# CSCI 5408 Data Analytics: DM and DW Tech (Week 9)

- **Ass4 Due: Mar 14**
  - Brightspace: Assignment 4, Tutorial slides, etc.
  - Help hours: Fri, 1:00-2:30PM, CS 233
- Write answers for review questions
  - Final Exam: Apr 20, 3:30-5:30 PM
- Reading: Lectures: 13-14, Text: Ch4 of 3$^{rd}$ edition, or Ch3 of 2$^{nd}$ edition

# 3. Data Warehouses and OLAP

(Textbook: Ch4 of 3$^{rd}$ edition, or Ch3 of 2$^{nd}$ edition)

- Objectives of DW/OLAP
- What is a DW?
- Multi-dimensional data space model
- DW schemas
- OLAP operations
- Aggregations
- DW architecture
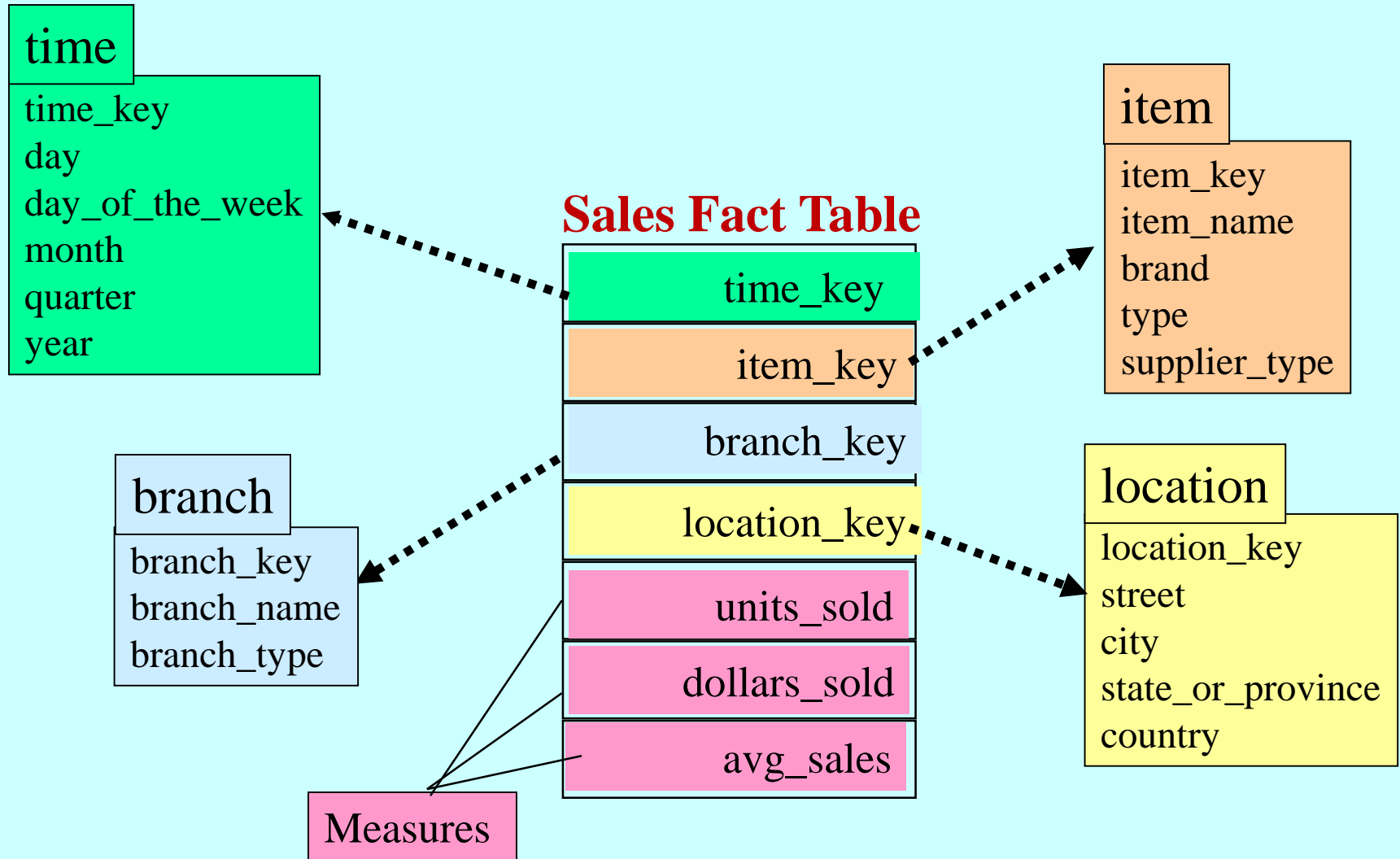- From DW to DM

# A Comparison: 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 | read only, 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 |

https://courses.cs.washington.edu/courses/csep573/01sp/lectures/class1/sld025.htm

3

# DW Schemas: Conceptual Models

- Basic DW Structure: dimensions & measures

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

- **Snowflake schema:** A refinement of star schema where some <u>dimensional hierarchy is normalized</u> into a set of smaller dimension tables, forming a shape similar to snowflake

- **Fact constellations:** <u>Multiple fact tables</u> share dimension tables, viewed as a collection of stars, therefore called galaxy schema or fact constellation
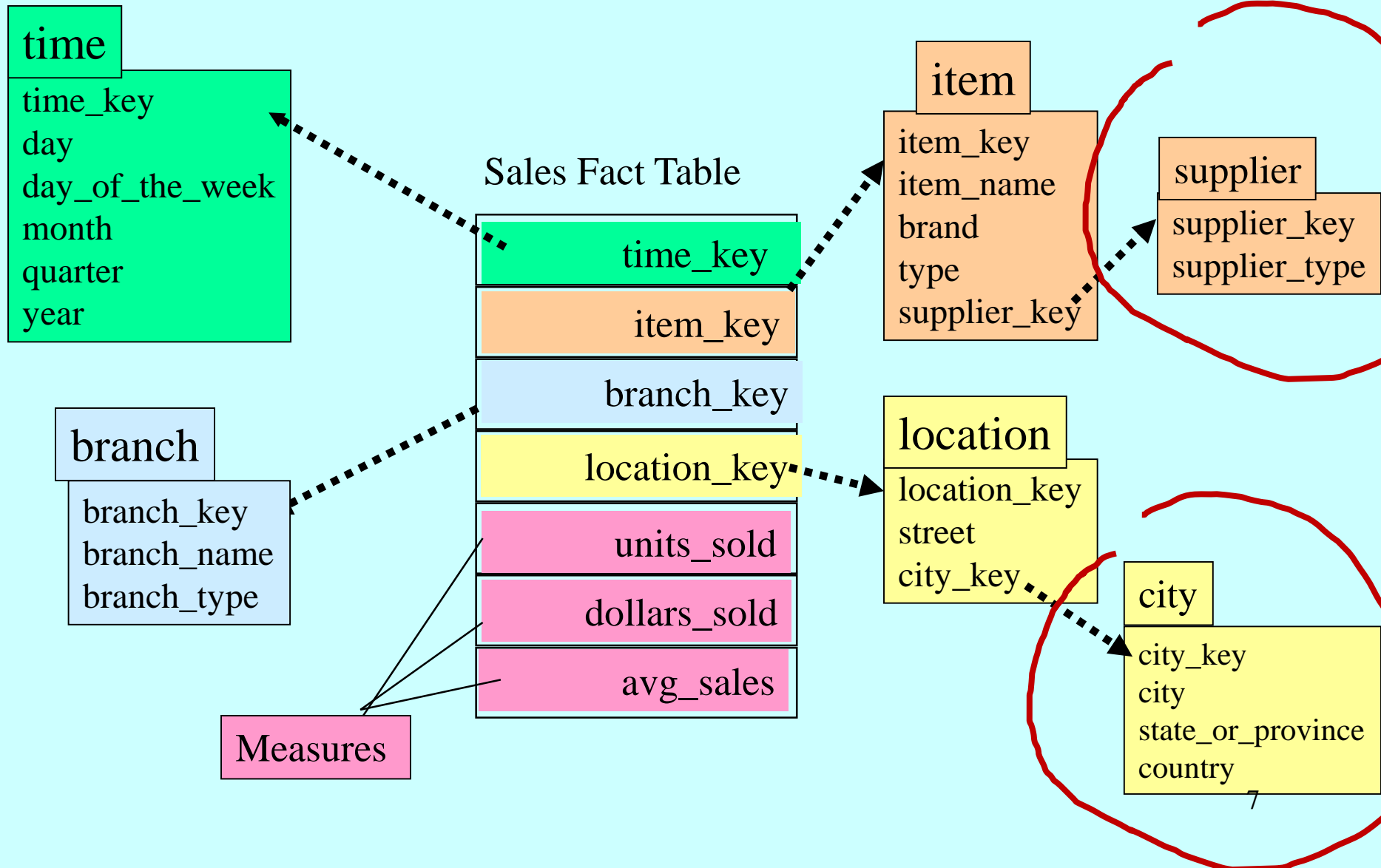
# Example of Star Schema
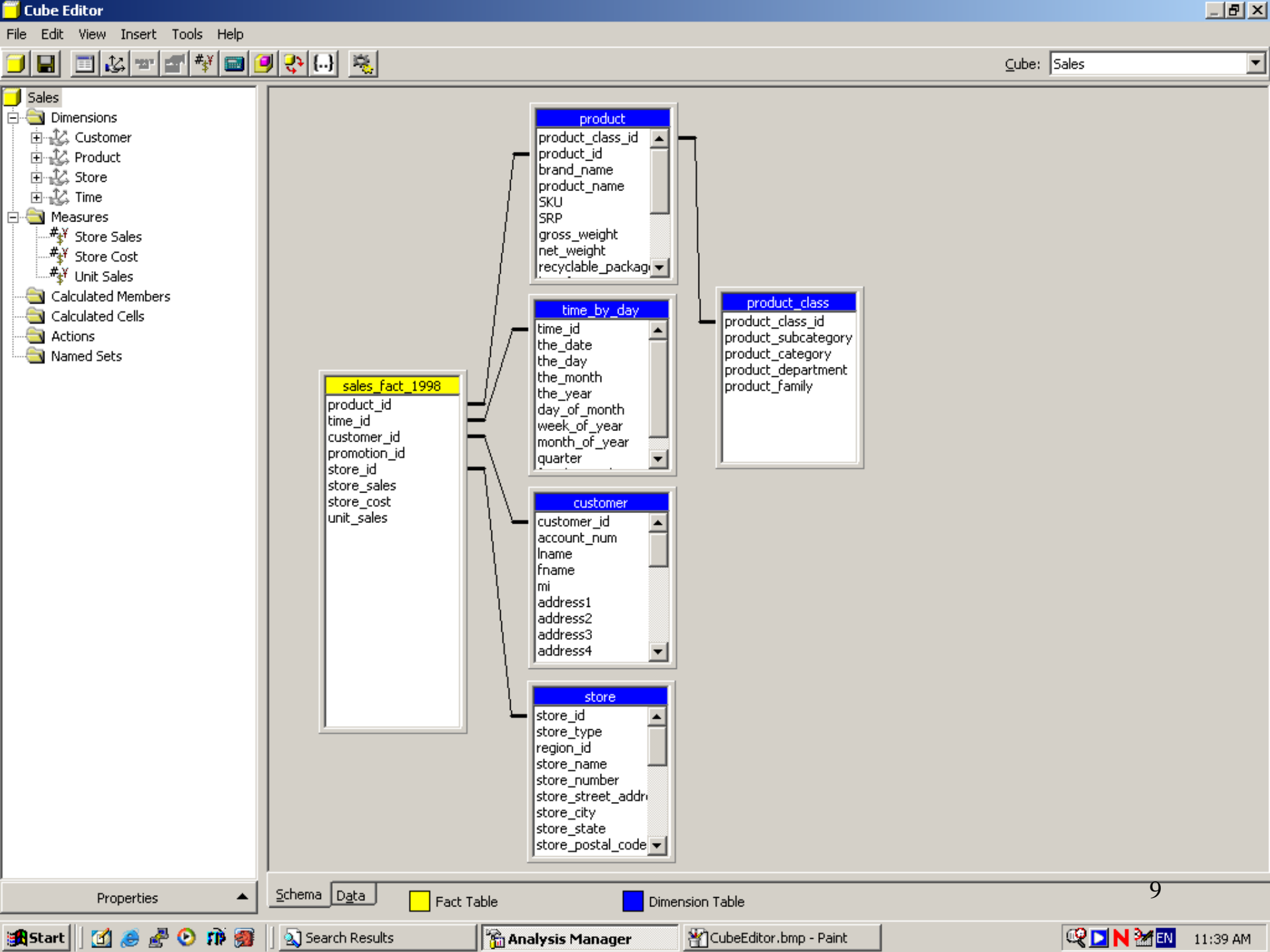
# Example of Star Schema (cont)

