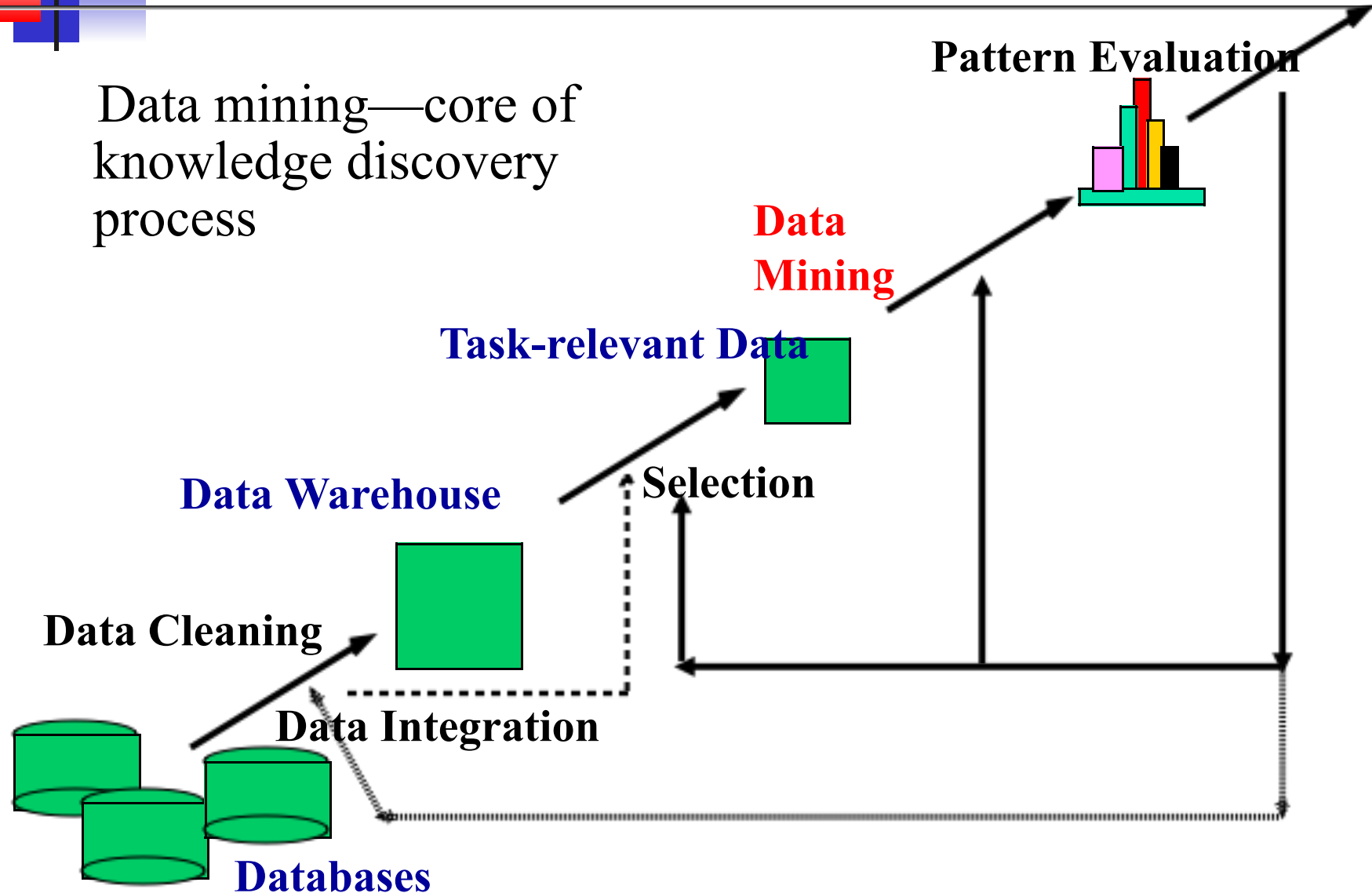
Fall 2010

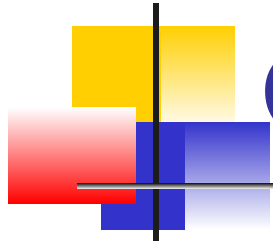


# Knowledge Discovery (KDD) Process

Knowledge

Data mining—core of knowledge discovery process

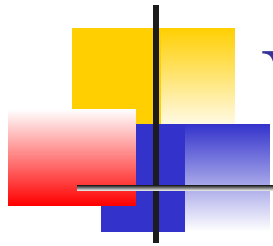




# Outline

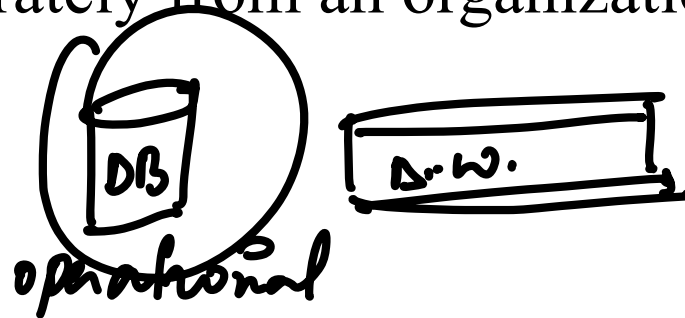
---

- What is Data Warehouse?
- Data Warehouse: A multidimensional data model
- Data Warehouse Architecture



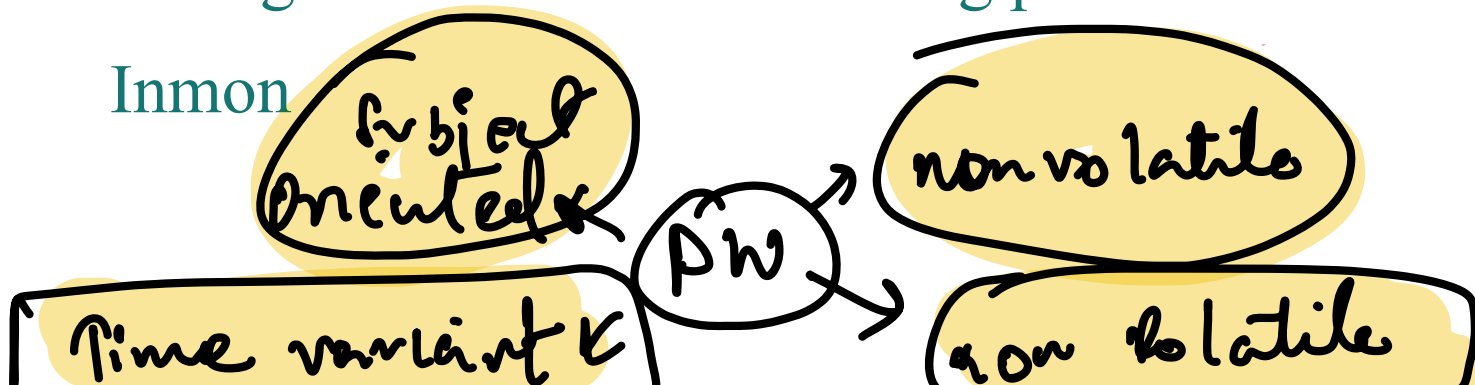
# What is Data Warehouse?

- Loosely speaking, a data warehouse refers to a database that is maintained separately from an organization's operational database



- Officially speaking:
- “A data warehouse is a subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process.”—William. H.

Inmon



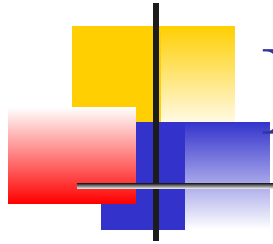


# Data Warehouse—Subject-Oriented

---

- Organized around major subjects, such as **customer, product, sales** *CRM data, order details, Billing details, SLA/KPI*
- 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

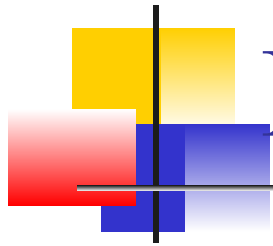
Multiple data sources.



