# Chapter 20: Data Analysis

**Database System Concepts, 6th Ed.**

©Silberschatz, Korth and Sudarshan
See www.db-book.com for conditions on re-use

# Decision-Support Systems: Overview

- **Data analysis** tasks are simplified by specialized tools and SQL extensions

    - Example tasks

        ‣ For each product category and each region, what were the total sales in the last quarter and how do they compare with the same quarter last year

        ‣ As above, for each product category and each customer category

- **Statistical analysis** packages (e.g., : S++) can be interfaced with databases

    - Statistical analysis is a large field, but not covered here

- **Data mining** seeks to discover knowledge automatically in the form of statistical rules and patterns from large databases.

- A **data warehouse** archives information gathered from multiple sources, and stores it under a unified schema, at a single site.

    - Important for large businesses that generate data from multiple divisions, possibly at multiple sites

    - Data may also be purchased externally

# Data Warehousing

- Data sources often store only current data, not historical data

- Corporate decision making requires a unified view of all organizational data, including historical data

- A **data warehouse** is a repository (archive) of information gathered from multiple sources, stored under a unified schema, at a single site

  - Greatly simplifies querying, permits study of historical trends

  - Shifts decision support query load away from transaction processing systems

# Data Warehouse vs. Operational DBMS

- OLTP (on-line transaction 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

- Distinct features (OLTP vs. OLAP):

    - Data contents: current, detailed vs. historical, consolidated

    - Database design: ER, Normalized design + application vs. star + subject

    - View: current, local vs. evolutionary, integrated

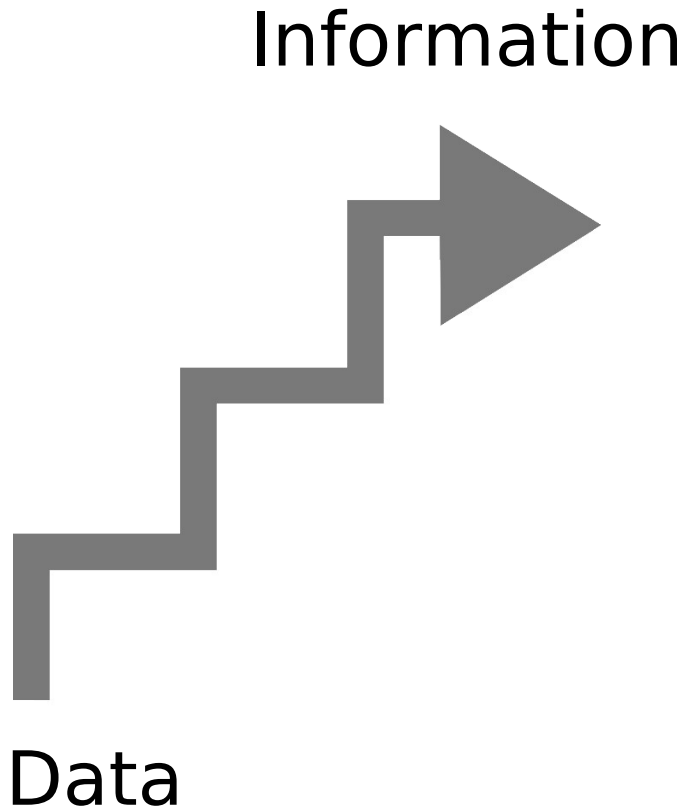    - Access patterns: update vs. read-only but complex queries

# OLTP vs. OLAP

|                    | 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 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

# What is Data Warehousing?

Information

A process of transforming data into information and making it available to users in a timely enough manner to make a difference

[Forrester Research, April 1996]

Data

# Data Warehouse?

- Different definitions -
  - 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

# Data Warehouse—Subject-Oriented

◼ Organized around major subjects.

[For example - 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

  ‣ "Interoperability"

- 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—Non-Volatile

- 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*.

