# Data Warehouse Modeling:
## Data Cube and OLAP

CS5483 Data Warehousing and Data Mining

# Motivation

- Suppose you want to know the sales information of a supermarket chain.
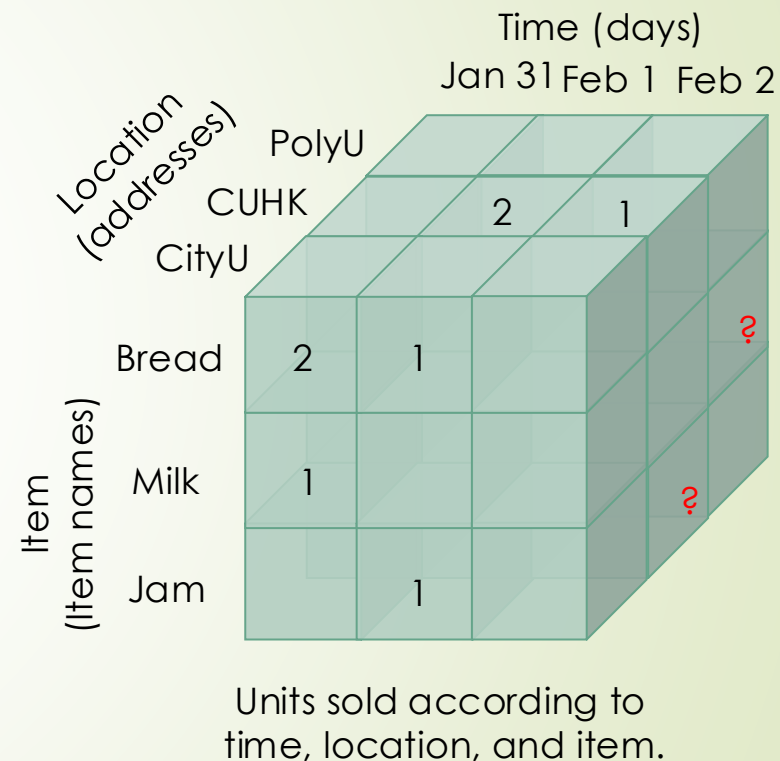
| TID | Time | Location | Item (quantity,unit price) |
|-----|------|----------|----------------------------|
| 1 | Jan 31 | CityU, Kln. | bread (2, HK$5), milk (1, HK$10) |
| 2 | Feb 1 | CUHK, N.T. | bread (2, HK$5) |
| 3 | Feb 1 | CityU, Kln. | bread (1, HK$5), jam (1, HK$5) |
| 4 | Feb 2 | CUHK, N.T. | bread (1, HK$5), jam (1, HK$5) |
| 5 | Feb 2 | PolyU, Kln. | milk (2, HK$10) |