# 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
  - 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)

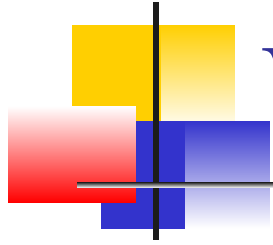


• do not require transactional processing & concurrency control mechanisms  
• update of data does not occur  
initial load + access of data

# Data Warehouse—Nonvolatile

- A physically separate store of data transformed from the operational environment
- 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*

Subject oriented, Time Variant, Nonvolatile  
& Integrated

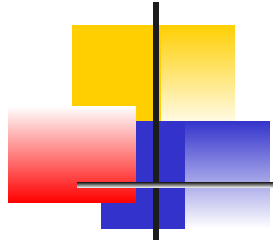


## What is Data Warehouse?

---

- In sum, data warehouse is a semantically consistent data store that serves as a physical implementation of a decision support data model and stores the information on which an enterprise needs to make strategic decisions

data on which enterprise needs to take strategic decisions



# Data Warehouse vs. Heterogeneous DBMS

---

- Traditional heterogeneous DB integration: A **query driven** approach
- Build wrappers/mediators <sup>?</sup> on top of heterogeneous databases
- Data warehouse: **update-driven**, high performance
- Information from heterogeneous sources is integrated in advance and stored in warehouses for direct query and analysis
- heterogeneous: query driven
- data warehouse: update driven

# Data Warehouse vs. Operational DBMS

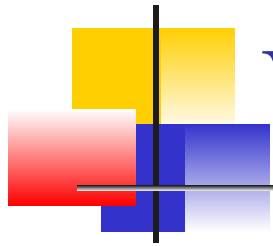
- OLTP (on-line transaction (query) processing)
  - Major task of traditional relational DBMS
  - Day-to-day operations: purchasing, inventory, banking, manufacturing, payroll, registration, accounting, etc.

- OLAP (on-line analytical processing)
  - Major task of data warehouse system
  - Data analysis and decision making



# Data Warehouse vs. Operational DBMS

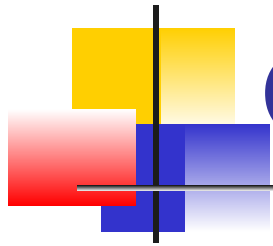
	OLTP	OLAP
<b>users</b>	clerk, IT professional	knowledge worker
<b>function</b>	day to day operations	decision support
<b>DB design</b>	application-oriented	subject-oriented
<b>data</b>	current, up-to-date detailed, flat relational isolated	historical, summarized, multidimensional integrated, consolidated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



# Why Separate Data Warehouse?

---

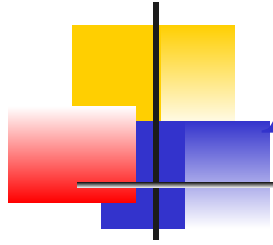
- Question: why not perform on-line analytical processing directly on such database instead of spending additional time and recourses to construct a separate data warehouse?
- High performance for both systems
  - DBMS— tuned for OLTP: access methods, indexing, concurrency control, recovery
  - Warehouse—tuned for OLAP: complex OLAP queries, multidimensional view, consolidation



# Outline

---

- What is Data Warehouse?
- Data Warehouse: A multidimensional data model
- Data Warehouse Architecture



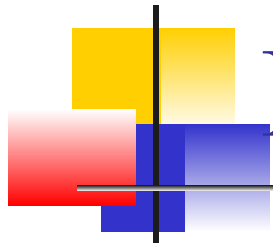
# A multi-dimensional data model

---

- A data warehouse is based on a multidimensional data model which views data in the form of a data cube

Multidimensional data Model → data  
↓  
data cube

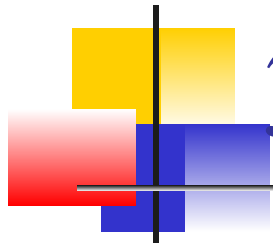




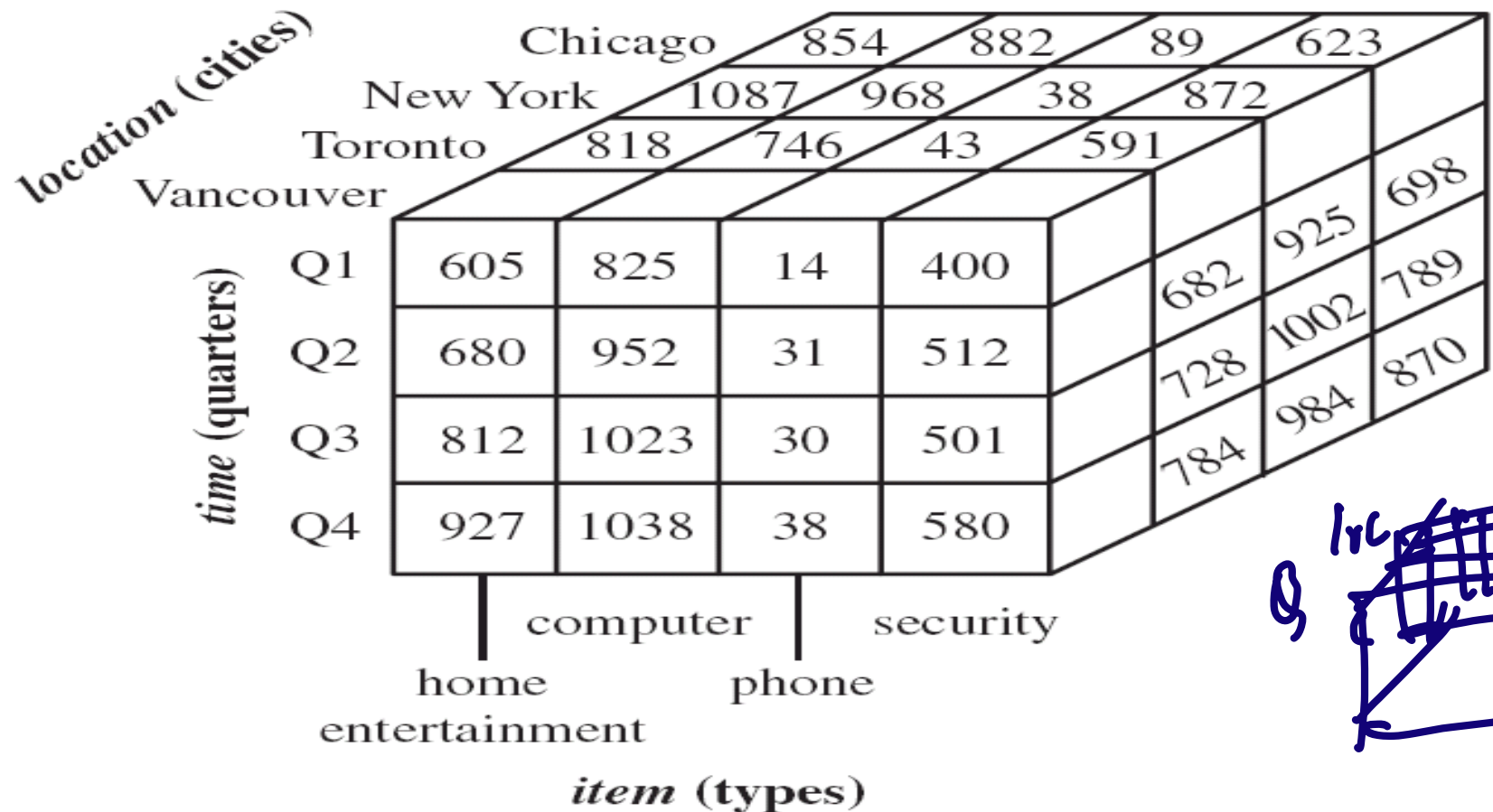
# Data cube

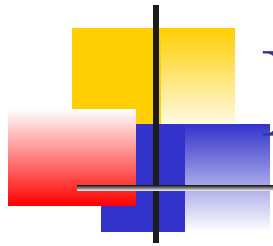
---

- A data cube, such as *sales*, allows data to be modeled and viewed in multiple dimensions
- Suppose ALLELETRONICS create a *sales* data warehouse with respect to dimensions
  - Time
  - Item
  - Location