# Data Warehouse vs. Heterogeneous DBMS

■ Traditional heterogeneous DB integration:

- Build wrappers/mediators on top of heterogeneous databases

- Query driven approach

  ‣ When a query is posed to a client site, a meta-dictionary is used to translate the query into queries appropriate for individual heterogeneous sites involved, and the results are integrated into a global answer set

  ‣ Complex information filtering, compete for resources

■ Data warehouse: update-driven, high performance

- Information from heterogeneous sources is integrated in advance and stored in warehouses for direct query and analysis

# Data Warehouse vs. Operational DBMS

- **OLTP (on-line transaction 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

- **Distinct features (OLTP vs. OLAP):**

  - User and system orientation: customer vs. market

  - Data contents: current, detailed vs. historical, consolidated

  - Database design: ER + application vs. star + subject

  - View: current, local vs. evolutionary, integrated

  - Access patterns: update vs. read-only but complex queries

# Why 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:**
  - <u>missing data</u>: Decision support requires historical data which operational DBs do not typically maintain
  - <u>data consolidation</u>:  DS requires consolidation (aggregation, summarization) of data from heterogeneous sources
  - <u>data quality</u>: different sources typically use inconsistent data representations, codes and formats which have to be reconciled

# Multi-dimensional Data Model – 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

- 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

# Example of Star Schema

**time**

time_key
day
day_of_the_week
month
quarter
year

**item**

item_key
item_name
brand
type
supplier_type

Sales Fact Table

time_key

item_key

branch_key

location_key

units_sold

dollars_sold

avg_sales

**branch**

branch_key
branch_name
branch_type

**location**

location_key
street
city
province_or_state
country

Measures

# Example of Snowflake Schema

**time**

time_key
day
day_of_the_week
month
quarter
year

**item**

item_key
item_name
brand
type
supplier_key

**supplier**

supplier_key
supplier_type

Sales Fact Table

| time_key |
| item_key |
| branch_key |
| location_key |
| units_sold |
| dollars_sold |
| avg_sales |