- Star schema is the most common used schema for OLAP applications, and used as **Data Mart** for department-level DW

- The star schema is simple <u>but some **redundancy** may occur for dimension tables</u>
  - E.g., Location {location_key, street, city, province, country}
    (102 St…, Vancouver, British Columbia, Canada)
    (206 St…, Vancouver, British Columbia, Canada)

# 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
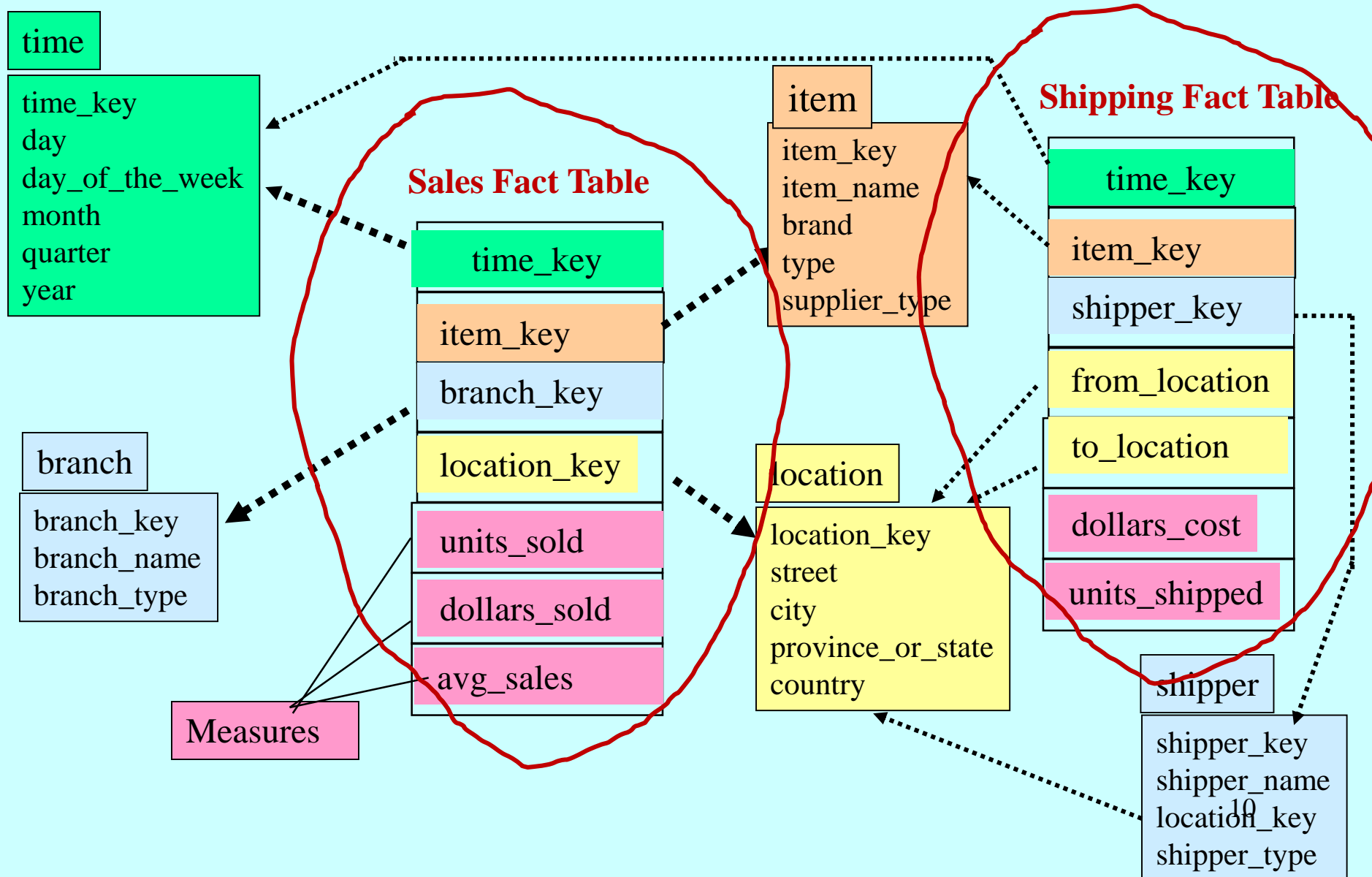- state_or_province
- country

Measures

# Example of Snowflake Schema (cont)

- A variation of star schema, in which the dimension tables are normalized

- Main purpose: saving space and for easier maintenance
    - Normalizing large dimension tables for saving storage space
    - However, it can reduce the effectiveness of browsing since more joints will be needed to execute a query
        - Keeping small dimension tables as it is for reducing the cost.
        - Performance degradation because of join operation on multiple tables.

# 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**

time_key

item_key

shipper_key

from_location

to_location

dollars_cost

units_shipped

**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

**shipper**

shipper_key
shipper_name
location_key
shipper_type

# Example of Fact Constellation (cont)

- Fact constellation schema is for sophisticate DW
  - It is for the design of a DW with <u>multiple subjects,</u> such as for a large corporation which needs information for quickly updating big pictures of entire organization, etc.