# 3D Data cube Example



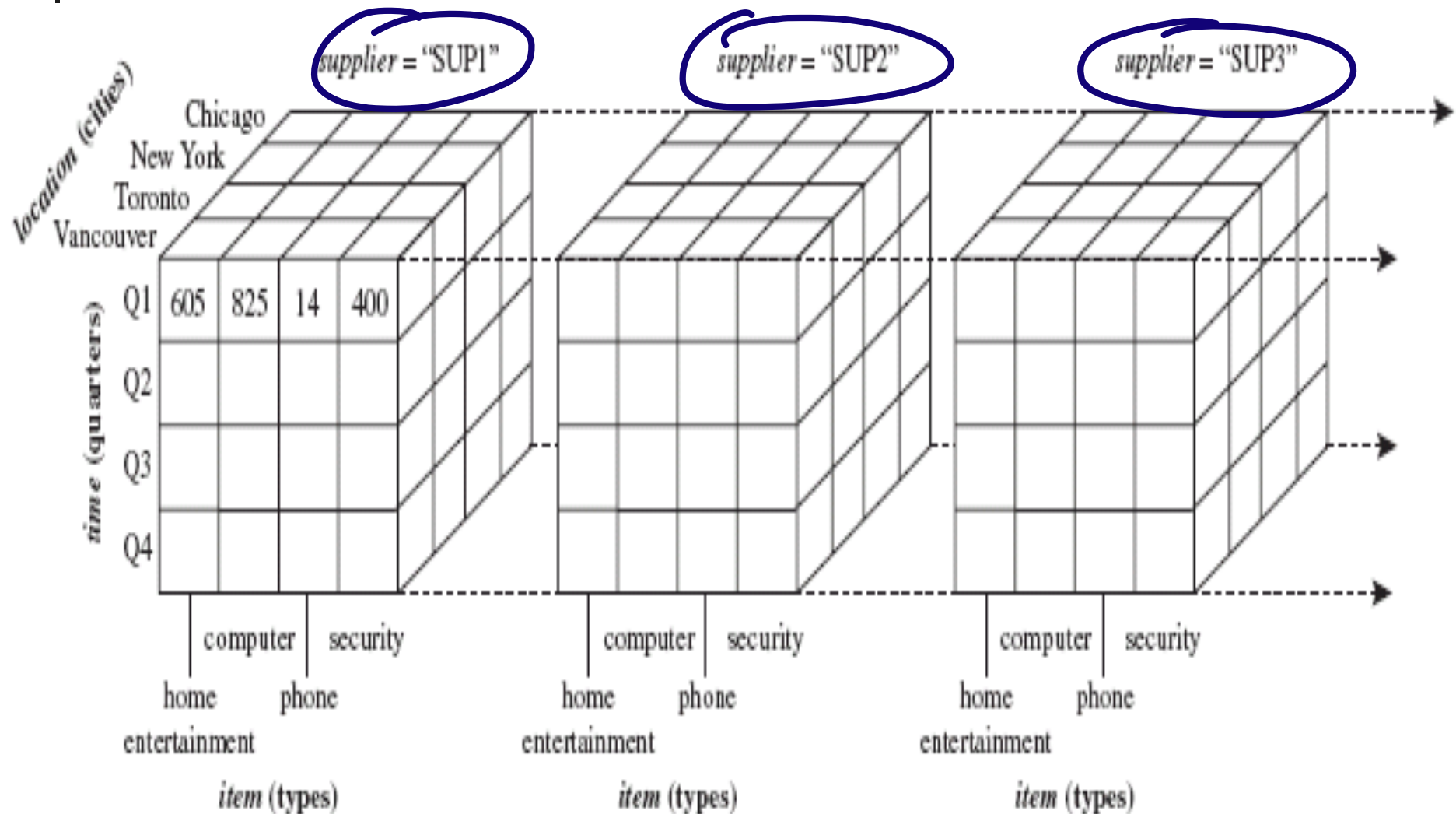


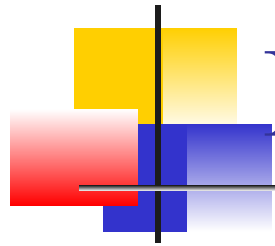
# Data cube

---

- A data cube, such as *sales*, allows data to be modeled and viewed in multiple dimensions
- Suppose ALLELETRONICS create a *sales* data warehouse with respect to dimensions
  - Time
  - Item
  - Location
  - Supplier

# 4D Data cube Example





# Practice Question

---

- What is a 5D cube looks like?



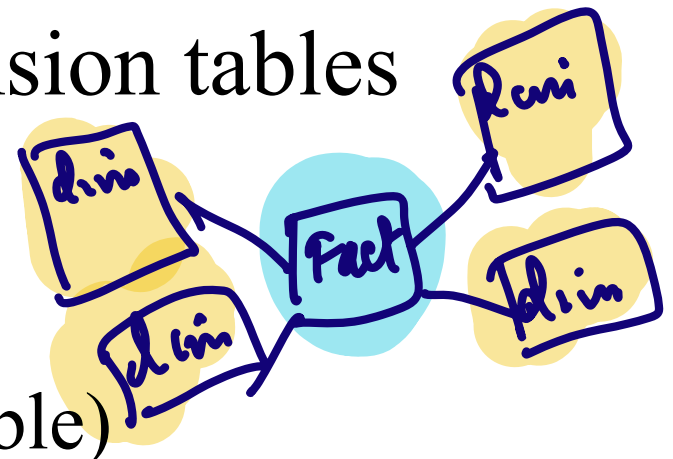
# Conceptual Modeling of Data Warehouses

---

- The most popular data model for a data warehouse is a multi-dimensional model
  - ↳ types → Star Schema, Snowflake, Fact constellation
- Such a model can exist in the form of:
  - Star schema
  - Snowflake schema
  - Fact constellations

# Conceptual Modeling of Data Warehouses

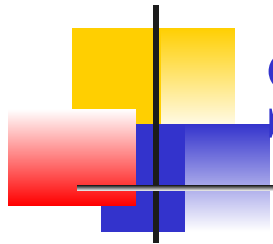
- **Star schema:** A fact table in the middle connected to a set of dimension tables



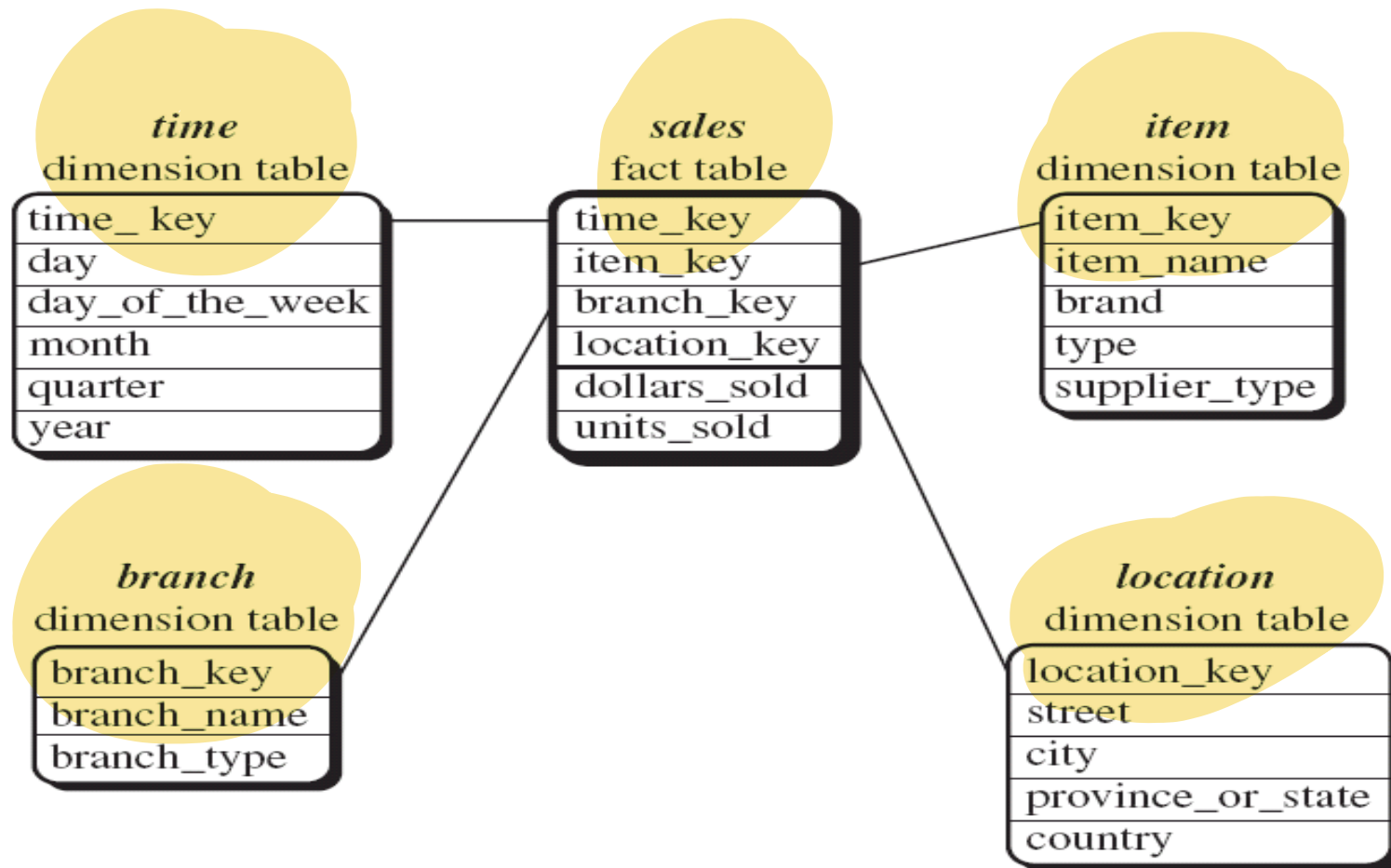
- **It contains:**
  - A large central table (fact table)
  - A set of smaller attendant tables (dimension table), one for each dimension

