

Data Mining:

Concepts and Techniques —

Unit II: Data Warehouse and OLAP Technology for Data Mining

UNIT – II: Data Warehouse and OLAP Technology for Data Mining Data Warehouse, Multidimensional Data Model, Data Warehouse Architecture, Data Warehouse Implementation

Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- From data warehousing to data mining
- Further Development of Data Cube Technology

What Is Data Mining?

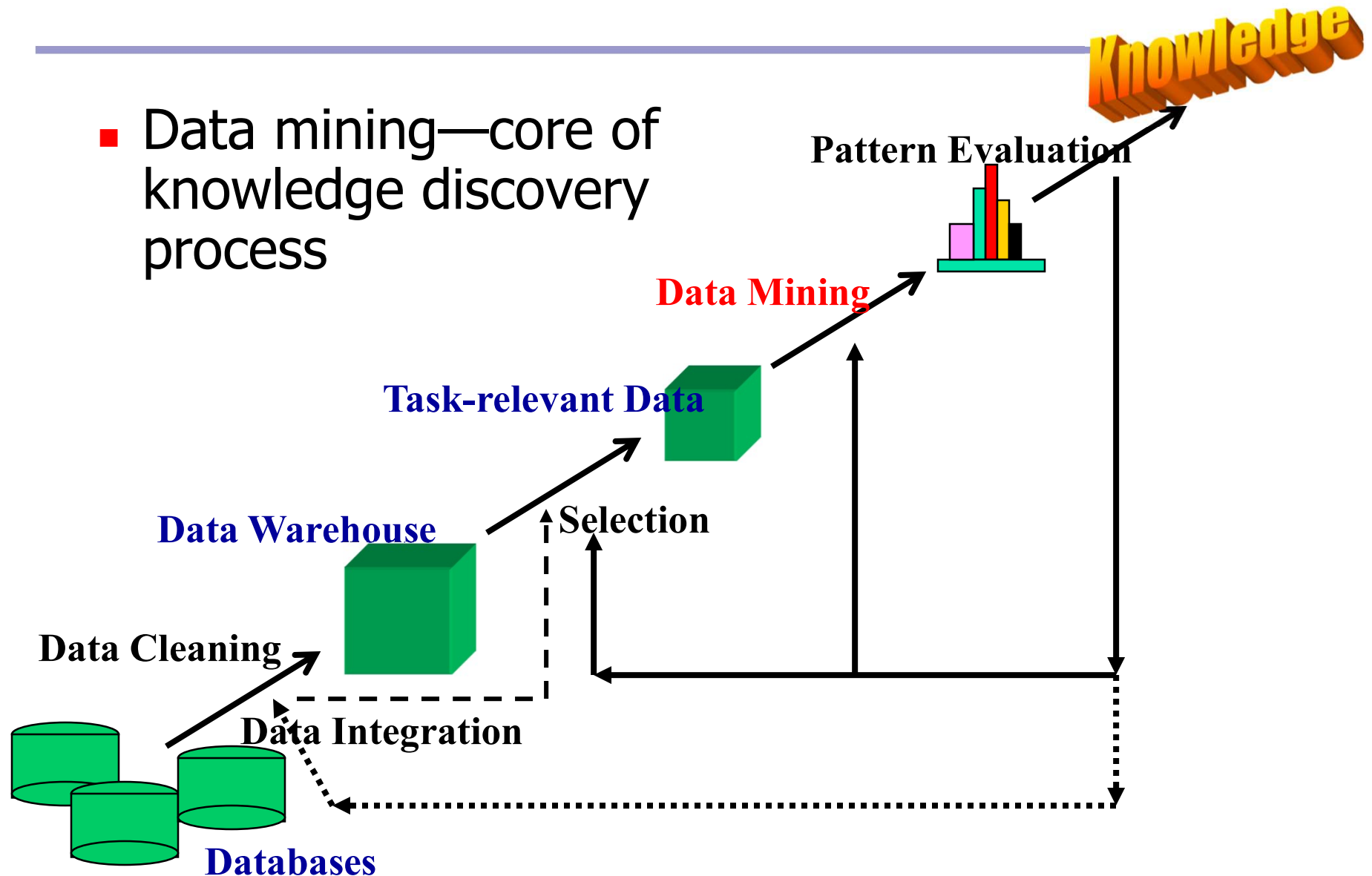


- Data mining (knowledge discovery from data)
 - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
 - Data mining: a misnomer?
- Alternative names
 - Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
- Watch out: Is everything “data mining”?
 - Simple search and query processing
 - (Deductive) expert systems