**branch**

branch_key
branch_name
branch_type

**location**

location_key
street
city_key

**city**

city_key
city
province_or_state
country

Measures

# Example of Fact Constellation

**time**

time_key
day
day_of_the_week
month
quarter
year

**item**

item_key
item_name
brand
type
supplier_type

Shipping Fact Table

Sales Fact Table

time_key

item_key

branch_key

location_key

units_sold

dollars_sold

avg_sales

Measures

**branch**

branch_key
branch_name
branch_type

**location**

location_key
street
city
province_or_state
country

time_key

item_key

shipper_key

from_location

to_location

dollars_cost

units_shipped

**shipper**

shipper_key
shipper_name
location_key
shipper_type

# Data Warehouse: A Multi-Tiered Architecture



**Data Sources**

Other sources

Operational DBs

Metadata

Monitor & Integrator

Extract
Transform
Load
Refresh

Data Warehouse

Data Marts

OLAP Server

Serve

Analysis
Query
Reports
Data mining

**Data Storage**

**OLAP Engine**

**Front-End Tools**

# Example of Star Schema

**date**

date_key
day
day_of_the_week
month
quarter
year

**item**

item_key
item_name
brand
type
supplier_type

Sales Fact Table

| date_key |
| item_key |
| location_key |
| units_sold |
| dollars_sold |

**location**

location_key
street
city
state_or_province
country
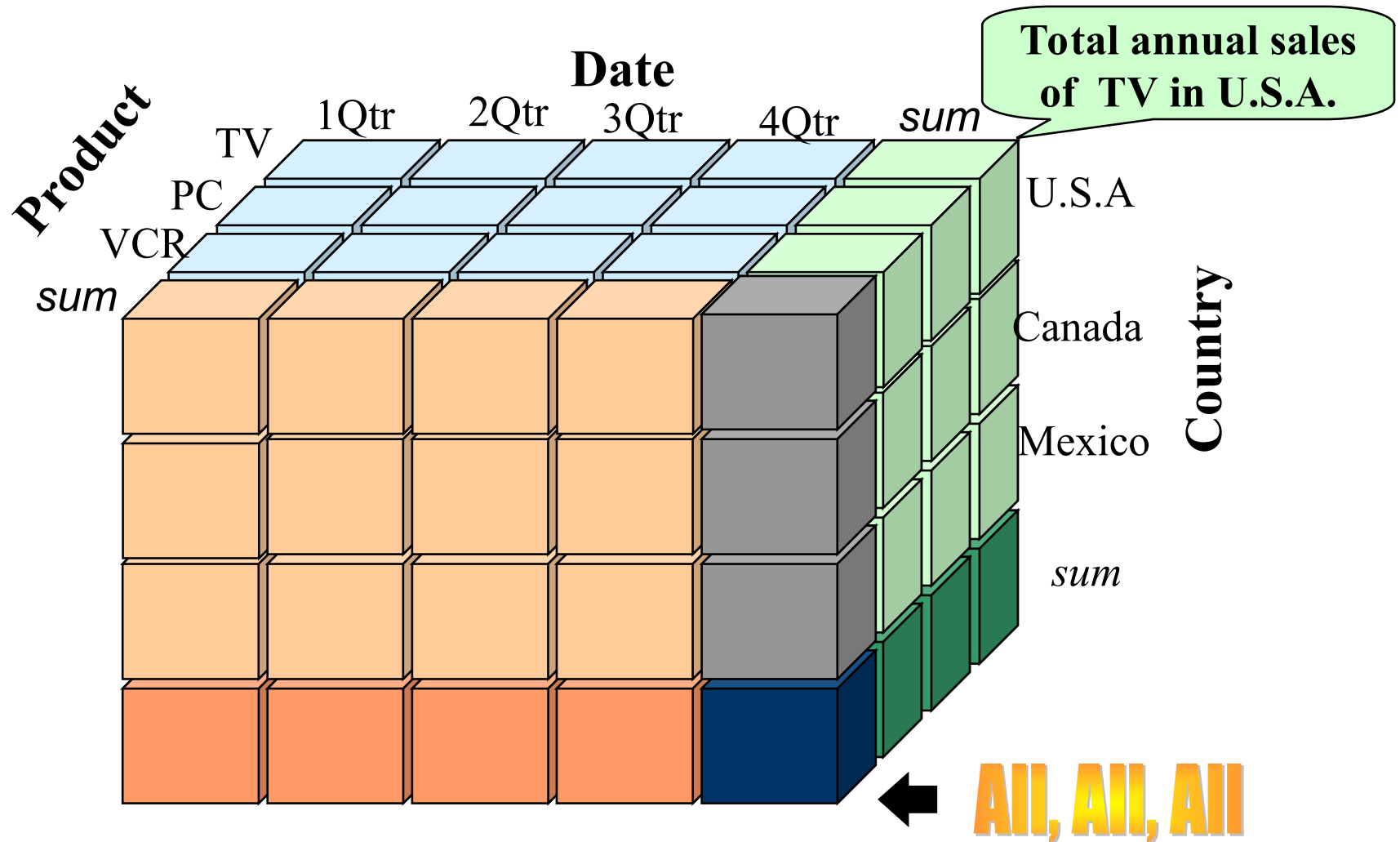
Measures

# From Tables 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 date(day, week, month, quarter, year) or location

  - Fact table contains measures (such as dollars_sold) and keys to each of the related dimension tables