- From the above transactional data, which store has the best sales performance?
  - A maximum of _____ units sold in _____.
  - A maximum of _____ dollars sold in _____.
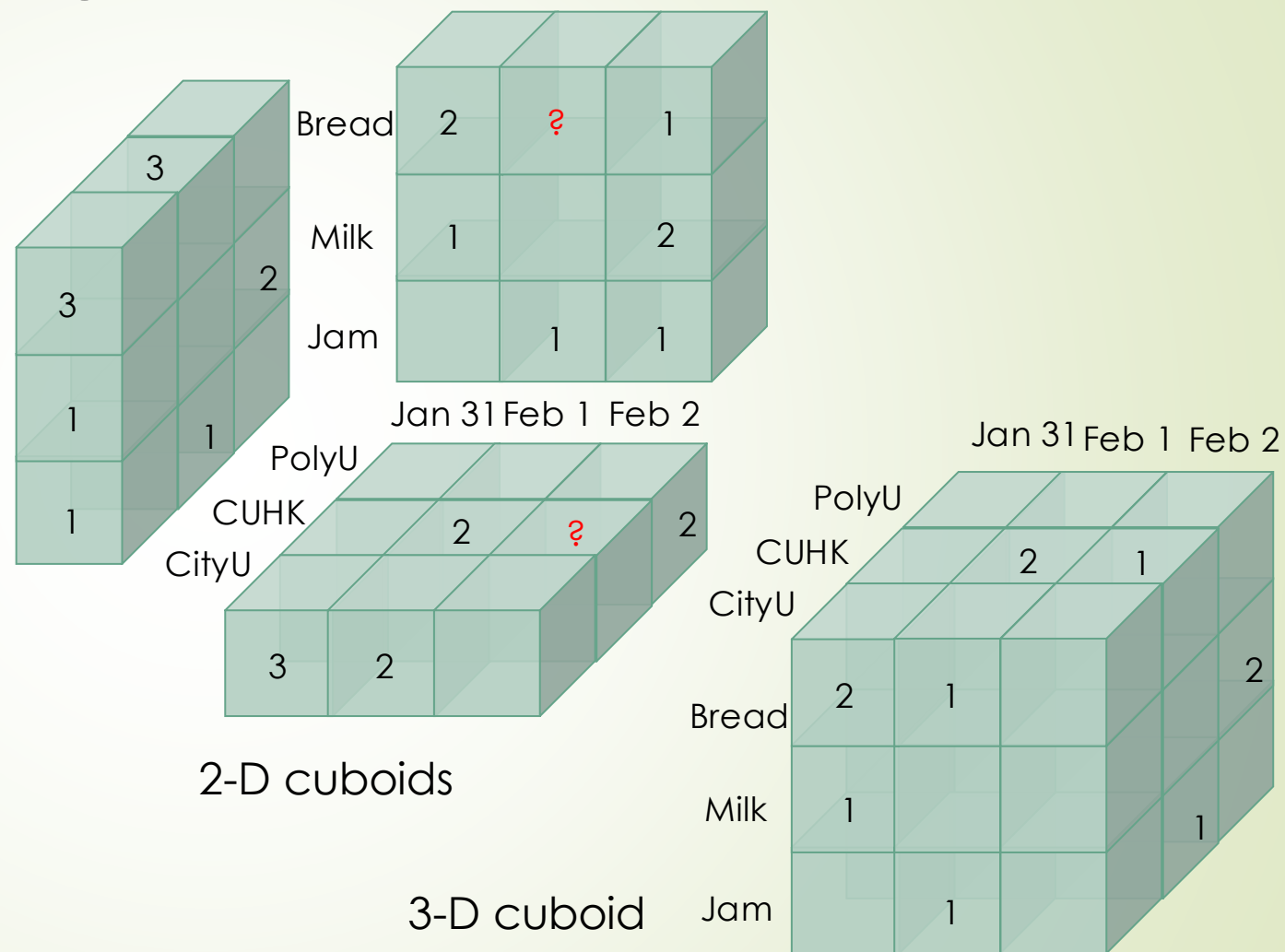