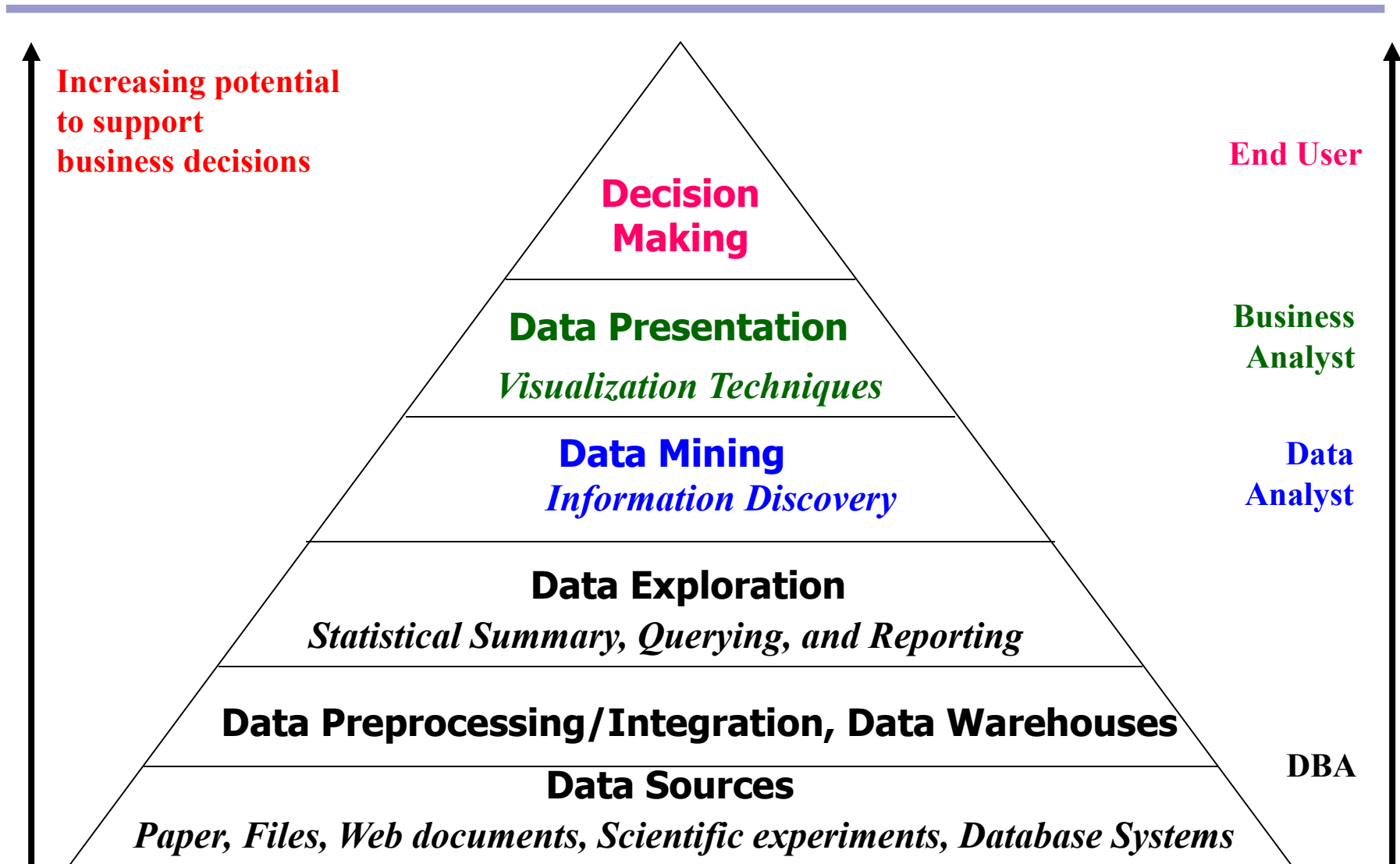


KDD Process

- Data mining—core of knowledge discovery process



Data Mining and Business Intelligence



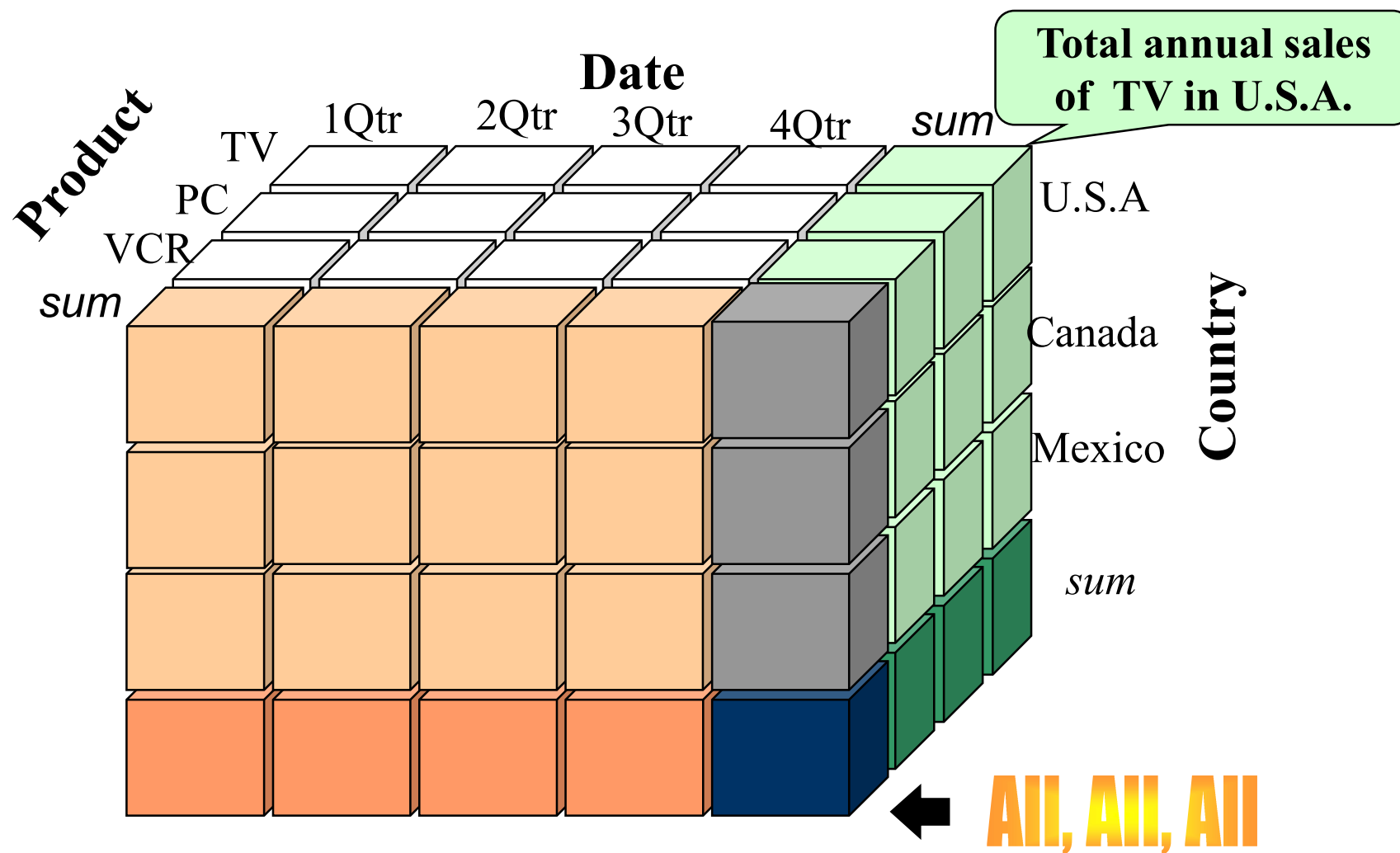
What is Data Warehouse?