# DW Schema Definition Language

- Cube Definition (Fact Table)

  **define cube** <cube_name> [<dimension_list>]:
  <measure_list>

- Dimension Definition ( Dimension Table )

  **define dimension** <dimension_name> **as**
  (<attribute_or_subdimension_list>)

- Special Case (Shared Dimension Tables)
  - First time as "cube definition"
  - **define dimension** <dimension_name> **as**
    <dimension_name_first_time> **in cube**
    <cube_name_first_time>

# E.g., Define Star Schema for "Sales":

**define cube** sales_star [time, item, branch, location]:

> dollars_sold = sum(sales_in_dollars),
>
> avg_sales = avg(sales_in_dollars),
>
> units_sold = count(*)

**define dimension** time **as** (time_key, day, day_of_week, month, quarter, year)

**define dimension** item **as** (item_key, item_name, brand, type, supplier_type)

**define dimension** branch **as** (branch_key, branch_name, branch_type)

**define dimension** location **as** (location_key, street, city, province_or_state, country)

# E.g. Define Snowflake Schema for "Sales"

**define cube** sales_snowflake [time, item, branch, location]:

      dollars_sold = sum(sales_in_dollars),

      avg_sales = avg(sales_in_dollars),

      units_sold = count(*)

**define dimension** time **as** (time_key, day, day_of_week, month, quarter, year)

**define dimension** Item **as** (item_key, item_name, brand, type, supplier(supplier_key, supplier_type))

**define dimension** branch **as** (branch_key, branch_name, branch_type)

**define dimension** location **as** (location_key, street, city(city_key, province_or_state, country))

Embedded table

# Define Fact Constellation Schema for "Sales" and "Shipping"

**define cube** Sales [time, item, branch, location]:

        dollars_sold = sum(sales_in_dollars), avg_sales = avg(sales_in_dollars), units_sold = count(*)

**define dimension** Time **as** (time_key, day, day_of_week, month, quarter, year)

**define dimension** Item **as** (item_key, item_name, brand, type, supplier_type)

**define dimension** Branch **as** (branch_key, branch_name, branch_type)

**define dimension** Location **as** (location_key, street, city, province_or_state, country)

**define cube** Shipping [time, item, shipper, from_location, to_location]:

        dollar_cost = sum(cost_in_dollars), unit_shipped = count(*)

**define dimension** Time **as** time in cube sales

**define dimension** Item **as** item in cube sales

**define dimension** Shipper **as** (shipper_key, shipper_name, location **as** location in cube sales, shipper_type)

**define dimension** from_location **as** location in cube sales

**define dimension** to_location **as** location in cube sales

# E.g., A DW instance: How data is materialized?

**STORE LOOKUP (dimension table)**

| Store ID | Store | Region | Company |
|---|---|---|---|
| 1 | Ridgewood | Northeast | B&B |
| 2 | Newbury | Northeast | B&B |
| 3 | Avon | Northeast | B&B |
| 4 | Francis | Midwest | B&B |
| 5 | Nikki's | Midwest | B&B |
| 6 | Roger's | Midwest | B&B |

**TIME LOOKUP (dimension table)**

| Month ID | Month | Quarter |
|---|---|---|
| 1 | Jan | 1 |
| 2 | Feb | 1 |
| 3 | Mar | 1 |
| 4 | Apr | 2 |
| 5 | May | 2 |
| 6 | Jun | 2 |

**PRODUCT LOOKUP (dimension table)**

| Prod. ID | Prod. name | Prod. type | Prod. Group |
|---|---|---|---|
| 1 | rosewater soap | soap | skin care |
| 2 | olive oil soap | soap | skin care |
| 3 | hypoaller. lotion | lotion | skin care |
| 4 | bookshelves | office | furniture |
| 5 | dividers | office | furniture |
| 6 | mattresses | home | furniture |

**BASE SALES DATA (fact table)**

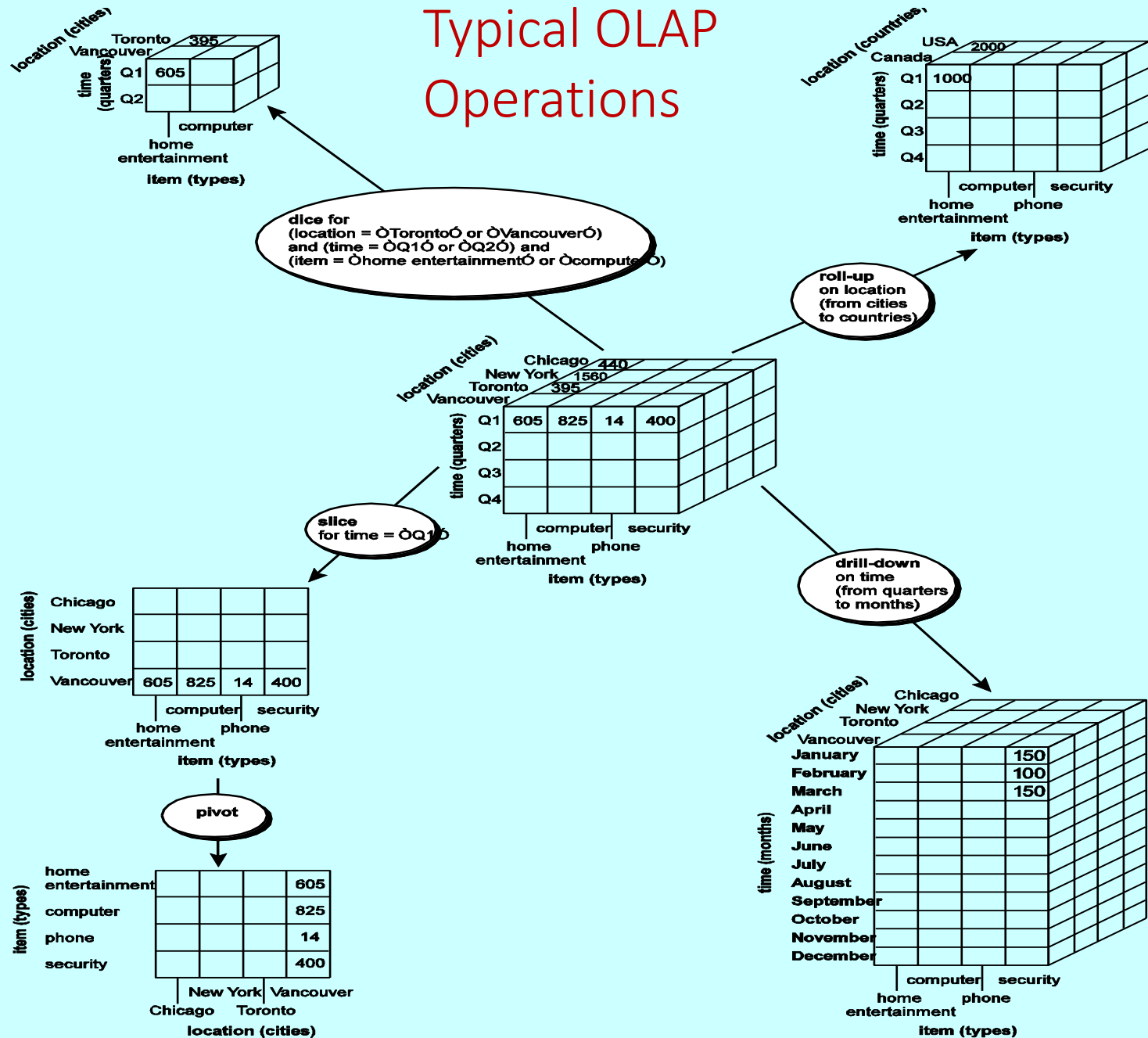| Store id | Month id | Prod. id | Scenario | Sales | Costs |
|---|---|---|---|---|---|
| 1 | 1 | 1 | actuals | 285 | 240 |
| 1 | 1 | 1 | plans | 280 | 230 |
| 1 | 1 | 2 | actuals | 270 | 260 |
| 1 | 1 | 2 | plans | 265 | 255 |
| 1 | 1 | 3 | actuals | 350 | 300 |
| 1 | 1 | 3 | plans | 300 | 280 |
| 1 | 1 | 4 | actuals | 220 | 230 |
| 1 | 1 | 4 | plans | 230 | 235 |
| 1 | 1 | 5 | actuals | 480 | 400 |
| 1 | 1 | 5 | plans | 450 | 380 |
| 1 | 1 | 6 | actuals | 380 | 370 |
| 1 | 1 | 6 | plans | 390 | 375 |
| 1 | 2 | 1 | actuals | 313 | 264 |
| 1 | 2 | 1 | plans | 308 | 253 |
| : | : | : | : | : | : |
| : | : | : | : | : | : |
| 6 | 12 | 6 | actuals | 1,199.28 | 1,168.14 |
| 6 | 12 | 6 | plans | 1,230.42 | 1,183.71 |

Measures

16

# OLAP Operations

- ## On-Line Analytical Processing (OLAP)

  - It is an approach to answering multi-dimensional analytical (MDA) queries swiftly in computing.