Fact: large central table

dim: Small tables, 1 for each dimensions



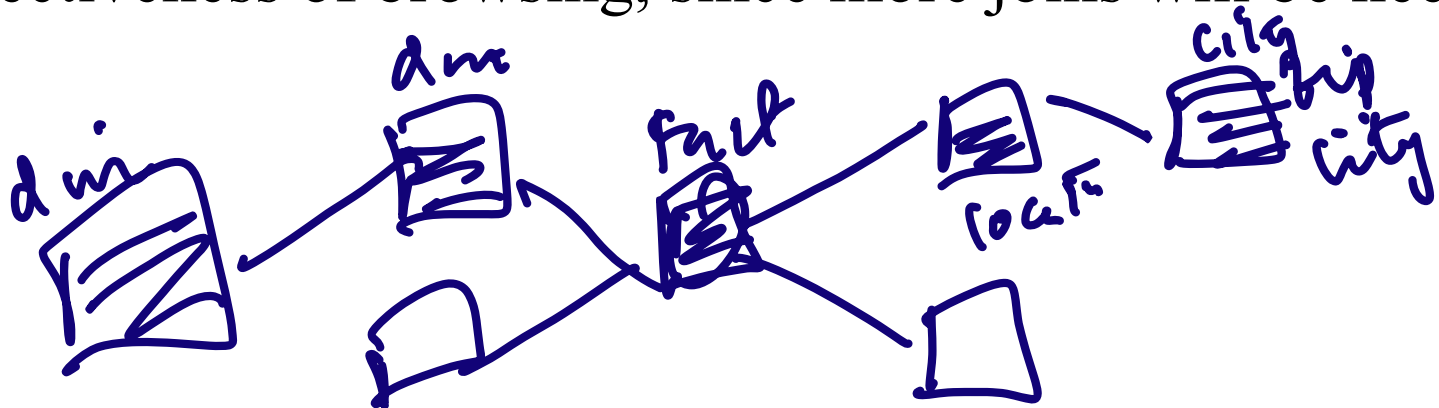
# Star schema

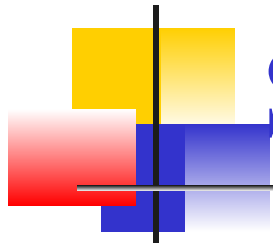




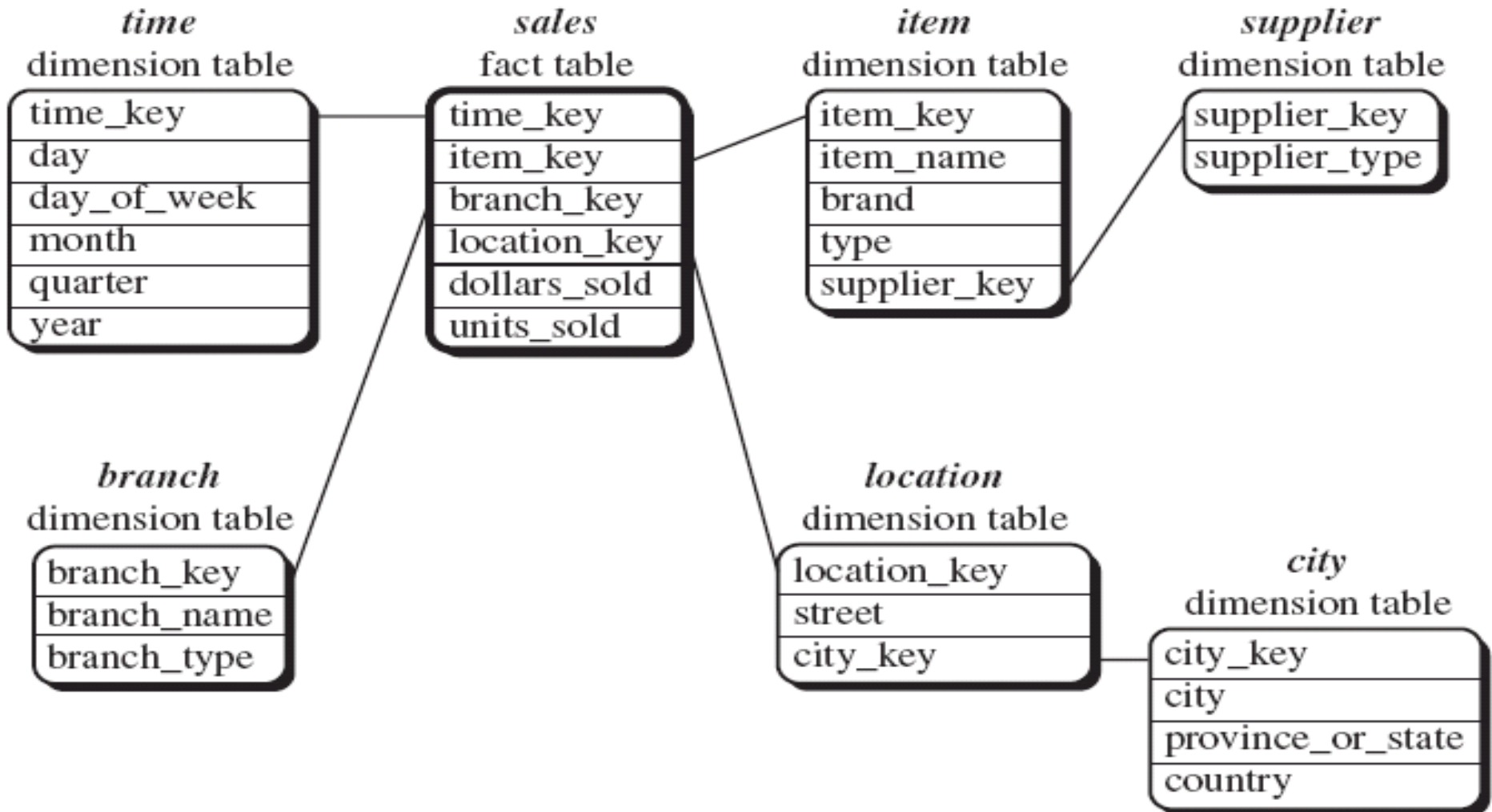
# Conceptual Modeling of Data Warehouses

- Snowflake schema: A refinement of star schema where some dimensional hierarchy is further splitting (normalized) into a set of smaller dimension tables, forming a shape similar to snowflake
- However, the snowflake structure can reduce the effectiveness of browsing, since more joins will be needed