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

Data Warehouse—Subject-Oriented

- Organized around major subjects, such as **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**

A Sample Data Warehouse (Data Cube) for SALES



Data Warehouse—Integrated

- Constructed by integrating **multiple, heterogeneous data sources**
 - relational databases, flat files, on-line transaction records

How?

- Data **cleaning and integration techniques** are applied.
 - **Ensure consistency in** naming conventions, encoding structures, attribute measures, etc. among different data sources
 - E.g., Hotel price: currency, tax, breakfast covered, etc.
 - 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—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*

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

5 Distinct features (OLTP vs. OLAP):

1. Data contents: current, detailed vs. historical, consolidated
2. Database design: ER + application vs. star + subject
3. View: current, local vs. evolutionary, integrated
4. Access patterns: update vs. read-only but complex queries
5. User and system orientation: customer vs. market

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, integrated, multidimensional, summarized, 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**

Why Separate Data Warehouse? (contd.)

- **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

Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- From data warehousing to data mining

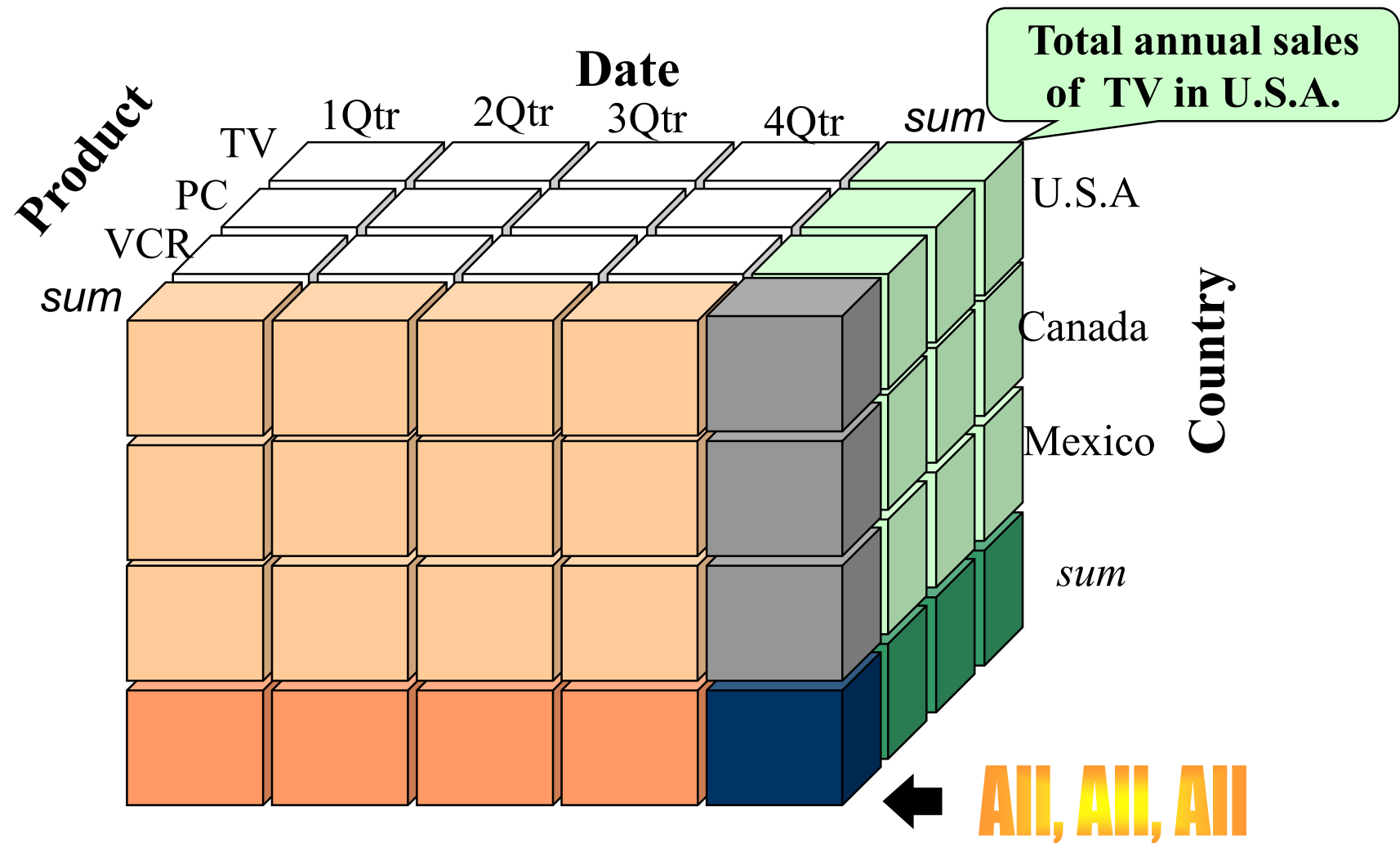
From Tables and Spreadsheets to Data Cubes