- ## OLAP Operators

  - They are the tools enable users to analyze multidimensional data interactively from multiple perspectives, consisting of four basic analytical operations: **roll-up** (consolidation), **drill-down**, and **slicing** and **dicing**
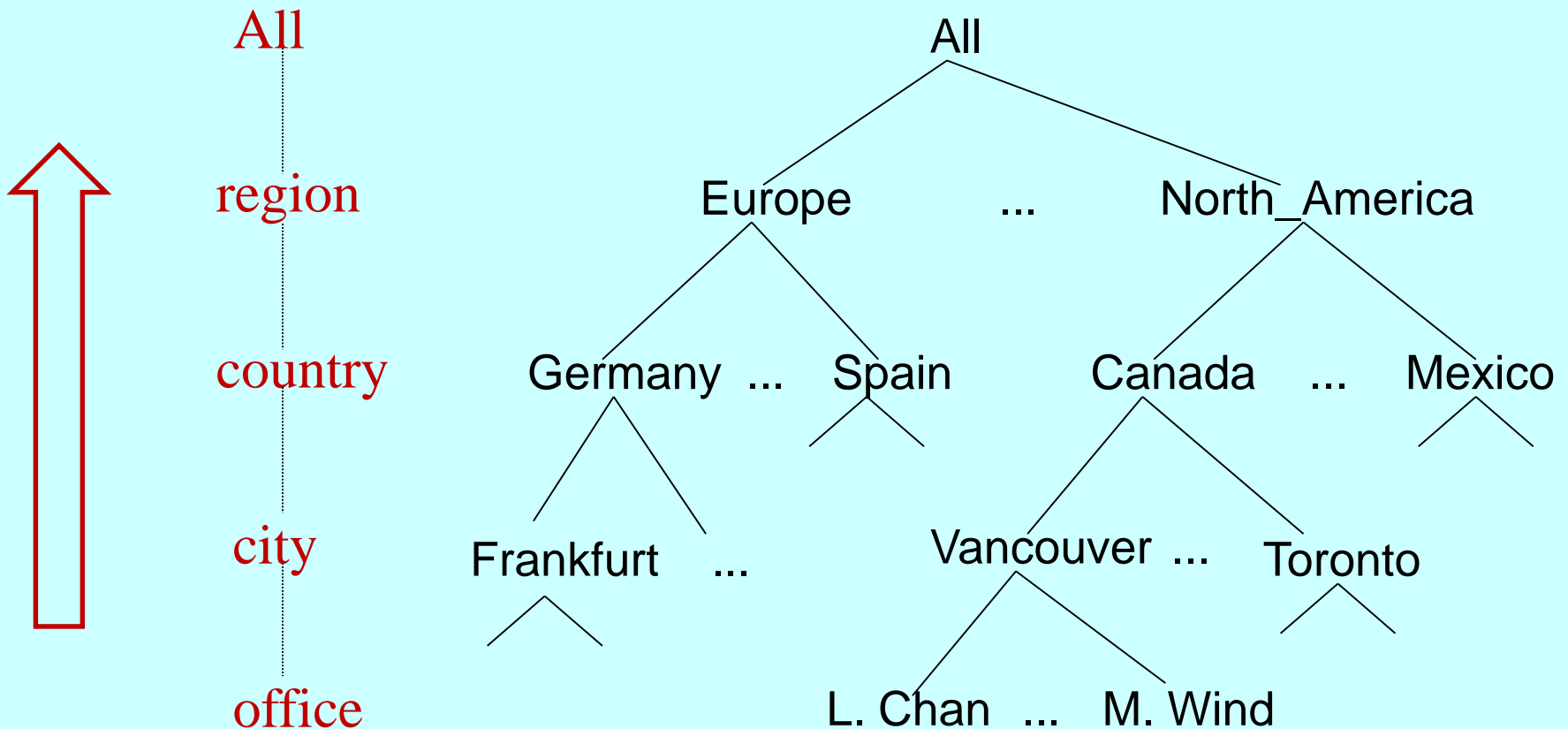
# OLAP Operations (cont)

- **Roll up:** (summarize data)
  - By climbing up hierarchy or by dimension reduction
- **Drill down:** (reverse of roll-up)
  - From higher level summary to lower level summary or detailed data, or introducing new dimensions
- **Slice:** (project)
  - By choosing a single value for one of its dimensions, creating a new cube with one fewer dimension
- **Dice:** (select)
  - By produces a subcube by allowing the analyst to pick specific values of multiple dimensions
- **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

# Aggregation by **rolling-up** along concept hierarchy

All

region

country

city

office

All

Europe ... North_America

Germany ... Spain Canada ... Mexico

Frankfurt ... Vancouver ... Toronto

L. Chan ... M. Wind

# Aggregation by **rolling-up** (or slicing by choosing "all") along dimensional lattice



**all**

0-D(apex) cuboid

product        time        country

1-D cuboids

product, time    product, country    time, country

2-D cuboids

product, time, country

3-D(base) cuboid

21

# Cube Aggregation by **roll-up**:
E.g., <u>Store</u> C sold # of <u>P</u> in <u>week</u> X.



How to compute sums?

**Week 2**

|  | c 1 | c 2 | c 3 |
|---|---|---|---|
| p 1 | 4 4 | 4 |  |

**Week 1**

|  | c 1 | c 2 | c 3 |
|---|---|---|---|
| p 1 | 1 2 |  | 5 0 |
| p 2 | 1 1 | 8 |  |

Along week

|  | c 1 | c 2 | c 3 |
|---|---|---|---|
| p 1 | 5 6 | 4 | 5 0 |
| p 2 | 1 1 | 8 |  |

Along P

|  | c1 | c2 | c3 |
|---|---|---|---|
| sum | 67 | 12 | 50 |

Along C

|  | sum |
|---|---|
| p1 | 110 |
| p2 | 19 |

129

**All, All, All**

⟶ rollup ⟶

⟵ drill-down ⟵

22

# Specify Roll Up Operations



**sale(c1,*,*)**

**sale(c2,p2,*)**

**sale(*,p2,*)**

**sale(*,*,*)**

**\* = All**

# Extended Cube

|   | c1 | c2 | c3 | * |
|---|----|----|----|----|
| p1 | 56 | 4 | 50 | 110 |
| p2 | 11 | 8 |  | 19 |
|   |   |   |   | 129 |

**Week 2**

|   | c1 | c2 | c3 | * |
|---|----|----|----|----|
| p1 | 44 | 4 |  | 48 |
|   |   |   |   | 48 |

**Week 1**

|   | c1 | c2 | c3 | * |
|---|----|----|----|----|
| p1 | 12 |  | 50 | 62 |
| p2 | 11 | 8 |  | 19 |
| * | 23 | 8 | 50 | 81 |

**sale(\*,p2,\*)**

# Data cube: **sales** (time, product, location)



**Time**

Product

1Qtr   2Qtr   3Qtr   4Qtr   *sum*

TV
PC
VCR
*sum*

Location

U.S.A

Canada

Mexico

*sum*

**Total annual sales of TV in U.S.A.**

All, All, All

25

# OLAP Query: A Star-Net Model



Customer Orders

Shipping Method

Customer

CONTRACTS

AIR-EXPRESS

TRUCK

ORDER

PRODUCT LINE

Time

Product

ANNUALY  QTRLY  DAILY

PRODUCT ITEM

PRODUCT GROUP

CITY

SALES PERSON

COUNTRY

DISTRICT

REGION

DIVISION

Location

Each circle is an abstraction level.

Promotion

Organization

26

# Ass4 hint on OLAP query generation:

1.  Each application query is an ad hoc business question (in English) about some specific analytical information on the subject.

2.  The English query is then translated into OLAP operation(s), each operation is to apply one of the four OLAP operators, defined by choosing appropriate dimensions/concepts/values for producing a report.

3.  Each report, i.e. the retrieved information, is represented in an analytical screen form (or in a graphical form).

# A DW Project Example for FCS Student Grade Analysis (Doc/Theses/MACSprojOu03.pdf)

- How to quickly get a <u>big picture of grades </u>in terms of *courses, terms, years, subject areas, student groups,* etc.

- How to quickly analyze the information stored in multiple data sources for answering complex queries in improving academic process management.

# Ad hoc analytical query examples

1. To answer why "Between 25-30 % of undergraduate students were dismissed from the Faculty of Computer Science each year over the past four years; most were first year students?" - *CS Faculty Retreat Report, 2002*

2. Do students fail on a particular course more than on other courses? Why?

3. If students fail on one course, will they also tend to fail on other particular courses?

4. Do students have better performances on the courses of a specific area than on the courses of other areas?

5. Is the class enrollment a factor affecting the students' grades?

**Challenges:** 1) How to get all data comparisons on all courses for a particular year, e.g. 2002-2003, and see the trend of the grades for different year levels? 2) How to handle multiple data sets?

# E.g., Multiple data sources: 11 CS student datasets.

| Name of the Source File | Format | Attribute Description | Time Base | Information | Sample |
|---|---|---|---|---|---|
| Students_200220.dat, Students_200230.dat, Students_200310.dat | Flat files | "BANNER, GENDER, NATIONALITY, AREA, PROVINCE, HIGHSCHOOL, DEGREE, ADMISSION AVERAGE" | 2001/2002 Winter, Summer and 2002/2003 Fall | Student biographical information | "B00XXXXXX, M, C, Halifax County, NS, NS3498, BSC, 94" |
| Grades_200220.dat, Grades_200230.dat, Grades_200310.dat | Flat files | "BANNER, DEGREE, MAJOR, TERM, SUBJECT, COURSE, GRADE" | 2001/2002 Winter, Summer and 2002/2003 Fall | Student grade information | "B00XXXXXX, BCSC, CSCI, 200310, CSCI, 2110, A+" |
| Courses_200220.xls, Courses_200230.xls, Courses_200310.xls | Excel spreadsheets | "TERM, SUBJECT, COURSE, SECTION, ENROLLMENT" | 2001/2002 Winter, Summer and 2002/2003 Fall | Course information | "200310, ANAT, 2160, 01, 41" |
| Csci_200010_200220_courses.xls | Excel spreadsheet | "TERM, TERM DESC, CRN, SUBJECT, CRSE NUMVER, SECTION, CREDIT HOURS, ENROLLMENT, INSTRUCTOR ID, INSTRUCTOR LAST NAME, INSTRUCTOR FIRST NAME" | 1999/2000 Fall ~ 2001/2002 Winter | Course information | "200010, 1999/2000 Fall, 10967, CSCI, 1100, 01, 3, 91, B00XXXXXX, Riordan, Denis" |
| CSCI _Averages.xls | Excel spreadsheet | "BANNER, DEGREE, AVERAGE" | 1999/2000 Fall ~ 2001/2002 Winter | Student average information | "B00XXXXXX, BCSC, 98" |