# Snowflake schema



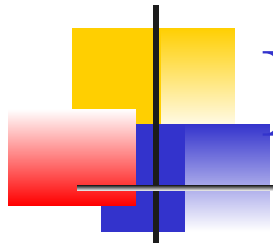


# Conceptual Modeling of Data Warehouses

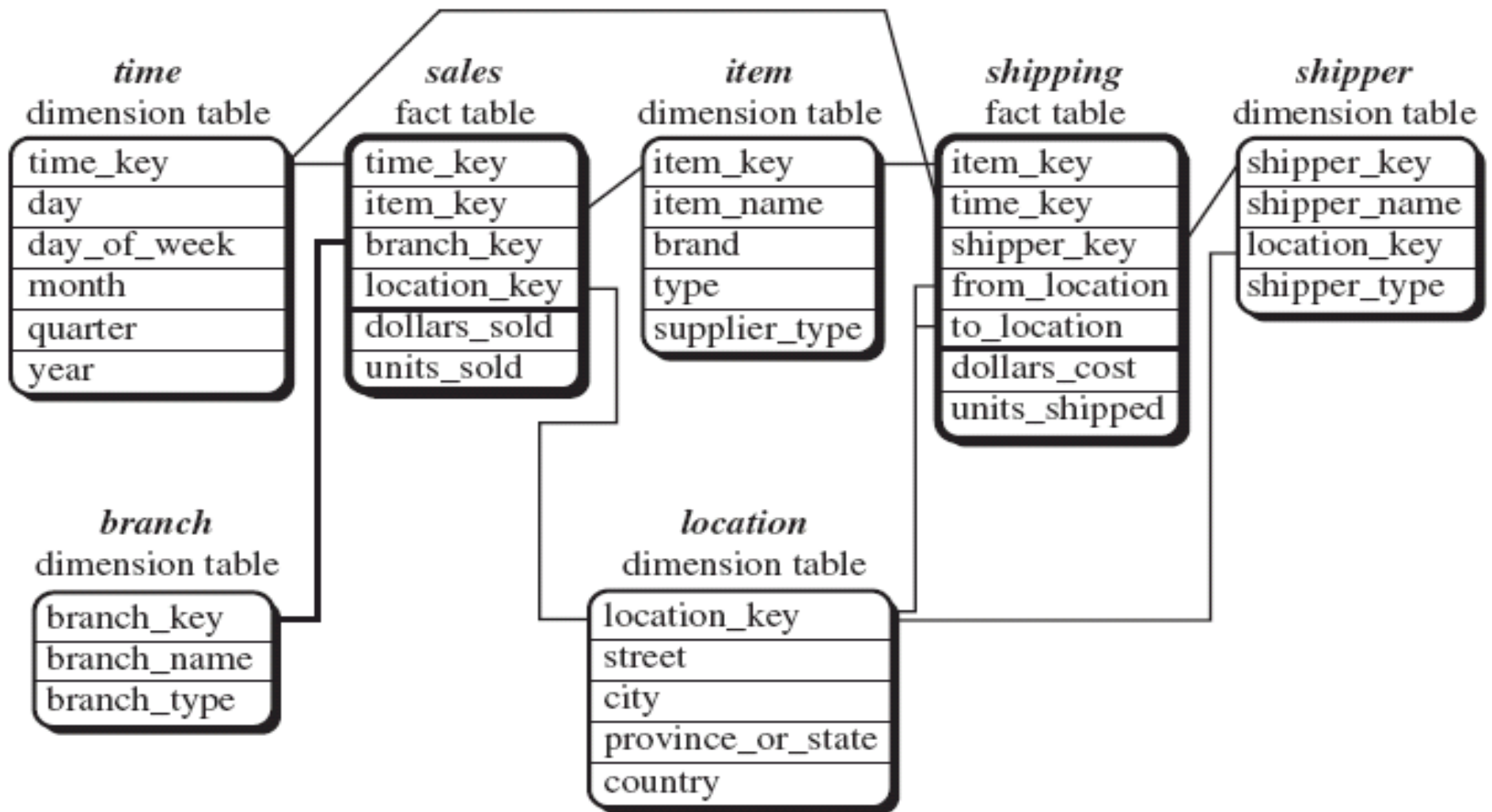
---

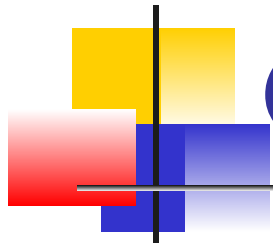
- Fact constellations: Multiple fact tables share dimension tables, viewed as a collection of stars, therefore called galaxy schema or fact constellation

→ Galaxy Schema



# Fact constellations



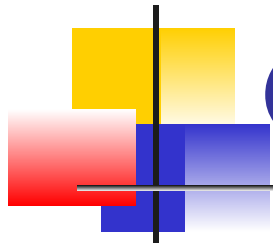


# Concept Hierarchies

---

- A **Concept Hierarchy** defines a sequence of mappings from a set of low-level concepts to high-level
- Consider a concept hierarchy for the dimension “**Location**”