- 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

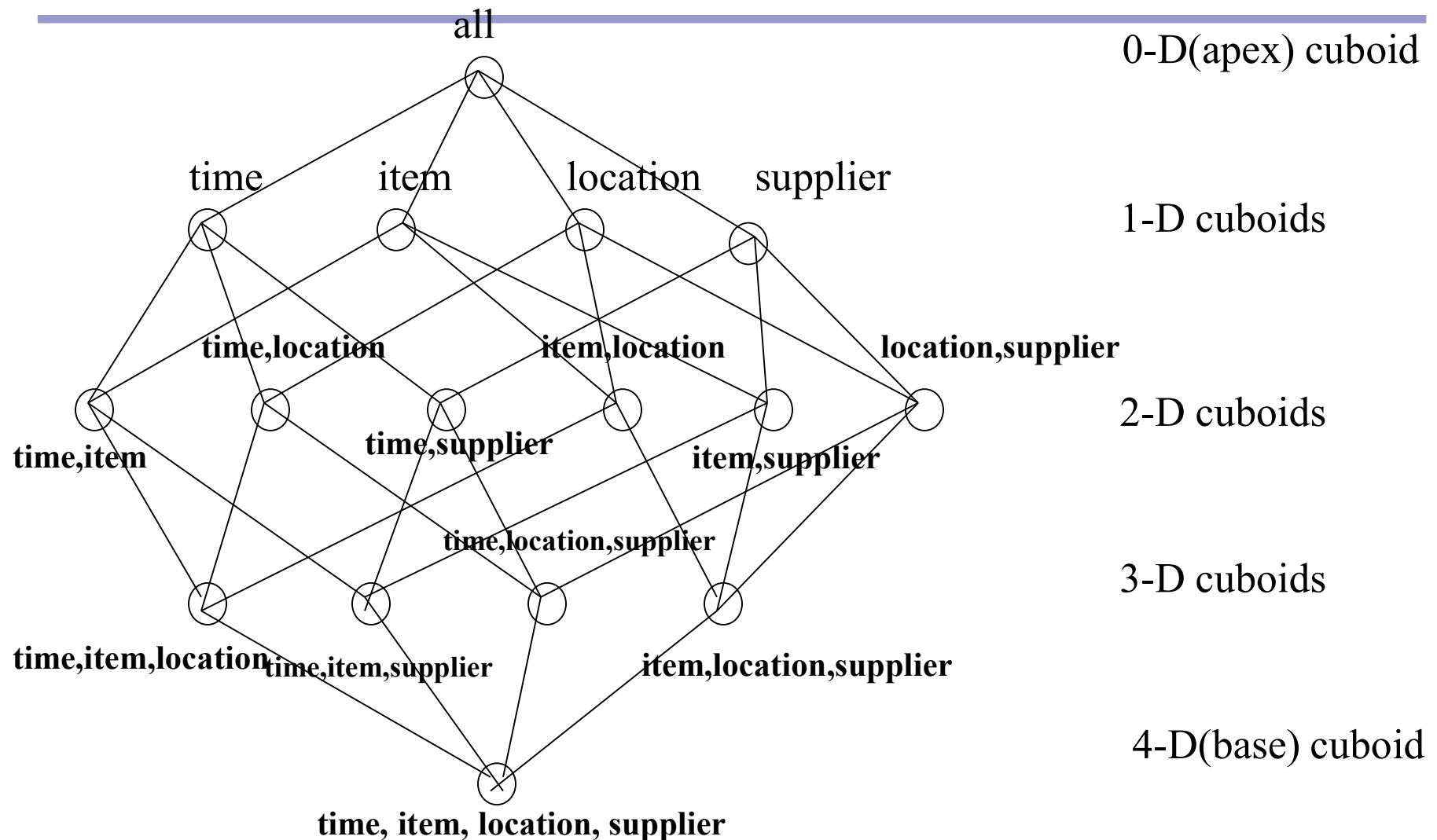
Data Cube- Cuboid Types

- 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**.