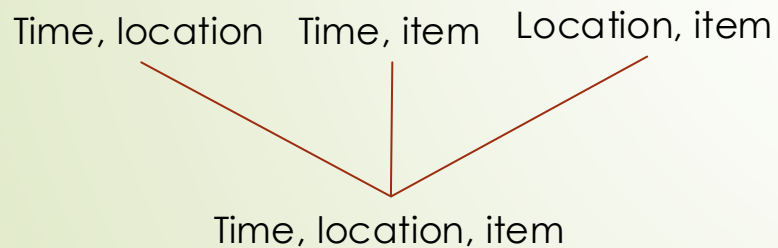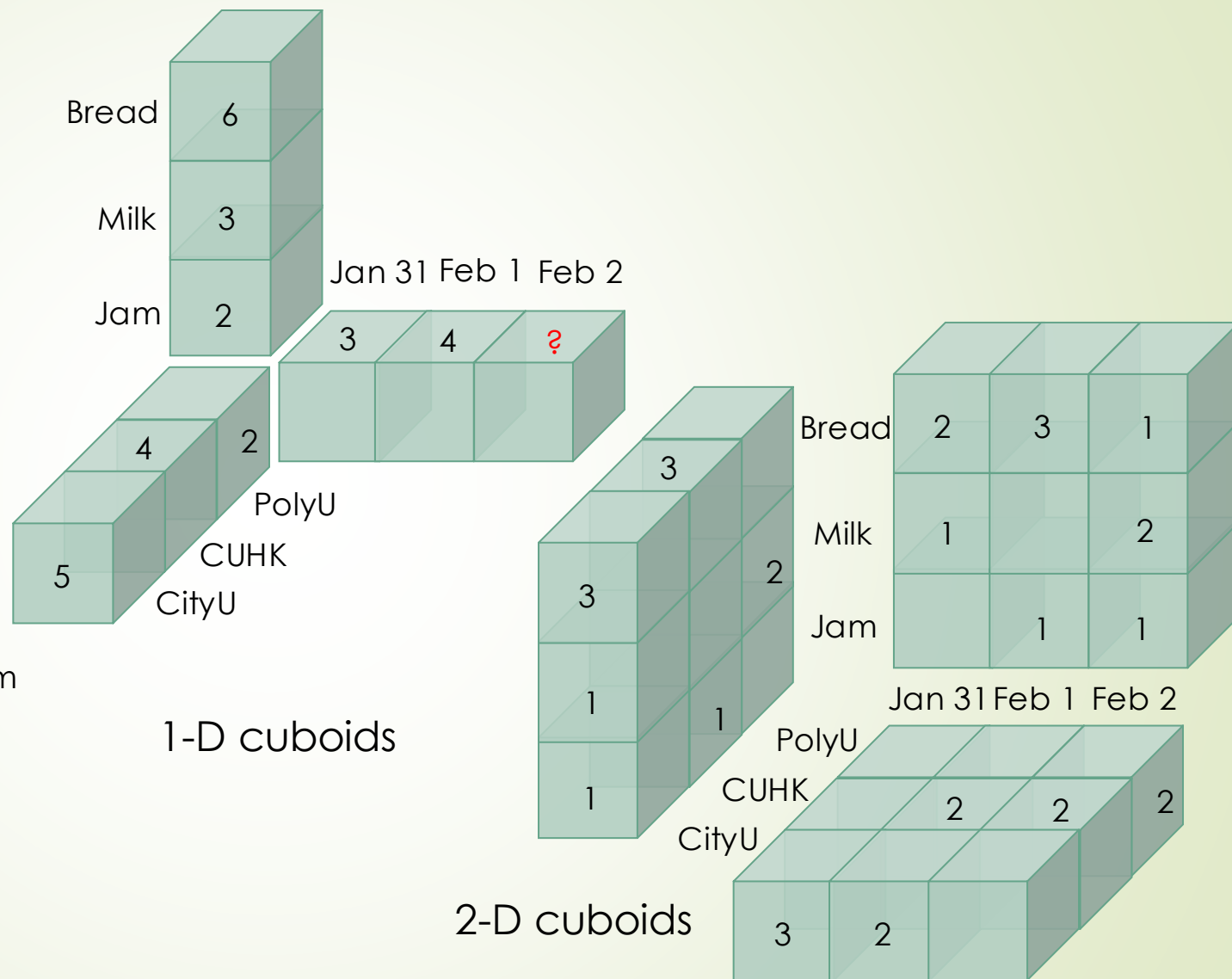- How to design a data warehouse for efficient analysis?
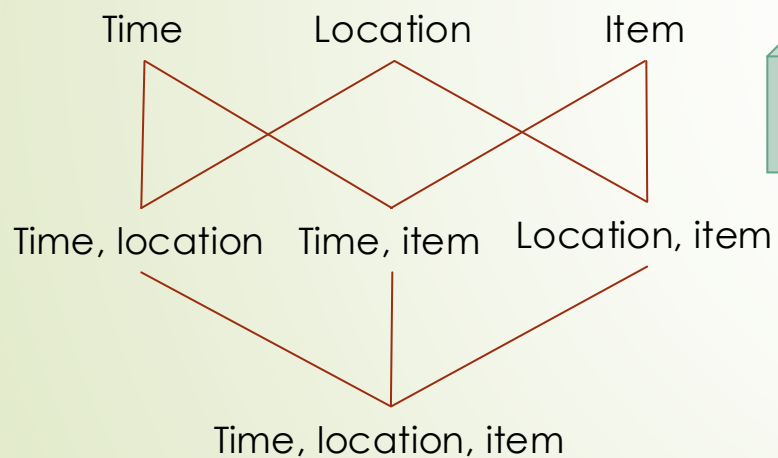
# Dimension modeling