# Concept Hierarchies

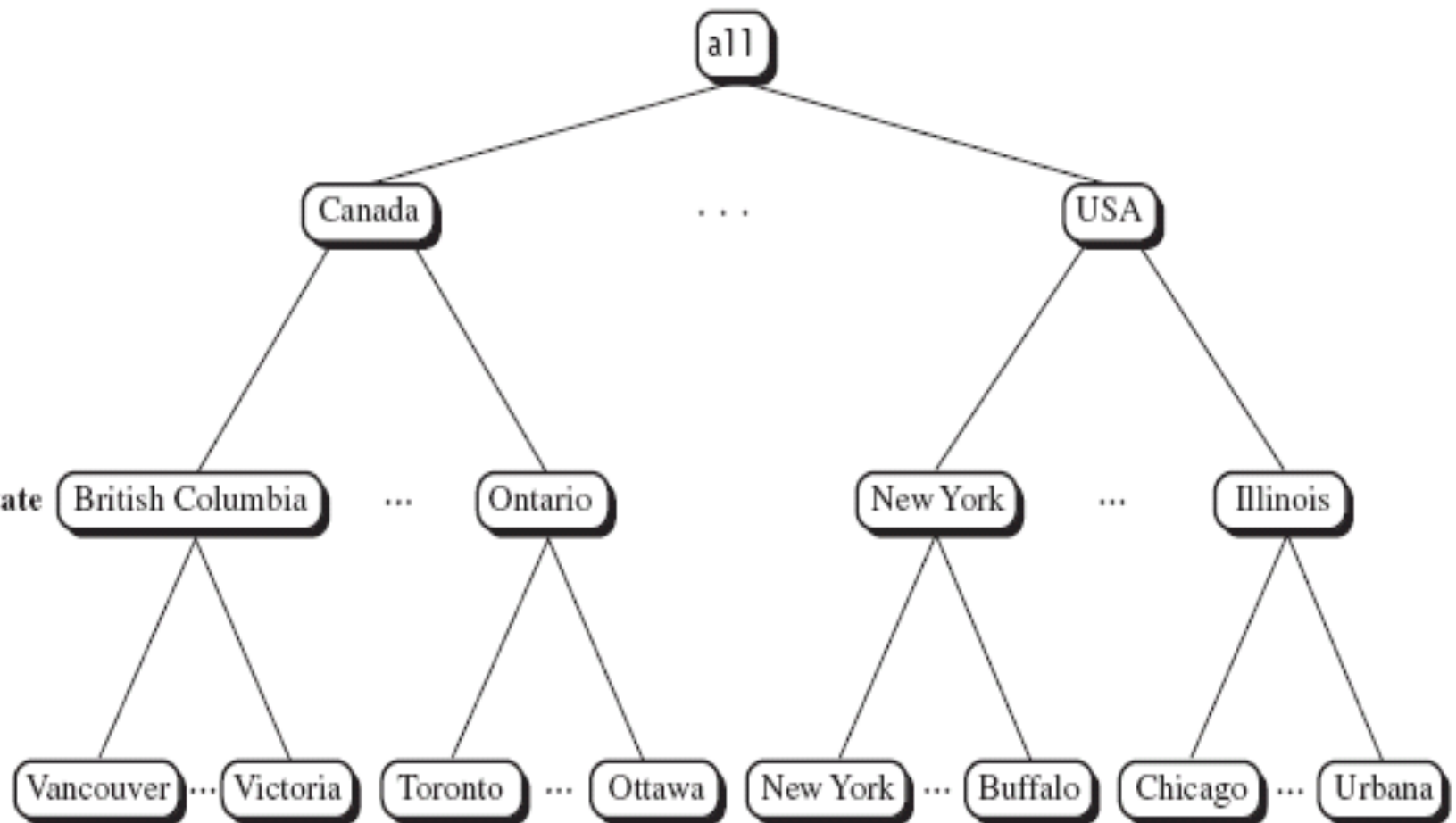
*location*

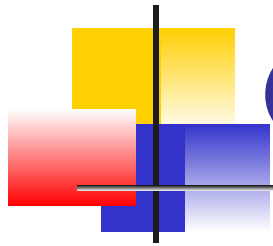
all

country

province\_or\_state

city

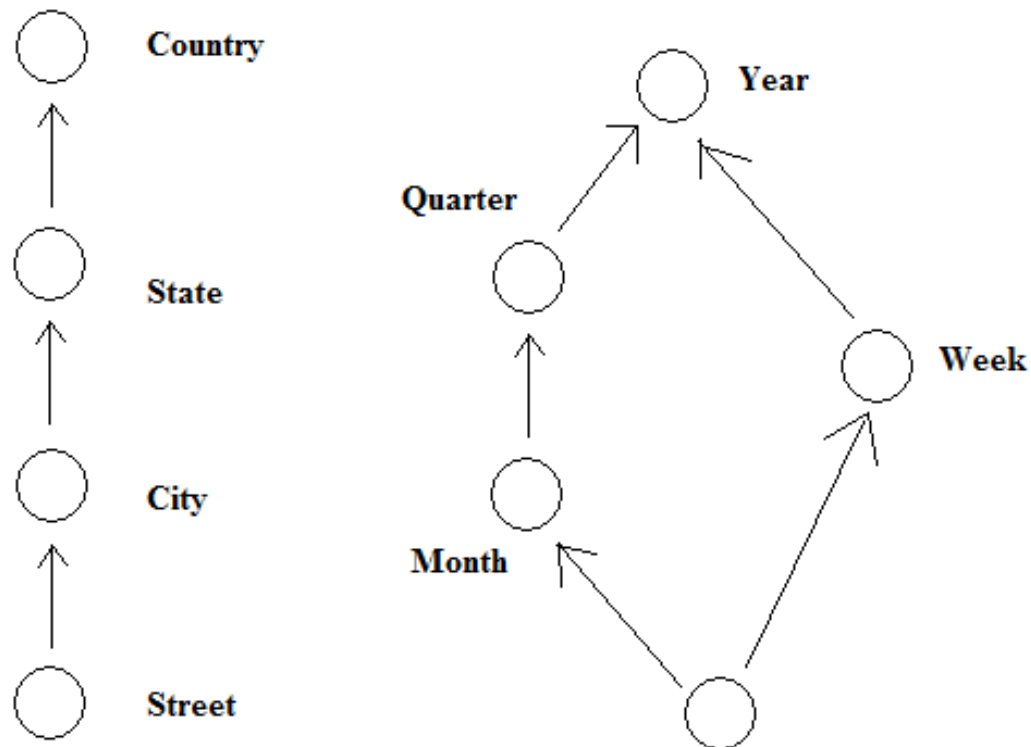


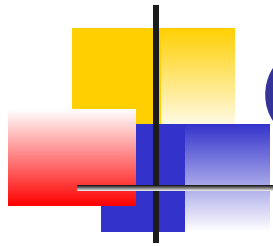


# Concept Hierarchies

---

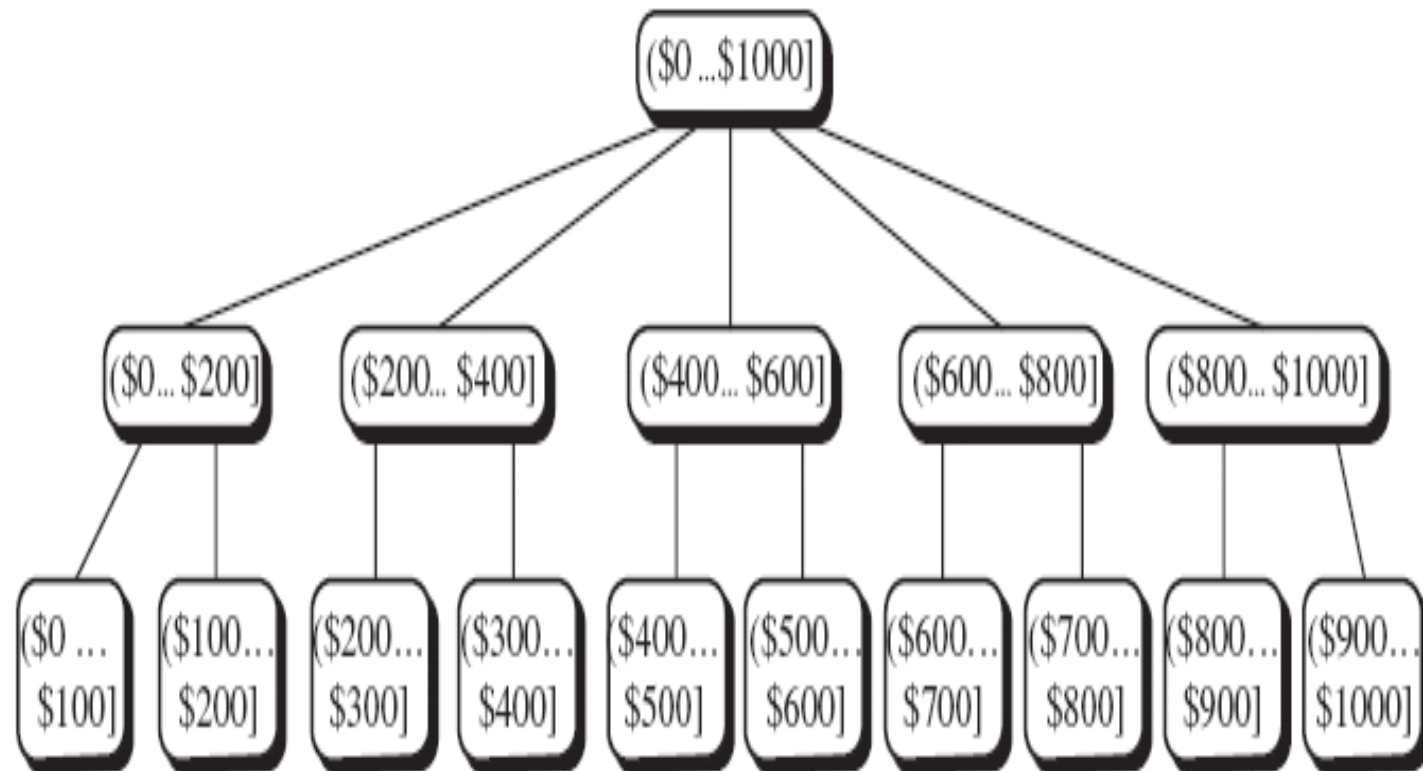
- Many concept hierarchies are implicit within the database system



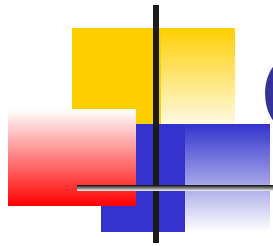


# Concept Hierarchies

- Concept hierarchies may also be defined by grouping values for a given dimension or attribute, resulting in a **set-grouping hierarchy**



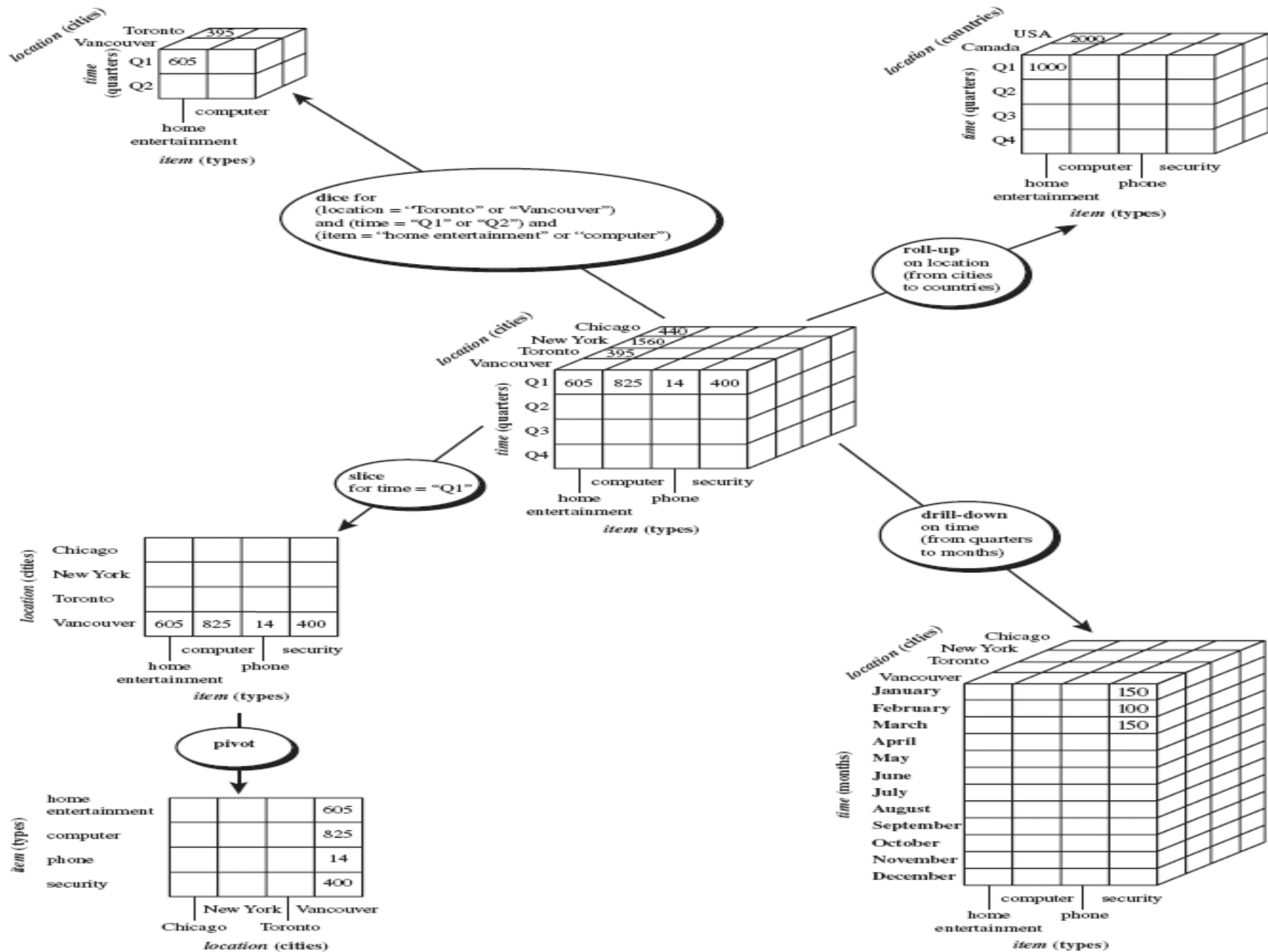


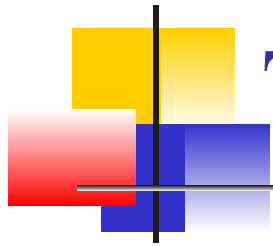


# OLAP Operation

---

- So, how are *concept hierarchies* useful in OLAP?
- In the multidimensional model, data are organized into multiple dimensions,
- And each dimension contains multiple levels of abstraction defined by concept hierarchies