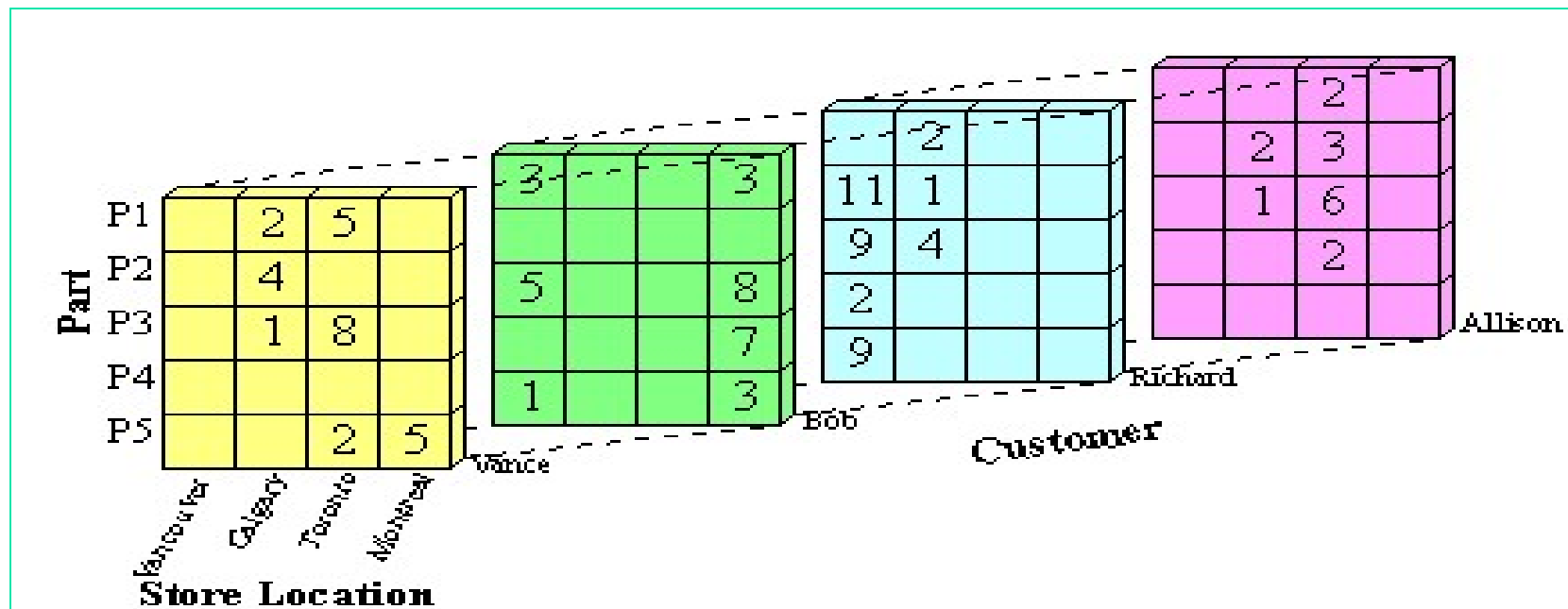
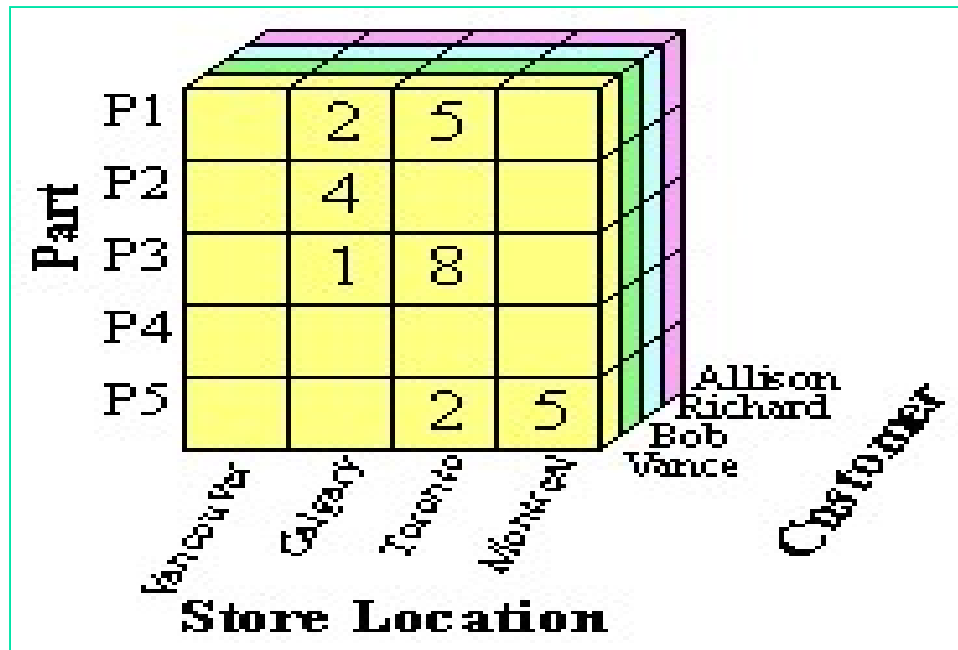
A Sample Data Cube