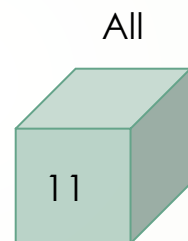| TID | Time | Location | Item (quantity,unit price) |
|-----|------|----------|----------------------------|
| 1 | Jan 31 | CityU, Kln. | bread (2, HK$5), milk (1, HK$10) |
| 2 | Feb 1 | CUHK, N.T. | bread (2, HK$5) |
| 3 | Feb 1 | CityU, Kln. | bread (1, HK$5), jam (1, HK$5) |
| 4 | Feb 2 | CUHK, N.T. | bread (1, HK$5), jam (1, HK$5) |
| 5 | Feb 2 | PolyU, Kln. | milk (2, HK$10) |

- **C**_____: Multi-dimensional array of **c**_____
  - containing the **f**_____ (units sold, dollars sold)
  - indexed using the **d**_____ (time, location, item).
- How to summarize?



Units sold according to time, location, and item.

# Dimension reduction

Time, location   Time, item   Location, item

Time, location, item

Bread   2   ?   1
Milk   1       2
Jam        1   1

Jan 31 Feb 1 Feb 2

PolyU
CUHK   2   ?   2
CityU
3   2

**2-D cuboids**

**3-D cuboid**

Jan 31 Feb 1 Feb 2

PolyU
CUHK   2   1
CityU           2
Bread   2   1
Milk   1       1
Jam        1

Bread 6
Milk 3
Jam 2

Jan 31  Feb 1  Feb 2
3    4    ?

4    2
PolyU
CUHK
5    CityU

**1-D cuboids**

Time      Location      Item

Time, location   Time, item   Location, item

Time, location, item

3
2
3    2
1    1
1    1
PolyU
CUHK
CityU

Bread  2  3  1
Milk   1     2
Jam       1  1

Jan 31 Feb 1 Feb 2

PolyU
CUHK
CityU

2  2  2
3  2

**2-D cuboids**

All

Time   Location   Item

Time, location   Time, item   Location, item

Time, location, item

All

11

0-D cuboid

Bread   6

Milk   3

Jam   2

Jan 31  Feb 1  Feb 2

3   4   4

PolyU

4   2

CUHK

CityU

5

1-D cuboids

# Lattice structure

All

11

**A____cuboid:**
Highest level of summarization

All

Time          Location          Item

Time, location   Time, item   Location, item

Time, location, item

Jan 31  Feb 1  Feb 2

PolyU
CUHK          2      1
CityU

Bread   2    1                        2

Milk    1                        1

Jam         1

**B____ cuboid:**
Lowest level of summarization

{Time, location, item}

{Time, location}  {Time, item}  {Location, item}

{Time}      {Location}      {Item}

Ø

Hasse diagram for the Boolean lattice

# Data Cube: The lattice of all cuboids



All

11

Bread 6
Milk 3
Jam 2

Jan 31 Feb 1 Feb 2
3 4 4

4 2
5
PolyU
CUHK
CityU

3
3 2
1 1
1
PolyU
CUHK
CityU

Bread 2 3 1
Milk 1 2
Jam 1 1
Jan 31 Feb 1 Feb 2

2 2 2
3 2
Jan 31 Feb 1 Feb 2
PolyU
CUHK
CityU

Jan 31 Feb 1 Feb 2
PolyU
CUHK
CityU
2 1
2

Bread 2 1
Milk 1
Jam 1
1

0-D cuboid

1-D cuboids

2-D cuboids

**Dimension reduction**

Any other operations of interest?

3-D cuboid

# Roll-up operations

- How many units of items was sold in different months and regions?