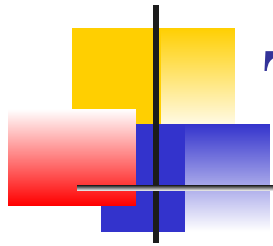




# Typical OLAP Operations

---

- **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

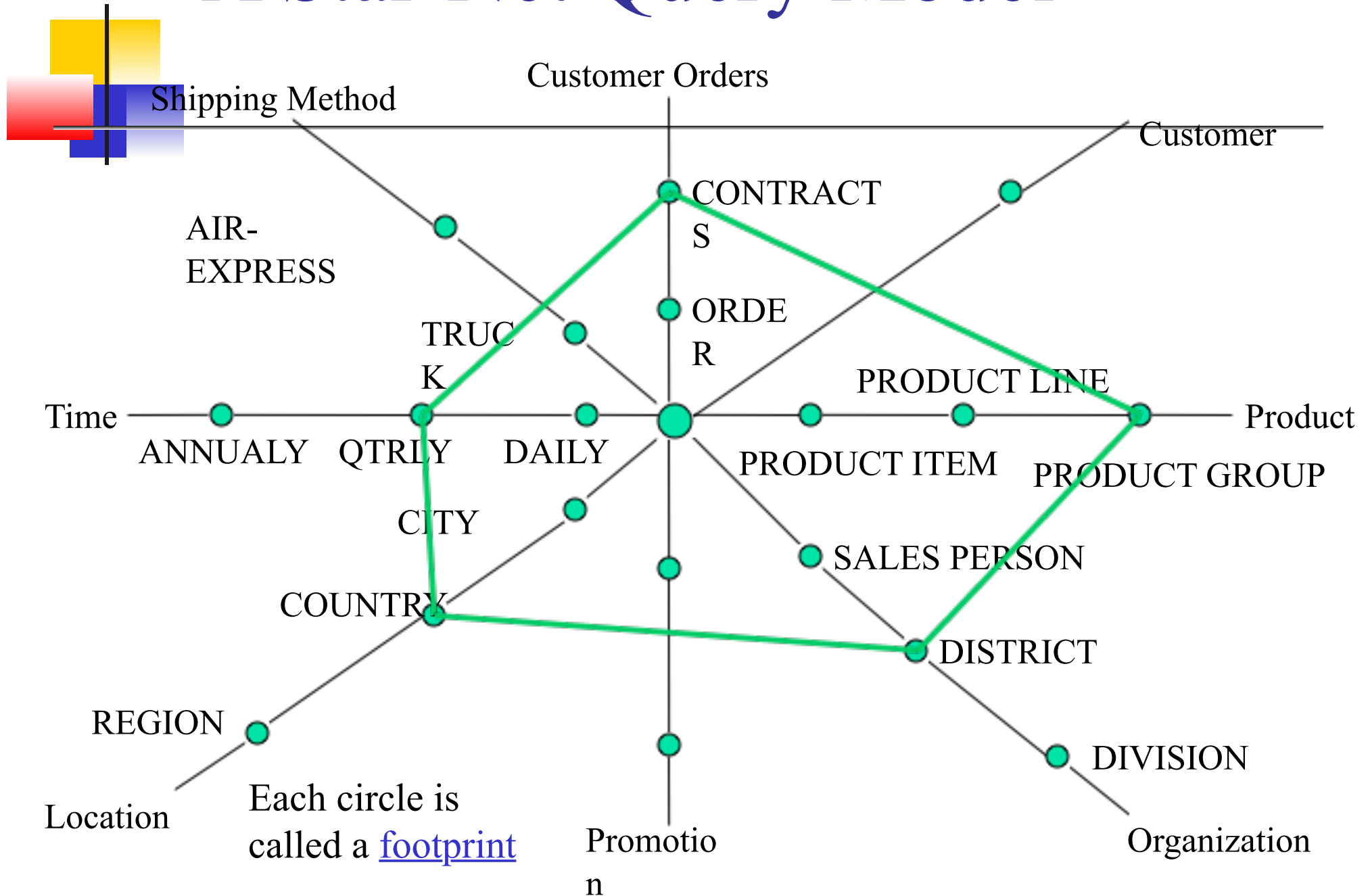


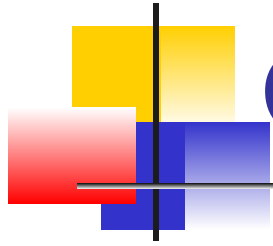
# Typical OLAP Operations

---

- **Slice and dice:** *project and select*
- **Pivot (rotate):**
  - *reorient the cube, visualization, 3D to series of 2D planes*

# A Star-Net Query Model

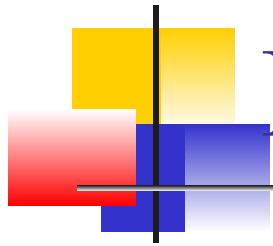




# Outline

---

- What is Data Warehouse?
- Data Warehouse: A multidimensional data model
- Data Warehouse Architecture