Table 3.2.1-1 Source Data Files from the Registrar's Office

# Solution: DW + OLAP + Visualization

- CS student data mart (star schema)

**Query 1:** *"Get a picture of the grade trends in terms of different levels of courses, and the grade distribution for year 2002-2003".*
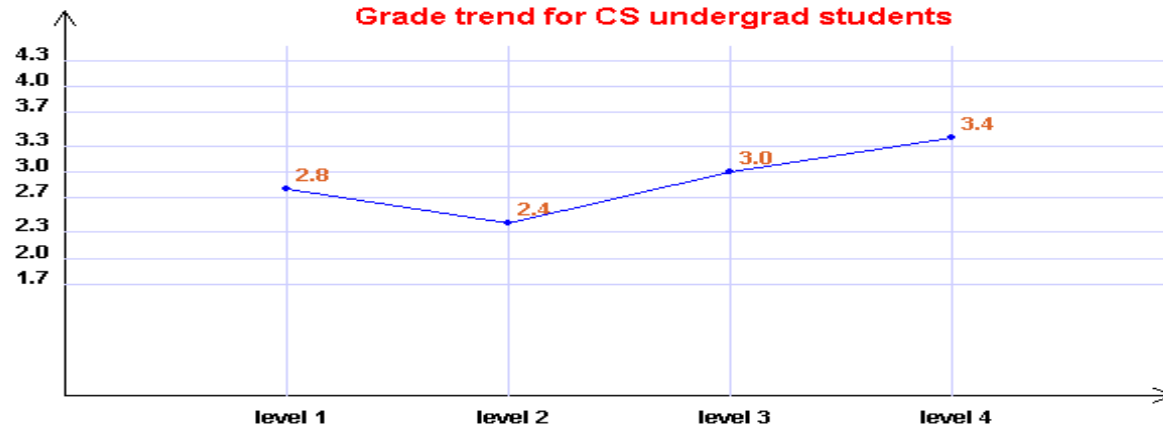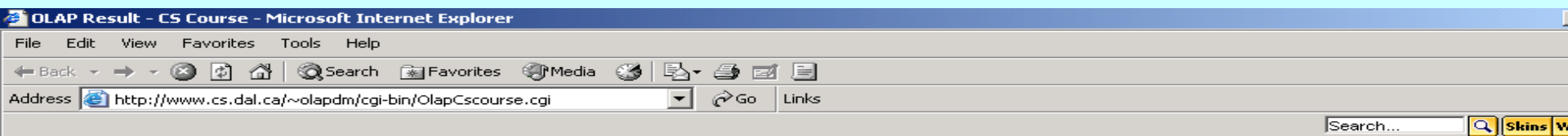
**OLAP:** **year=**2002-2003, **course=** All

# Query 2: *"Compare grade distributions between the second year courses and all other courses data for year 2002/2003".*

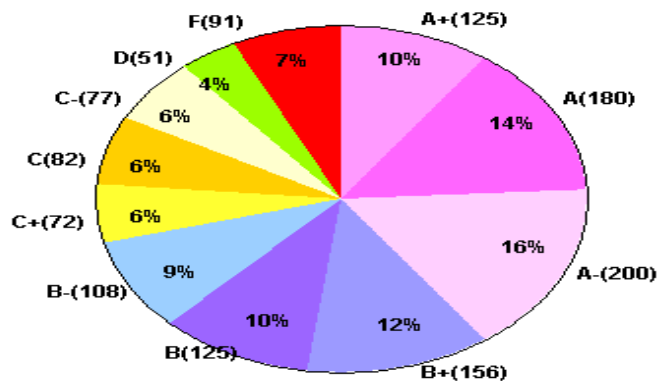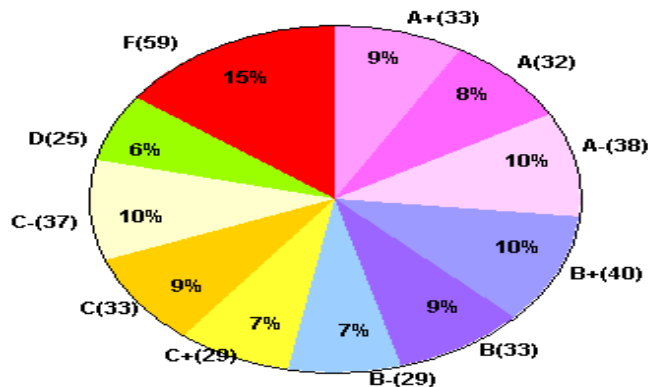## OLAP:  year = 2002-2003, course = All, courses = All level 2
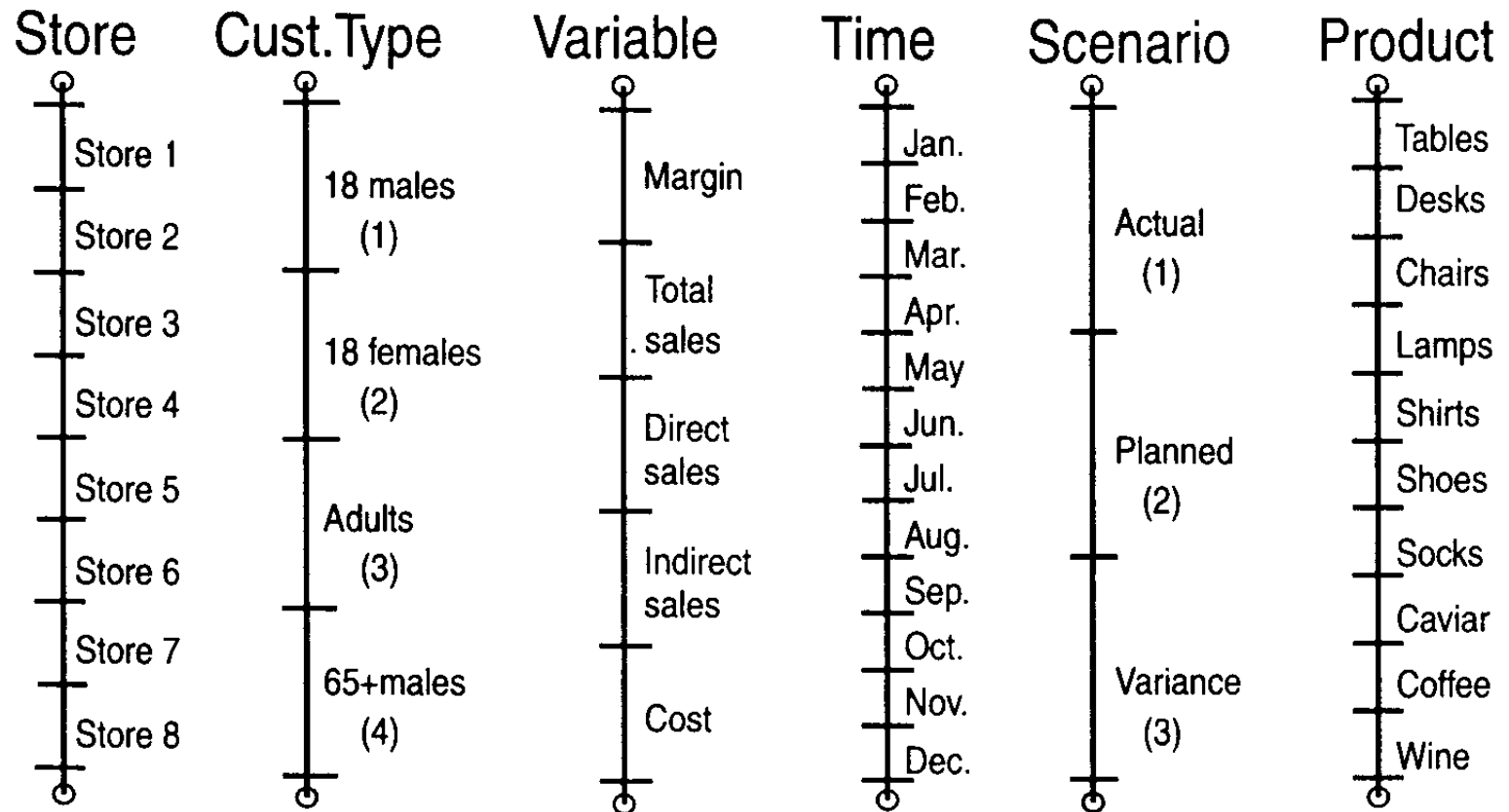
# DW/OLAP Research Project Examples

1.  Tayyaba Sharif, MACS project: **"Integrated data cube for Web content accessing pattern analysis",** 2006

    –   Doc/Theses/MACSprojTayyaba06.pdf

2.  Nariman Amiri, HINF Master thesis: **"Designing a framework of intelligent information process on dentistry administration data",** 2005