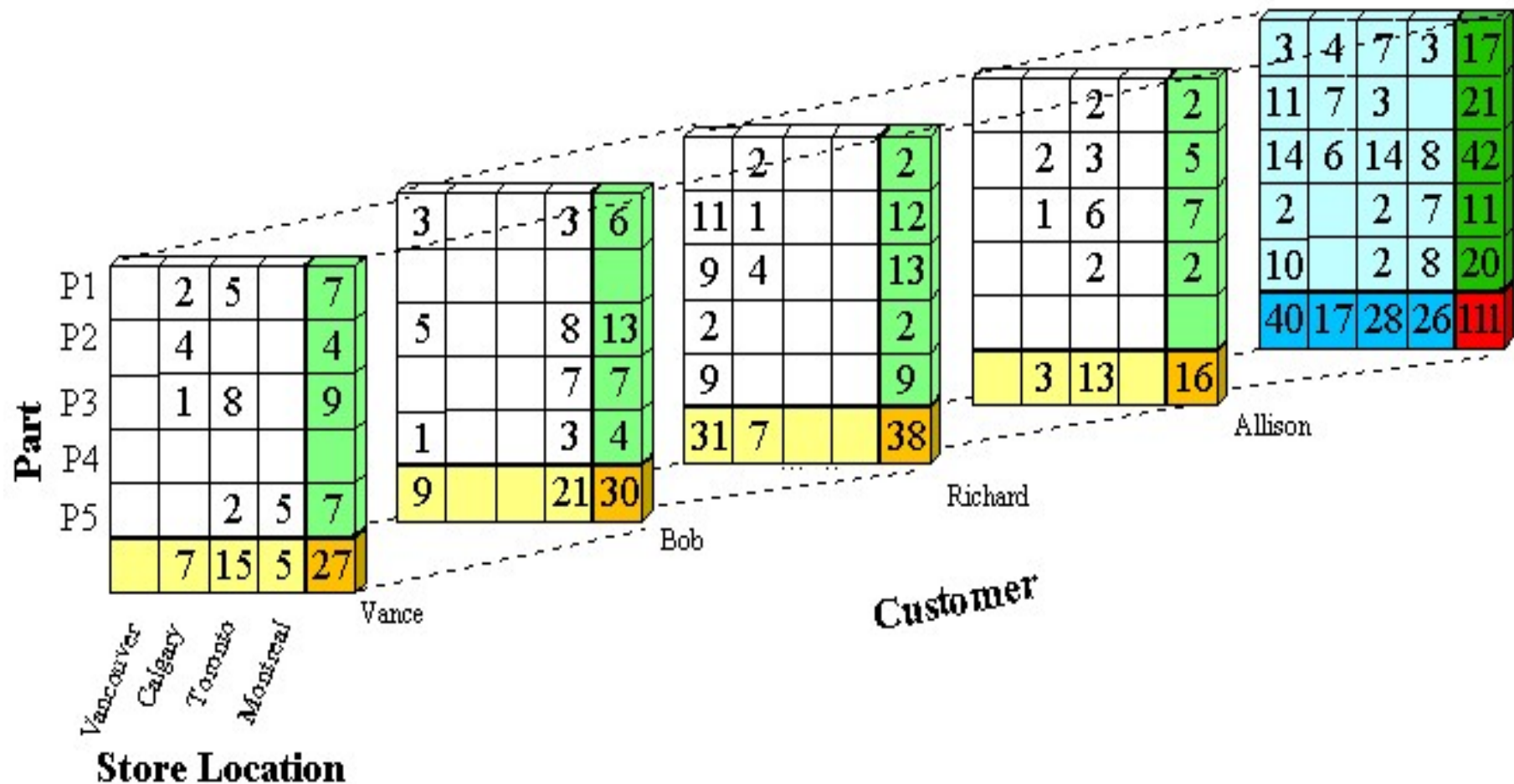
Cube: A Lattice of Cuboids



Stored Data Cubes



Computed Data Cubes



Ordered Set representation of a Data Cube

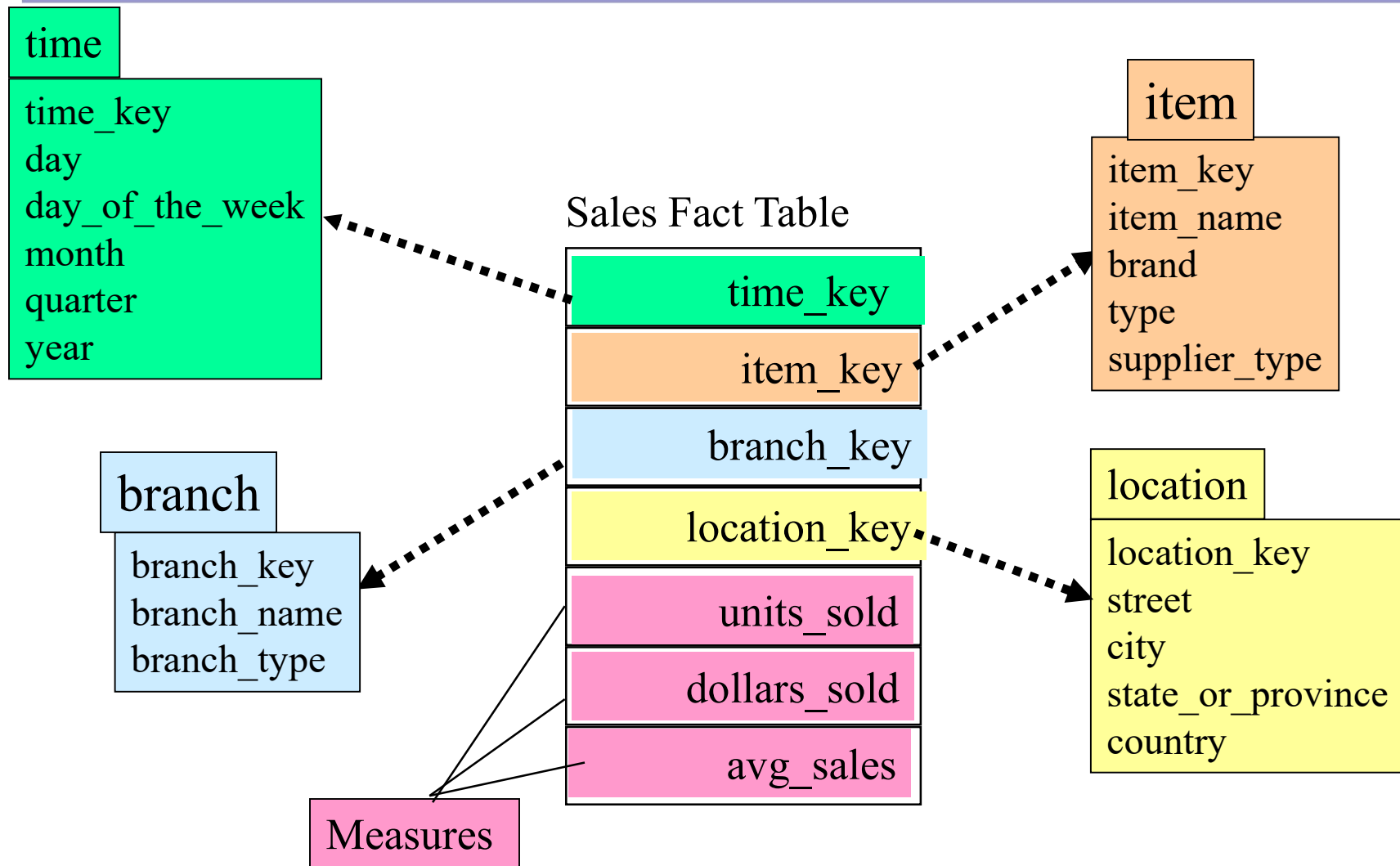
Combination	Count
{P1, Calgary, Vance}	2
{P2, Calgary, Vance}	4
{P3, Calgary, Vance}	1
{P1, Toronto, Vance}	5
{P3, Toronto, Vance}	8
{P5, Toronto, Vance}	2
{P5, Montreal, Vance}	5
{P1, Vancouver, Bob}	3
{P3, Vancouver, Bob}	5
{P5, Vancouver, Bob}	1
{P1, Montreal, Bob}	3
{P3, Montreal, Bob}	8
{P4, Montreal, Bob}	7
{P5, Montreal, Bob}	3
{P2, Vancouver, Richard}	11