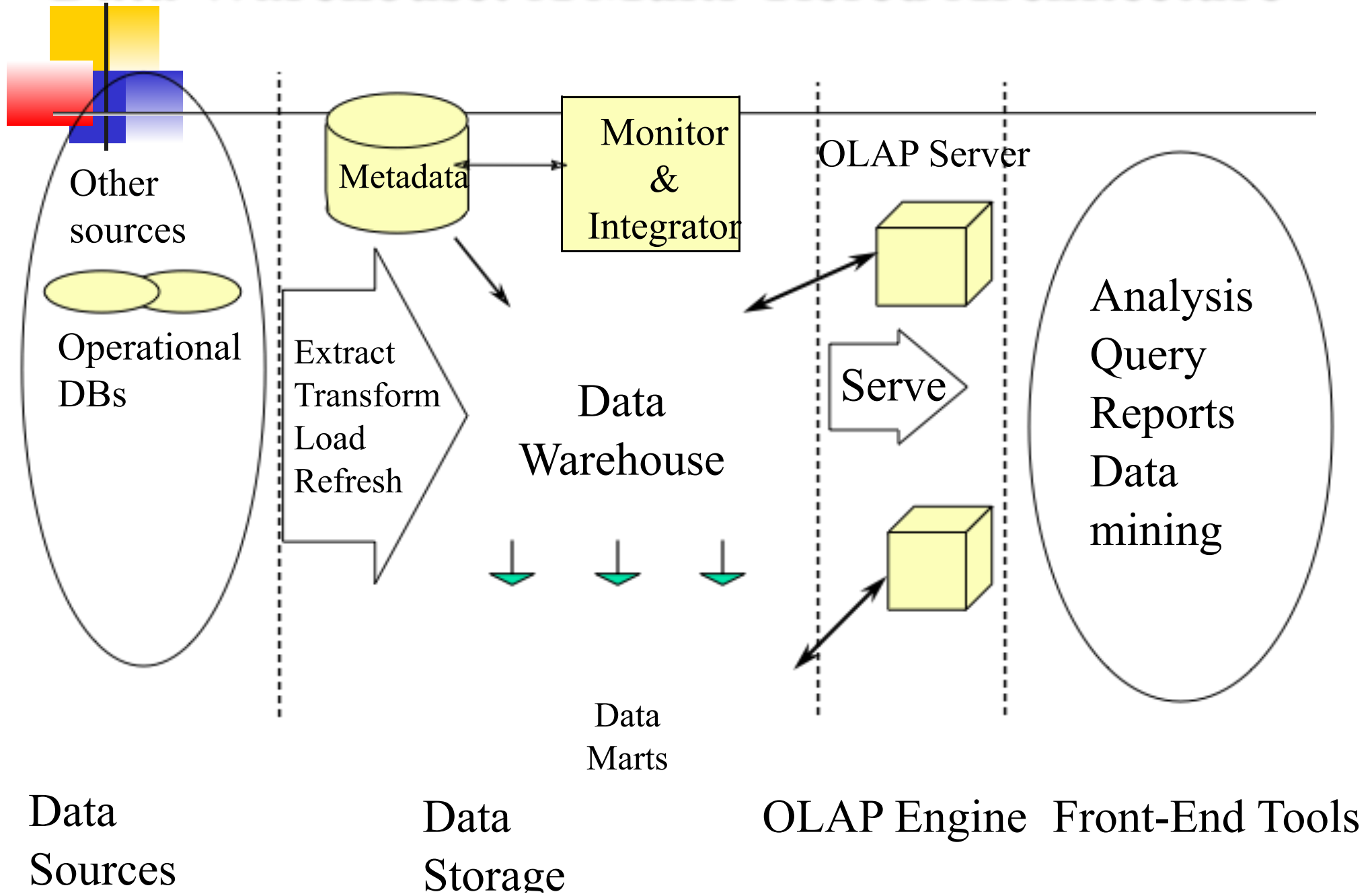


# Design of Data Warehouse

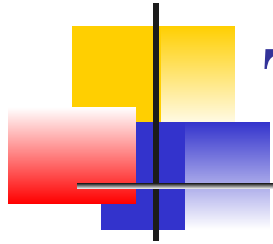
---

- Four views regarding the design of a data warehouse MUST be conserved
  - **Top-down view**
    - allows selection of the relevant information necessary for the data warehouse
  - **Data source view**
    - exposes the information being captured, stored, and managed by operational systems
  - **Data warehouse view**
    - consists of fact tables and dimension tables
  - **Business query view**
    - sees the perspectives of data in the warehouse from the view of end-user

# Data Warehouse: A Multi-Tiered Architecture





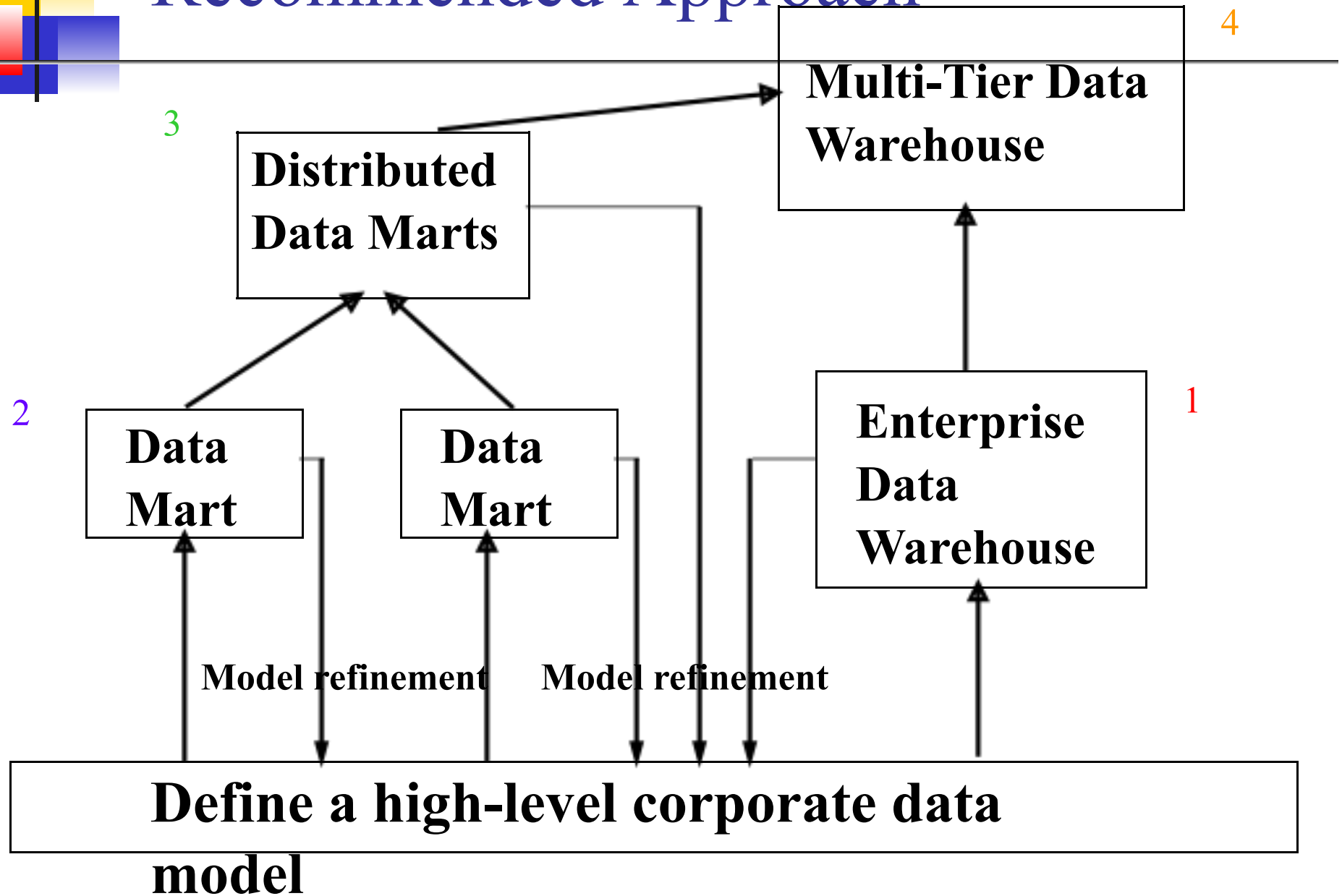
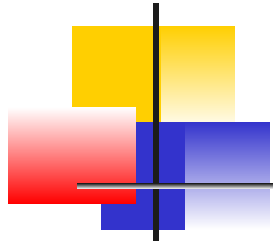


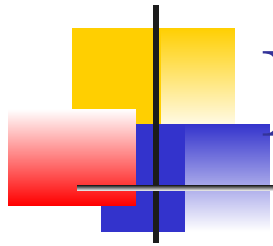
# Three Data Warehouse Models

---

- **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.

# Data Warehouse Development: A Recommended Approach





# 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



# Efficient Processing OLAP Queries

---

- Determine which operations should be performed on the available cuboids
- Determine which materialized cuboid(s) should be selected for OLAP op.

# Efficient Processing OLAP

## Queries

---

- Suppose we define a data cube of the form of
  - Day<month<quarter<year
  - Item\_name<brand<type
  - Street<city<province or state<country



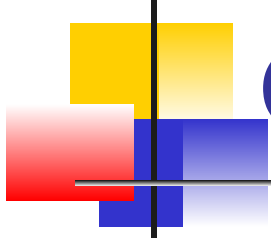
# Efficient Processing OLAP Queries

---

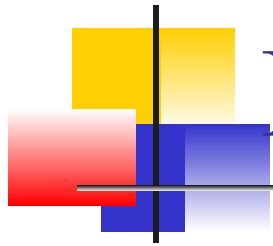
- Let the query to be processed be on {brand, province\_or\_state} with the condition “year = 2004”, and there are 4 materialized cuboids available:
  - 1) {year, item\_name, city}
  - 2) {year, brand, country}
  - 3) {year, brand, province\_or\_state}
  - 4) {item\_name, province\_or\_state} where year = 2004Which should be selected to process the query?

# Efficient Processing OLAP Queries

---



- Cuboid 1,3,4 can be applied
- Finer granularity data cannot be generated from coarser-granularity data

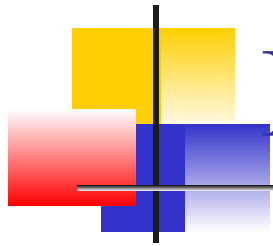


# Data Warehouse Usage

---

- Three kinds of data warehouse applications
  - **Information processing**
    - supports querying, basic statistical analysis, and reporting using crosstabs, tables, charts and graphs
  - **Analytical processing**
    - multidimensional analysis of data warehouse data
    - supports basic OLAP operations, slice-dice, drilling, pivoting





# Data Warehouse Usage

---

- **Data mining**
- knowledge discovery from hidden patterns
- supports associations, constructing analytical models, performing classification and prediction, and presenting the mining results using visualization tools