    –   Doc/Theses/MHINthesisNariman05.pdf

3.  Jie Ou, MACS project: **"Web-based OLAP and data mining for CS student database",** 2003

    –   Doc/Theses/MACSprojOu03.pdf

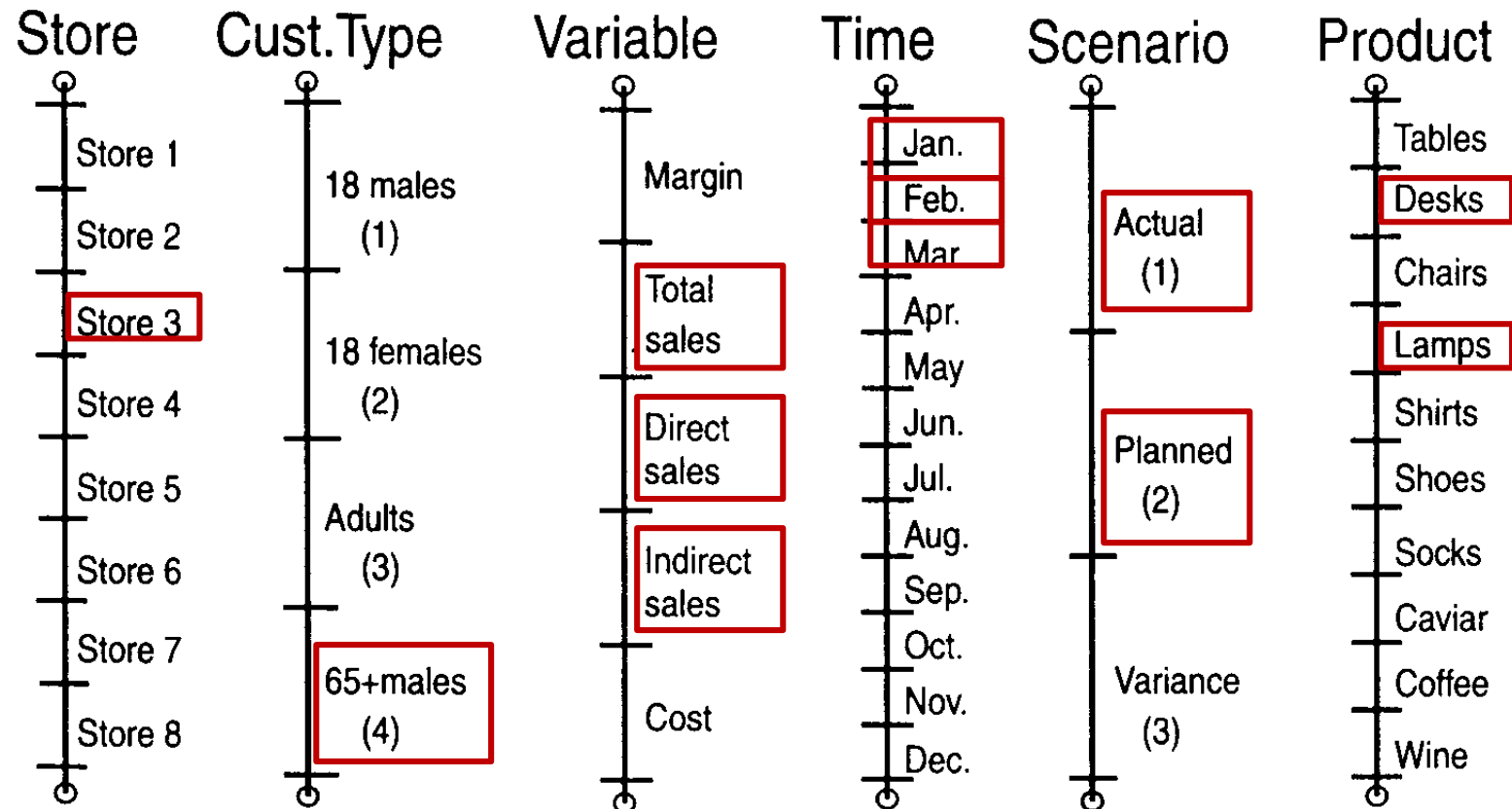# How to Present Multi-dimensional Cuboid (OLAP report) on a 2-D screen?

- How to map multiple logical dimensions onto a single computer screen?
  - Issue: Use 2-D screen to see the M-D space
- Different metaphors:
  - <u>Physical dimension metaphor</u>, e.g. 3-D graphics: **virtual camera**
  - <u>Logic dimension metaphor</u>, e.g. table-based OLAP report: **?**

| Store | Cust.Type | Variable | Time | Scenario | Product |
|-------|-----------|----------|------|----------|---------|
| Store 1 | 18 males (1) | Margin | Jan. | Actual (1) | Tables |
| Store 2 | | | Feb. | | Desks |
| Store 3 | | Total sales | Mar. | | Chairs |
| Store 4 | 18 females (2) | | Apr. | | Lamps |
| Store 5 | | Direct sales | May | Planned (2) | Shirts |
| Store 6 | Adults (3) | Indirect sales | Jun. | | Shoes |
| Store 7 | | | Jul. | | Socks |
| Store 8 | 65+males (4) | Cost | Aug. | Variance (3) | Caviar |
| | | | Sep. | | Coffee |
| | | | Oct. | | Wine |
| | | | Nov. | | |
| | | | Dec. | | |

**A six-dimensional MDS.**

**Query:** What are the sales of desk and lamp in Jan, Feb and Mar for Store 3 purchased by male senior citizens only?



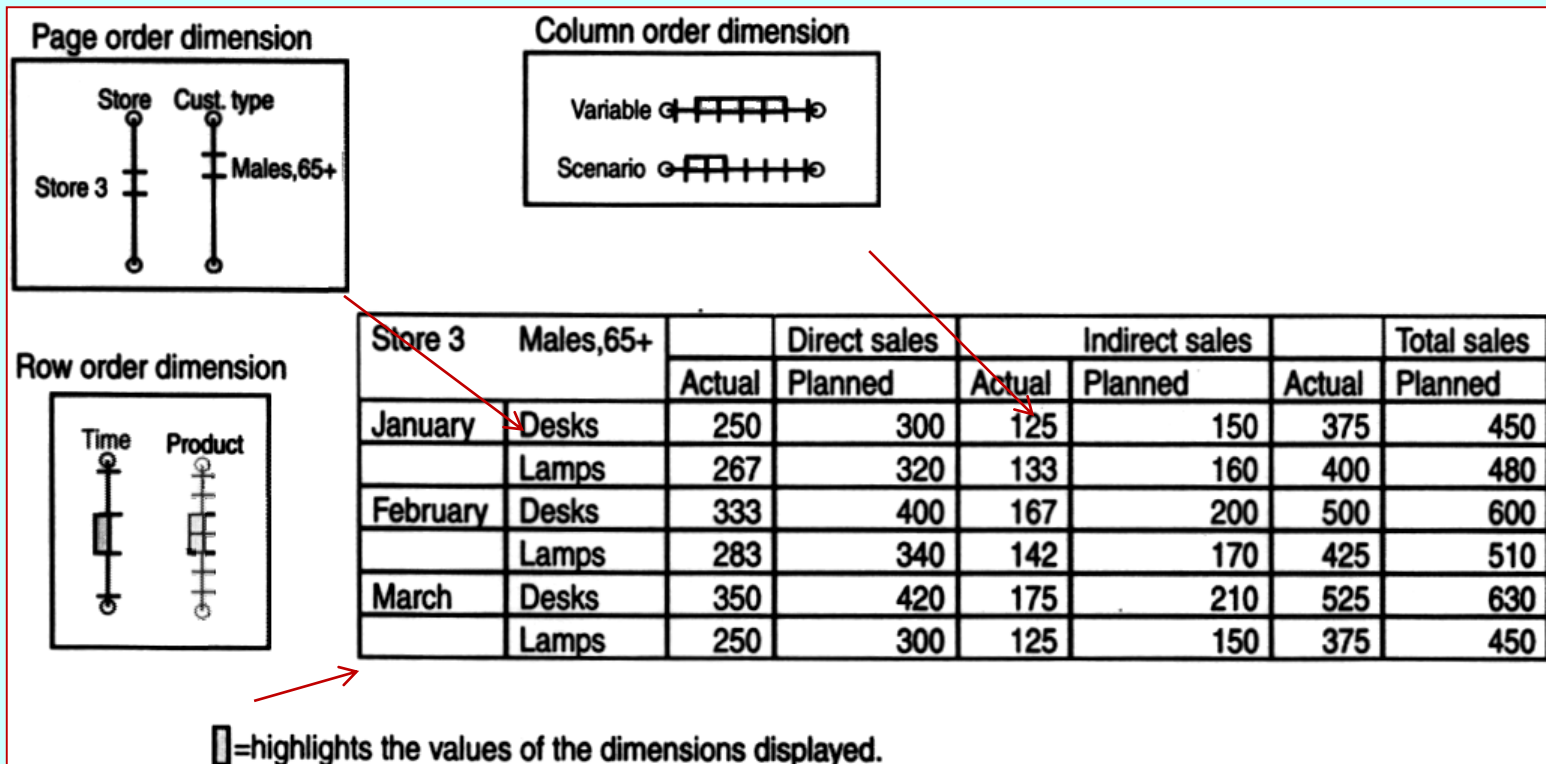A six-dimensional MDS.