{P3, Vancouver, Richard}	9
{P4, Vancouver, Richard}	2
{P5, Vancouver, Richard}	9
{P1, Calgary, Richard}	2
{P2, Calgary, Richard}	1
{P3, Calgary, Richard}	4
{P2, Calgary, Allison}	2
{P3, Calgary, Allison}	1
{P1, Toronto, Allison}	2
{P2, Toronto, Allison}	3
{P3, Toronto, Allison}	6
{P4, Toronto, Allison}	2

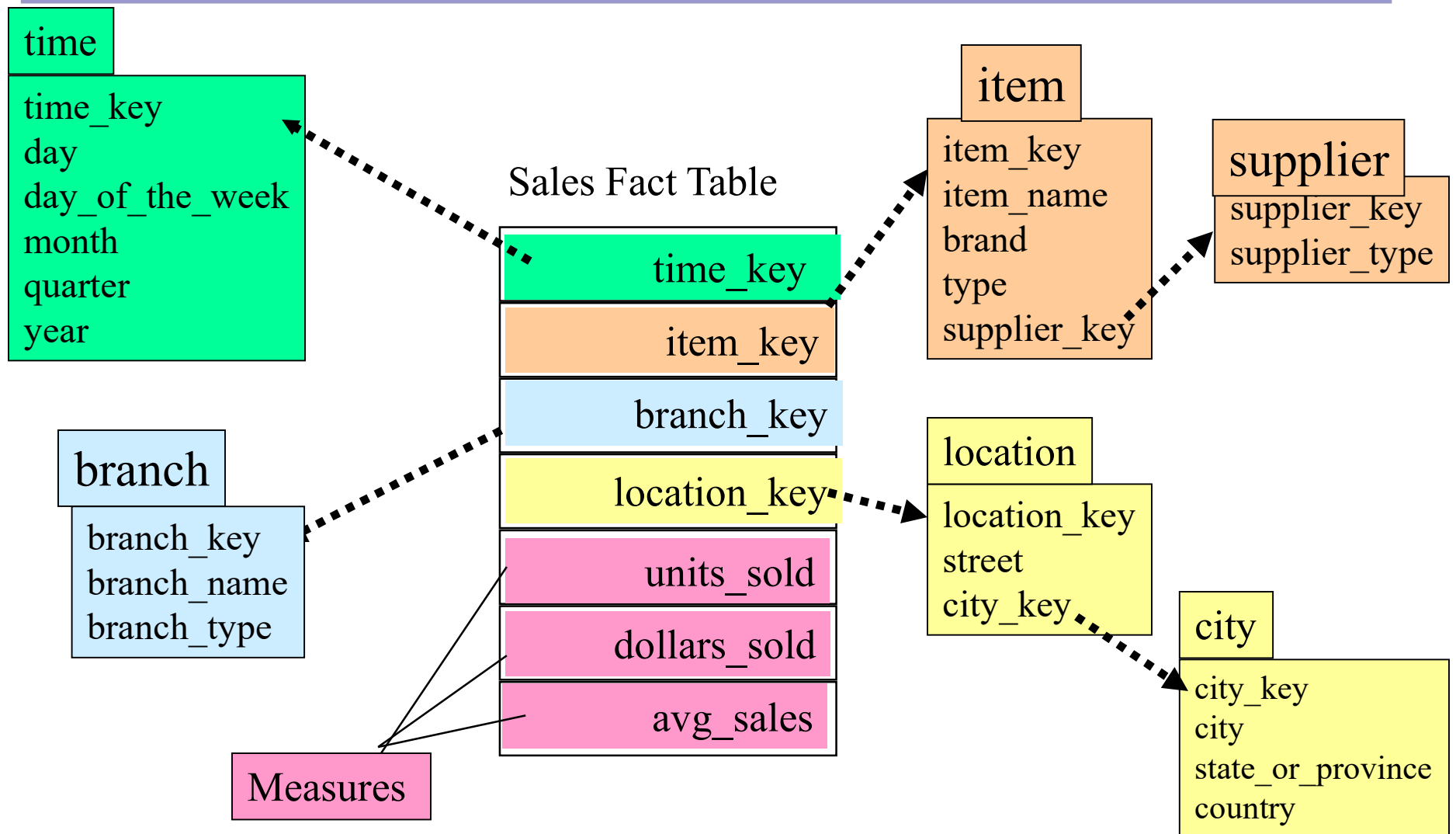
Conceptual Modeling of Data Warehouses

- 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