- 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.
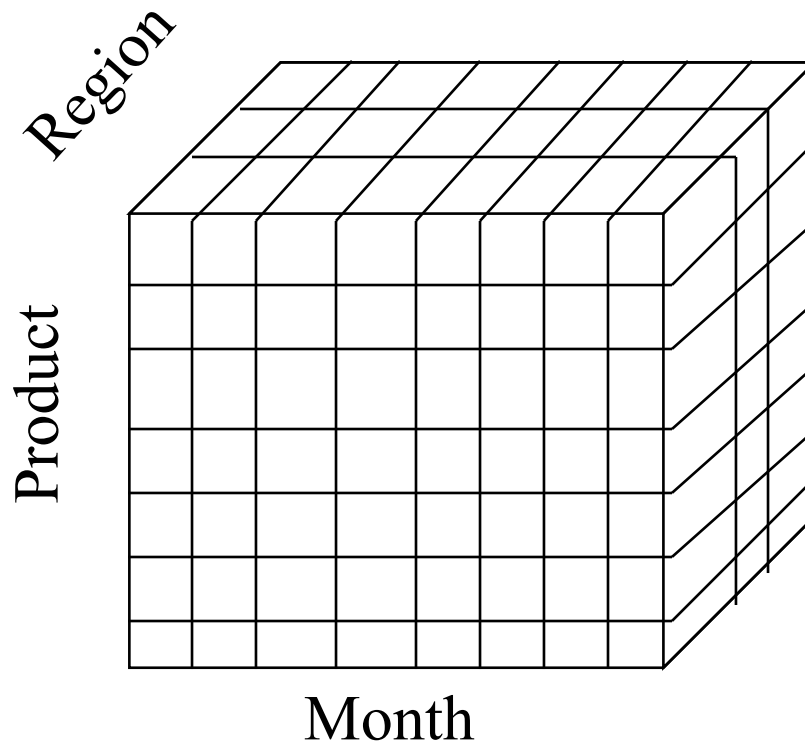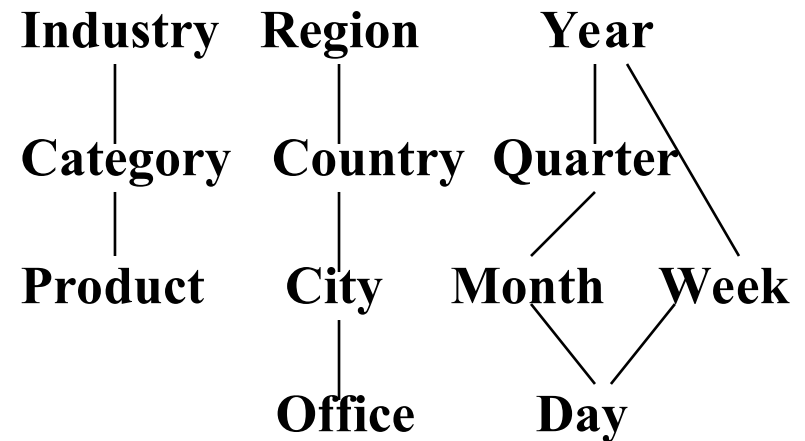
# A Sample Data Cube

# Multidimensional Data

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

**Dimensions: Product, Location, Time**
**Hierarchical summarization paths**

Region

Product

Month

| Industry | Region | Year |
| --- | --- | --- |
| Category | Country | Quarter |
| Product | City | Month    Week |
| | Office | Day |

# A Concept Hierarchy: Dimension (location)

all ·········· region ·········· country ·········· city ·········· office

```
                                all
                 ┌───────────────┴───────────────┐
              Europe           ...          North_America
            ┌────┴────┐                    ┌────────┴────────┐
        Germany  ...  Spain             Canada   ...   Mexico
          ┌──┴──┐      ⋀               ┌────┴────┐        ⋀
     Frankfurt  ...                Vancouver ... Toronto
        ⋀                          ┌────┴────┐     ⋀
                              L. Chan ... M. Wind
```