# Summary of Logical Dimensions

- As distinguished from physical dimensions, which are based on angles and limited to three, logical dimensions have no such limits.

- Two types of dimensions of a data cube

  - **Identifier dimensions:** they are logical factors or identifying attributes of measurable events or things that we track.

  - **Variable dimension:** they identify what we track in a situation.

- MDS software enables <u>multi-dimensions of information to be combined onto each **row, column**, and **page** axis</u> of a screen, thus making it possible to visualize and understand a multi-dimensional data set in terms of information presented on flat screen.

- The ability of MDS software to model multidimensional information and to handle the user representation of the information makes it better suited for <u>working with complex datasets</u> than either SQL databases or traditional spreadsheets.

# Present M-D Cuboid Data in 2D Screen

- Logic dimension metaphor: (for table-based report):

  – Analytical Screen

- Solution:

  – To combine multiple logical dimensions within the same display dimensions: Row, Column, and Page

    ➢ Each dimension of the vertical bar can be connected to either a row, column, or page axes

**Query:** What are the sales of desk and lamp in Jan, Feb and Mar for Store 3 purchased by male senior citizens only?

## A six dimensional display:



Page order dimension

Column order dimension

Row order dimension

| Store 3 | Males,65+ | | Direct sales | | | Indirect sales | | | Total sales | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Actual | Planned | | Actual | Planned | | Actual | Planned |
| January | Desks | 250 | 300 | | 125 | 150 | | 375 | 450 |
| | Lamps | 267 | 320 | | 133 | 160 | | 400 | 480 |
| February | Desks | 333 | 400 | | 167 | 200 | | 500 | 600 |
| | Lamps | 283 | 340 | | 142 | 170 | | 425 | 510 |
| March | Desks | 350 | 420 | | 175 | 210 | | 525 | 630 |
| | Lamps | 250 | 300 | | 125 | 150 | | 375 | 450 |

☐=highlights the values of the dimensions displayed.

**A different six-dimensional data display of the same MDS.**

# Use analytical screen in different ways

- There are different ways that the same model dimensions can be mapped onto **row**, **column**, and **page** axes
  - The ability to easily view the same data by reconfiguring how dimensions are displayed is one of the great benefits of MDS
    - The reason is due to the separation of data structure, as represented in the vertical logical dimensions, from data display, as represented in the multi-dimensional grid.

# The issue of using analytical screen

- The more screen space is consumed for displaying dimension members, the less space is left for displaying data.

- The less space left for displaying data, the more scrolling you need to do between screens to see the same data.

- The more scrolling you need to perform, the harder it is  to understand what you are looking for.

# Make optimal use of analytical screen

- To maximize the degree to which everything on the screen is relevant, try keeping dimensions along pages <u>unless you know you need to see more than one member at a time.</u>

- Ask yourself "What do I want to look at?", or "What am I  trying to compare?" before deciding how to display information on the screen.

**Query:** Look at and compare <u>actual costs across **stores** (1, 2 and 3) and **time** (first four months)</u>, for some product (e.g. shoes) and customer type (e.g. type 2).

A six dimensional display:

Page:

product: shoes

variable: cost

scenario: actual

customer type: 2

|  | January | February | March | April |
|---|---|---|---|---|
| Store 1 | 1250 | 1700 | 1570 | 1140 |
| Store 2 | 2000 | 1950 | 1290 | 1570 |
| Store 3 | 1360 | 1580 | 1320 | 1440 |

Arranging data to compare costs across stores **and time.**

# Summarization: Aggregate Measures

- **What do I want to look at? What am I trying to compare?**
  - **Define a grouping** (i.e. determine a cuboid of the data cube)
  - **Choose measures** about the cuboid (<u>for pre-calculated values, or invoke online aggregate functions to the grouping</u>)

  - **OLAP query:** *cuboid-value* pairs (or *Dimension-value* pairs).

  E.g., "What is the total sales of computers in Halifax for the first quarter?"
  *cuboid: <time ="Q1", location ="Halfax", item ="Computer">*

  *value: sales = sum(the data set of the cuboid)*

- **Aggregate functions** are statistics models of data summarization, such as *sum, count, average, maximum, minimum, variance, standard deviation, median, mode, rank,* etc.

# Categories of Aggregate Dunctions

- **Distributive:** If the result derived by applying the function to n aggregate values is <u>the same as that derived by applying the function on all the data without partitioning</u>.
  - E.g., **count( ), sum( ), min(), max( ),** etc.
- **Algebraic:** If it can be computed by an algebraic function with the arguments which are obtained by applying distributive aggregate functions
  - E.g., **avg( )** = sum( ) / count( ), **variance( ), standard_deviation( ),** etc.
- **Holistic:** If it needs repeated search and comparison on the selected data set
  - E.g., **rank( ), median( ), mode( ),** etc.

Distributive and algebraic functions are most frequently used and suitable for on-line aggregation compuation.

# Aggregates, e.g.

- "Add up amounts for day 1"
- In SQL:  SELECT sum(amt) FROM SALE
  WHERE date = 1

| sale | prodId | storeId | date | amt |
|------|--------|---------|------|-----|
|      | p1     | c1      | 1    | 12  |
|      | p2     | c1      | 1    | 11  |
|      | p1     | c3      | 1    | 50  |
|      | p2     | c2      | 1    | 8   |
|      | p1     | c1      | 2    | 44  |
|      | p1     | c2      | 2    | 4   |

81

# Aggregates, e.g.

- "Add up amounts by day"
- In SQL:  SELECT date, sum(amt) FROM SALE
                         GROUP BY date

| sale | prodId | storeId | date | amt |
|------|--------|---------|------|-----|
|      | p1     | c1      | 1    | 12  |
|      | p2     | c1      | 1    | 11  |
|      | p1     | c3      | 1    | 50  |
|      | p2     | c2      | 1    | 8   |
|      | p1     | c1      | 2    | 44  |
|      | p1     | c2      | 2    | 4   |

| ans | date | sum |
|-----|------|-----|
|     | 1    | 81  |
|     | 2    | 48  |

# OLAP Server Architectures

- **Multi-dimensional OLAP (MOLAP)**
  - Implemented as a large multidimensional array
  - Fast indexing to pre-computed summarized data (with built-in indexing)
  - A fully materialized MOLAP array can contain an enormous number of empty cells, could result in un-acceptable storage requirements

- **Relational OLAP (ROLAP)**
  - Implemented as a collection of relational tables
  - Can be processed and queried with traditional RDBMS technology (i.e. indexing, grouping and join etc.)
  - Greater scalability
  - No "built-in" indexing

- **Hybrid OLAP (HOLAP)**
  - User flexibility, e.g., low level: relational, high-level: array
  - MS SQL Server

# Efficient Data Cube Computation

- **Review the concept of data cube lattice:**
  - Data cube can be viewed as a lattice of cuboids
    - The bottom-most cuboid is the base cuboid
    - The top-most cuboid (apex) contains only one cell
  - How many cuboids in an n-dimensional cube with L levels?

$$T = \prod_{i=1}^{n} (L_i + 1)$$

- **Materialization options of data cube**
  - Materialize <u>every</u> (cuboid) (full materialization), <u>none</u> (no materialization), or <u>some</u> (**partial materialization**)
  - Selection of which cuboids to materialize
    - Based on size, sharing, access frequency, etc.

# Multi-Tiered Architecture



Data Sources     Data Storage     OLAP Engine     Front-End Tools

# Data Warehouse Usage

Three kinds of data warehouse applications

- **Directed information processing:**
  - Deals directly with the stored aggregation information
  - Supports simple ad hoc queries, basic statistical analysis, and reporting using crosstabs, tables, charts and graphs

- **Analytical processing:**
  - supports more sophistic ad hoc queries based on a set of OLAP operations, such slice-dice, drilling, pivoting, etc.

- **Data mining:**
  - Supports user interactive exploration for finding hidden patterns
    - ➢ Visual data mining based on OLAP operations (aggregations), e.g. the trend analysis of "the wealth and health of the 200 countries over 200 years".
  - Finds hidden patterns from a defined data set in DW
  - Supports associations, constructing analytical models, performing classification and prediction, and presenting the mining results using visualization tools.

# E.g., Visual data mining based on OLAP operations (aggregations).

- Discover the trend of Wealth and Health of 200 Countries, over 200 Years:

  (http://www.youtube.com/watch?v=jbkSRLYSojo).

  In this application, by visualizing 120000 aggregation values represented in the space of income, life expectancy, and year, one can easily and clearly perceive what is the trend (the animation helps to view slices of the cube in motion along time dimension).

# From On-Line Analytical Processing to On Line Analytical Mining (OLAM)

- **Why online analytical mining?**
  - High quality of data in data warehouses
    - DW contains integrated, consistent, cleaned data
  - Available information processing structure surrounding data warehouses
    - ODBC (Open Data Base Connectivity), Web accessing, service facilities, reporting and OLAP tools
  - OLAP-based exploratory data analysis
    - mining with drilling, dicing, pivoting, etc
  - On-line selection of data mining functions
    - integration and swapping of multiple mining functions, algorithms, and tasks
- **Architecture of OLAM**

# An OLAM Architecture

**Mining query**

**Mining result**

**User Interface**

**User GUI API**

| OLAM Engine | | OLAP Engine |

**Layer3**

**OLAP/OLAM**

**Data Cube API**

**MDDB**

**Meta Data**

**Layer2**

**MDDB**

**Filtering & Integration**

**Database API**

**Filtering**

**Data cleaning**

**Databases**

**Data integration**

**Data Warehouse**

**Layer1**

**Data Repository**

55

# A case study:
## "Integrated data cube for Web content accessing pattern analysis"
(Doc/Theses/MACSprojTayyaba06.pdf)

- **Motivation**

  – Get bigger and clear pictures on web usages for supporting web site design and user group behavior analysis.