**Time**

Quarter

Month    Week

Date

**Location**

Region

Address

- **C_____ hierarchies**



Rolling up to regions

Dimension reduction

Rolling up from days to months
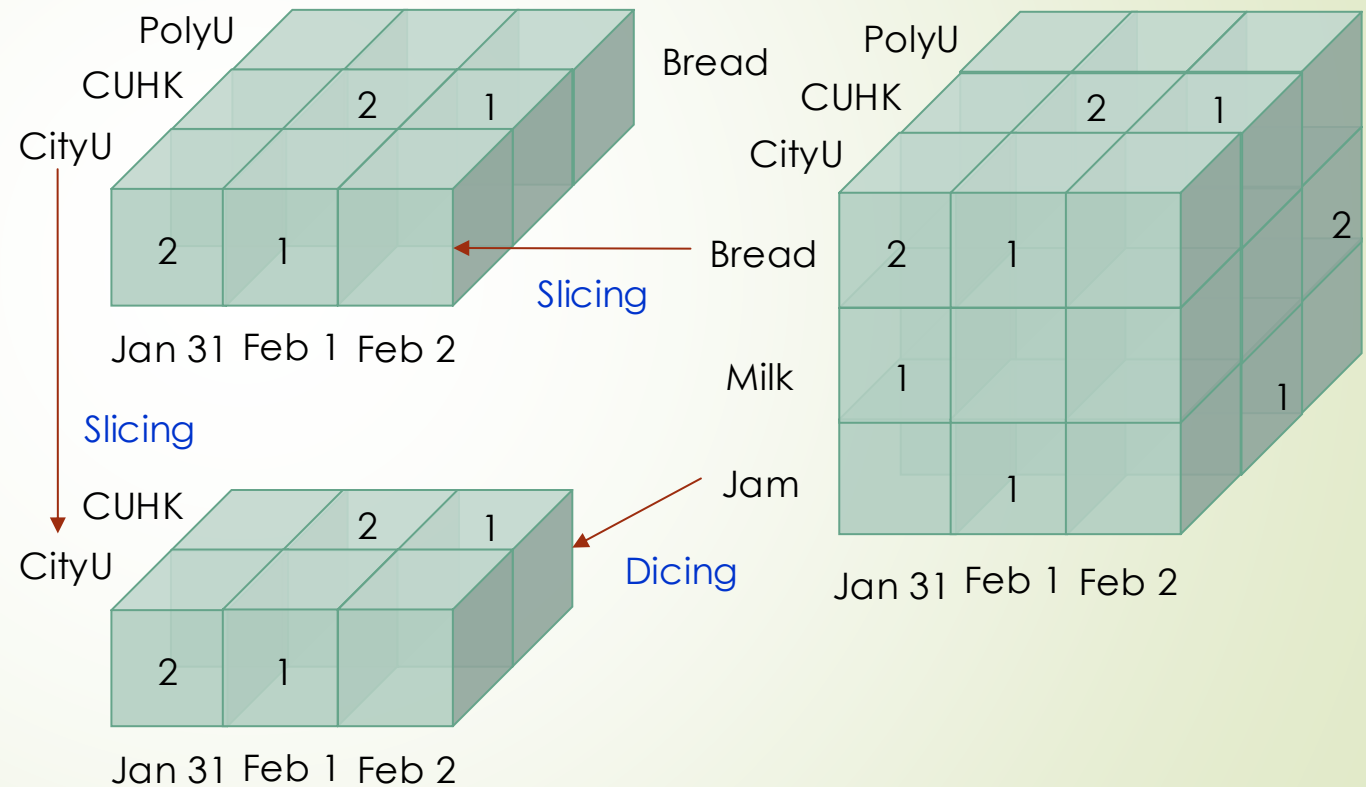
# Selection operations

- How many units of bread was sold in CityU or CUHK on different days?

- **Slice**:
  - Selection on 1 dimension.

- **Dice**:
  - Selection on multiple dimensions and/or with multiple values.

# OLAP operations

- **O_____ a_____ processing**
  - Roll-up
    - Dimension reduction
    - Climbing up a concept hierarchy
    - Reverse operation: **Drill-down**
  - Slice and Dice
    - Selection on one or more dimensions
  - Others (optional): **Pivot**/rotate, **drill-across**, **drill-through**
- In contrast with **o_____ t_____ processing (OLTP)**: insert, update, delete.

# The curse of dimensionality

- With $d$ dimensions, how many cuboids can be obtained by dimension reduction? _____

- If dimension $i$ has $L_i$ levels of concepts, how many cuboids can be obtained by the roll-up operation? _____

- How to store and compute the data cube?

  - **F_____ materialization**: Compute and store all the cuboids.

  - **N_____ materialization**: Store only the base cuboid and compute other cuboids on the fly.

  - **P_____ materialization**: Compute and store some parts of the data cube.

- How to store the base cuboid for efficient computation of other cuboids?

# References

- 4.2 Data Warehouse Modeling: Data Cube and OLAP
- Optional
  - Hands-on tutorial of ETL using Pentaho