# Cube: A Lattice of Cuboids

all — 0-D(apex) cuboid

time    item    location    supplier — 1-D cuboids

**time,location**    **item,location**    **location,supplier**

**time,item**    **time,supplier**    **item,supplier** — 2-D cuboids

**time,location,supplier** — 3-D cuboids

**time,item,location**    **time,item,supplier**    **item,location,supplier**

**time, item, location, supplier** — 4-D(base) cuboid
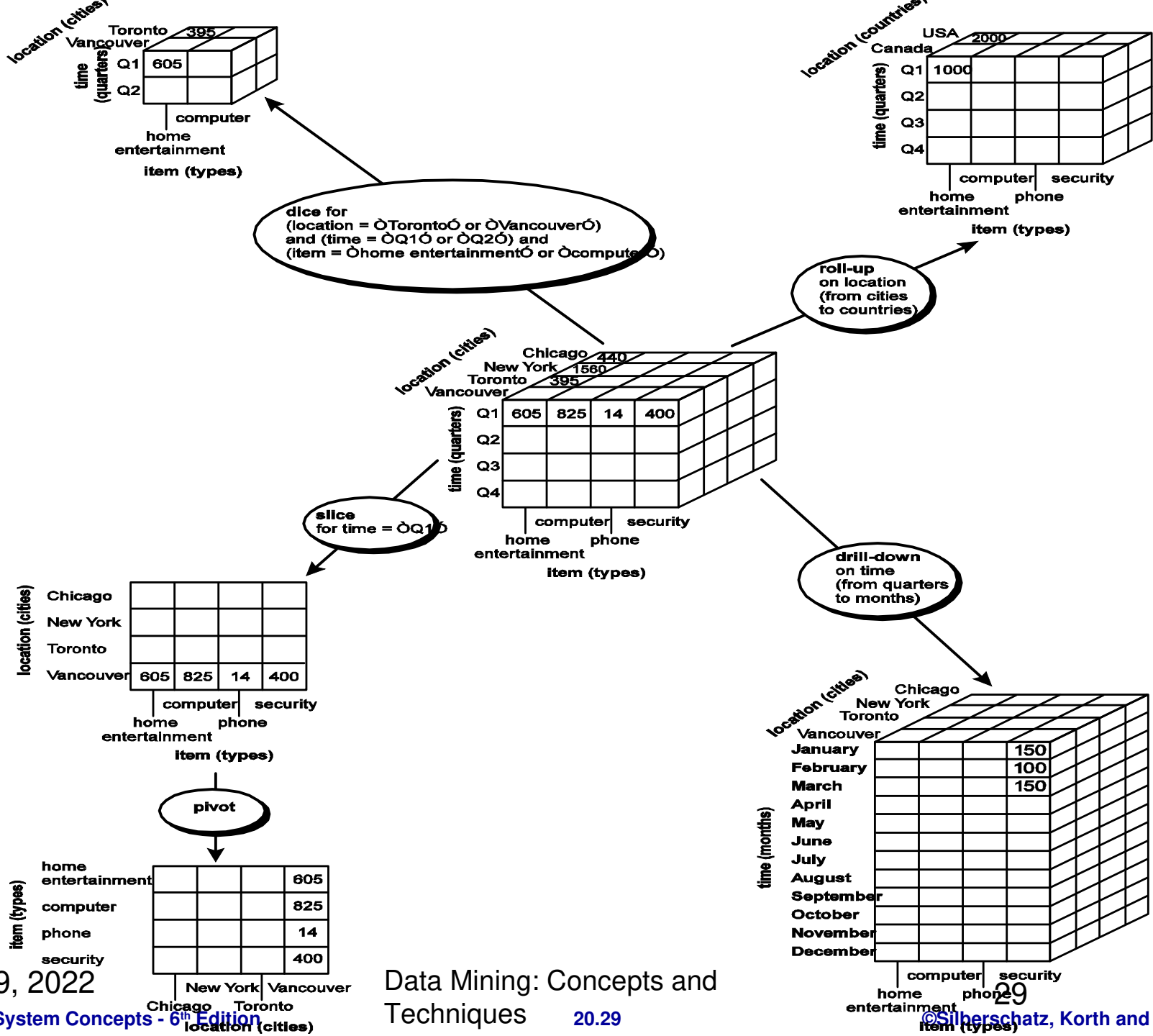
# Typical OLAP Operations

- **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)*

# Important Instructions

- Read these slides and also the paper on Association Rule Mining

# End of Chapter

**Database System Concepts, 6<sup>th</sup> Ed.**

Wait, I must use LaTeX not HTML sup.

# End of Chapter

**Database System Concepts, 6$^{th}$ Ed.**

# End of Chapter

**Database System Concepts, 6$^{\text{th}}$ Ed.**

**©Silberschatz, Korth and Sudarshan**
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# Example of Snowflake Schema

**time**
time_key
day
day_of_the_week
month
quarter
year

**item**
item_key
item_name
brand
type
supplier_key

**supplier**
supplier_key
supplier_type

**branch**
branch_key
branch_name
branch_type

Sales Fact Table

time_key

item_key

branch_key

location_key

units_sold

dollars_sold

avg_sales

Measures

**location**
location_key
street
city_key

**city**
city_key
city
state_or_province
country

# Example of Fact Constellation

**time**

time_key
day
day_of_the_week
month
quarter
year

**item**

item_key
item_name
brand
type
supplier_type

Shipping Fact Table

Sales Fact Table

time_key

item_key

branch_key

location_key

units_sold

dollars_sold

avg_sales

Measures

**branch**

branch_key
branch_name
branch_type

**location**

location_key
street
city
province_or_state
country

time_key

item_key

shipper_key

from_location

to_location

dollars_cost

units_shipped

**shipper**

shipper_key
shipper_name
location_key
shipper_type

# Data Warehousing



data warehouse