  – The results obtained by web log analyzers are limited in performance and lack depth of analysis, such as lack of content info.

- **Objective**

  – Integrate website usage and content data together and build data warehouse for conducting in-depth OLAP based analysis by generating multi-dimensional views of the comprehensive data cube.

# System Architecture

**Web Log Files**

**Web Server**

**Clean, Reduce,
Integrate, transform**

Usage data

URLs

**Data Warehouse**

Content data

**SQL Server
2000**

**Clusters
Information**

**Document
Clustering Using
WordStat**

**Retrieve
Web Pages**

**WWW**

Integrated Data
Warehouse Containing
Usage and Content
Information

**Web Pages
in Text Format**

**SQL Server 2000
Analysis Services**

**SQL server 2000
Analysis Services
Cube Browser &
MS Excel**

**Log Data Cube
With integrated
content information**

**User Interface for
Query and Analysis**

# Data Cube Design

# "Date" Dimension

# Snapshot of Dim_Date

**Data in Table 'Dim_Date' in 'logfiles' on '(local)'**

| Date_ID | The_Date | The_Year | The_Month | The_Day |
|---|---|---|---|---|
| 10904 | 01/Sep/2004 | 2004 | Sep | 1 |
| 11004 | 01/Oct/2004 | 2004 | Oct | 1 |
| 20904 | 02/Sep/2004 | 2004 | Sep | 2 |
| 21004 | 02/Oct/2004 | 2004 | Oct | 2 |
| 30904 | 03/Sep/2004 | 2004 | Sep | 3 |
| 31004 | 03/Oct/2004 | 2004 | Oct | 3 |
| 40904 | 04/Sep/2004 | 2004 | Sep | 4 |
| 41004 | 04/Oct/2004 | 2004 | Oct | 4 |
| 50904 | 05/Sep/2004 | 2004 | Sep | 5 |
| 51004 | 05/Oct/2004 | 2004 | Oct | 5 |
| 60904 | 06/Sep/2004 | 2004 | Sep | 6 |
| 61004 | 06/Oct/2004 | 2004 | Oct | 6 |
| 70904 | 07/Sep/2004 | 2004 | Sep | 7 |
| 71004 | 07/Oct/2004 | 2004 | Oct | 7 |
| 80904 | 08/Sep/2004 | 2004 | Sep | 8 |
| 81004 | 08/Oct/2004 | 2004 | Oct | 8 |
| 90904 | 09/Sep/2004 | 2004 | Sep | 9 |
| 91004 | 09/Oct/2004 | 2004 | Oct | 9 |
| 100904 | 10/Sep/2004 | 2004 | Sep | 10 |
| 101004 | 10/Oct/2004 | 2004 | Oct | 10 |
| 110904 | 11/Sep/2004 | 2004 | Sep | 11 |
| 111004 | 11/Oct/2004 | 2004 | Oct | 11 |
| 120904 | 12/Sep/2004 | 2004 | Sep | 12 |
| 121004 | 12/Oct/2004 | 2004 | Oct | 12 |
| 130904 | 13/Sep/2004 | 2004 | Sep | 13 |

# "Content" Dimension

Primary Key

**Dim_Content**
- 🔑 Content_ID
- Generalized_Cluster_Content
- Specialized_Cluster_Content
- Keywords

1. External Web Pages
2. FCS Info
3. Help
4. News
5. Research Topics
6. Seminars & CS Society Events
7. Tools & Software
8. Unavailable Pages

**Data in Table 'Dim_Content' in 'logfiles' on '(local)'**

| Content_ID | Generalized_Cluster_Content | Specialized_Cluster_Content | Keywords |
|---|---|---|---|
| 13 | Tools & Software | Eclipse and Glossary tools | Glossary, eclipse, tools, presence, reading, beneficial, universal, platform, |
| 14 | Tools & Software | Navigational Tool | Navigating, breadcrumbs, exploratory, subjects, websites, location, exper |
| 15 | Research Topics | E-Privacy | Ontology, privacy, web, PIPEDA (Personal Information Protection and Elec |
| 16 | Research Topics | Business and Association Rules | Trading, ILOG, rules, AI (artificial Intelligence), business, sampling, genetic |
| 17 | Research Topics | Speech Recognition and Healthcare | Practice, reflect, computerization, format, clinical, DAT (Dementia of Alzhei |
| 18 | Research Topics | Spatial data and ranking | Wavelet, transform, spatial, geo, CWT (Continuous Wavelet Transform), e |
| 19 | Research Topics | Data Cube | Cubes, partial, scheduling, computation, dimension, greedy, OLAP (On Line |
| 20 | Tools & Software | Teaching Aids for students | Chat, circles, engagement, instructional, features, interface, platform, vis |
| 21 | News | News Archives | Info, dean, examination, house, appearing, webmail, admission, professor |
| 22 | Seminars & CS Society Events | In-house Conference and Students Orientation | DCSI, conference, orientation, in-house, BBQ, auditorium, poster, submiss |
| 23 | Seminars & CS Society Events | Programming Competition | Competition, contest, teams, programming, APICS(Atlantic Provinces Cour |
| 24 | News | Hurricane News | Hurricane, power, afternoon, interruptions, aftermath, Juan, restored, ca |
| 25 | FCS info | CS Building and Information session | TUNS, syncrude, session, wings, houses, featured, campus, appearing, ar |
| 26 | Seminars & CS Society Events | Christmas party | Party, Christmas, gift, children, ages, wrapped, tag, bringing, reception, f |
| 27 | FCS info | Available positions and programs | Master, informatics, health, persons, referees, electronic, equity, teaching |
| 28 | FCS info | GINIus | GINIus, partners, contracts, institute, business, telecom, academic, netwo |
| 29 | FCS info | Faculty of Computer Science Index Pages | Indispensable, computer science, mentor, teachers, mission, conduct, high |
| 30 | Seminars & CS Society Events | Distinguished Speakers Series | Distinguished, speaker, ADT, watch, series, AST, engineering, radiation, N |
| 31 | News | Canadian Connectedness and Jobs | Job, jobpress, audit, sector, government, public, skills, herald, workers, br |
| 32 | Research Topics | Estimating System | Estimating, infarction, myocardial, acute, survival, AMI (Acute Myocardial I |
| 33 | FCS Info | Program and Course Description/Requirement | CSCI, bachelor, course, year, honors, requirements, fall, MWF, elective, c |
| 34 | FCS Info | Awards and Scholarships | Scholarships, award, recipients, eligibility, year, deadline, basis, study, ap |
| 35 | Help | Email and Technical Support | FAQ, account, locutus, password, mail, torch, novell, workstation, pine, us |
| 36 | Help | Help Desk and contact information | Desk, help, responsible, cshelp, policies, computing, student, center, facili |

# "HTTP Code" Dimension

Primary Key

## Dim_HTTPCode

🔑 Code_ID
HTTP_Status_Code
HTTP_Status_Desc

| HTTP Status Code | HTTP Status Code Description |
|------------------|------------------------------|
| 1XX | Information |
| 2XX | Success |
| 3XX | Redirection |

200, 206, 301, 302 and 304

# OLAP Interface

# Result Demonstration

**Query 1:** Suppose the website is redesigned in September, and the webmaster wants to know whether or not it effected the number of visitors and the ways they accessed.

# Result Demonstration (Cont.)

**Query 2:** What is the browsing distribution of different research topics in Oct and Sept, 2004?

# Result Demonstration (Cont.)

**Query 3:** How many times Web pages related to 'Photos and Web cam' are accessed by the network 141.211 in October morning which were displayed successfully?

# Summary

- **Why data warehousing?**

  - Major differences between OLAP and OLTP

- **Multi-dimensional model and data warehouse schemas**

  - Logic dimensions and space: virtual data cube (dimensions & measures)

  - Star schema, snowflake schema, fact constellations

- **OLAP operations**

  - Structures for supporting cuboids manipulation: lattice of cuboids & concept hierarchies

  - Typical operators: drilling, rolling, slicing, dicing and pivoting

  - A complex ad hoc query may be partitioned into multiple OLAP operations

- **DW and DW server implementation issue**s (only the concepts discussed in class)

  - Options of data cube materialization: full and partial materialization

  - Options of DW server for OLAP processing: ROLAP, MOLAP and HOLAP

- **Data warehouse architecture**

- From OLAP to OLAM (on-line analytical mining)

# Review Questions

1. What are the 3 typical DW schemas (describe each)?

2. What are the main differences between OLAP and OLTP process operations for answering users' queries?

3. What is the difference between logic dimensional space for OLAP analysis and the physical space for 3D computer graphics?

4. How to define OLAP queries based on Starnet query model?

5. What is the visualization metaphor for displaying OLAP result of a given multiple logic space (data cube), and how to best use it?

6. Use the demo example of "Wealth vs. Health" analysis (http://www.youtube.com/watch?v=jbkSRLYSojo) to explain how DW/OLAP may be able to support DM applications.