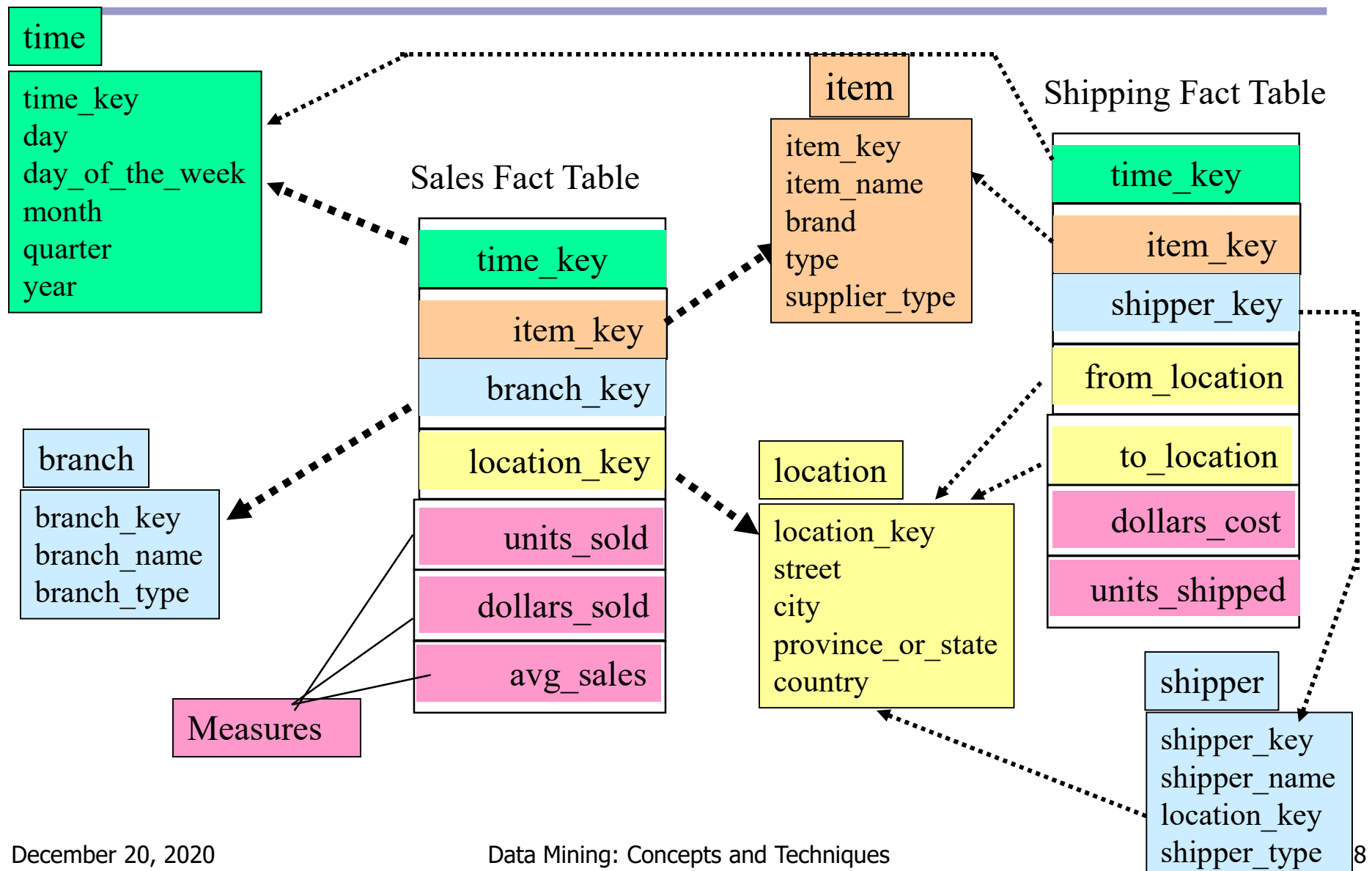
Example of Star Schema



Example of Snowflake Schema



Example of Fact Constellation



Cube Definition Syntax (BNF) in DMQL

- 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>

Defining Star Schema in DMQL

```
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)
```

Defining Snowflake Schema in DMQL

```
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))
```

Defining Fact Constellation in DMQL

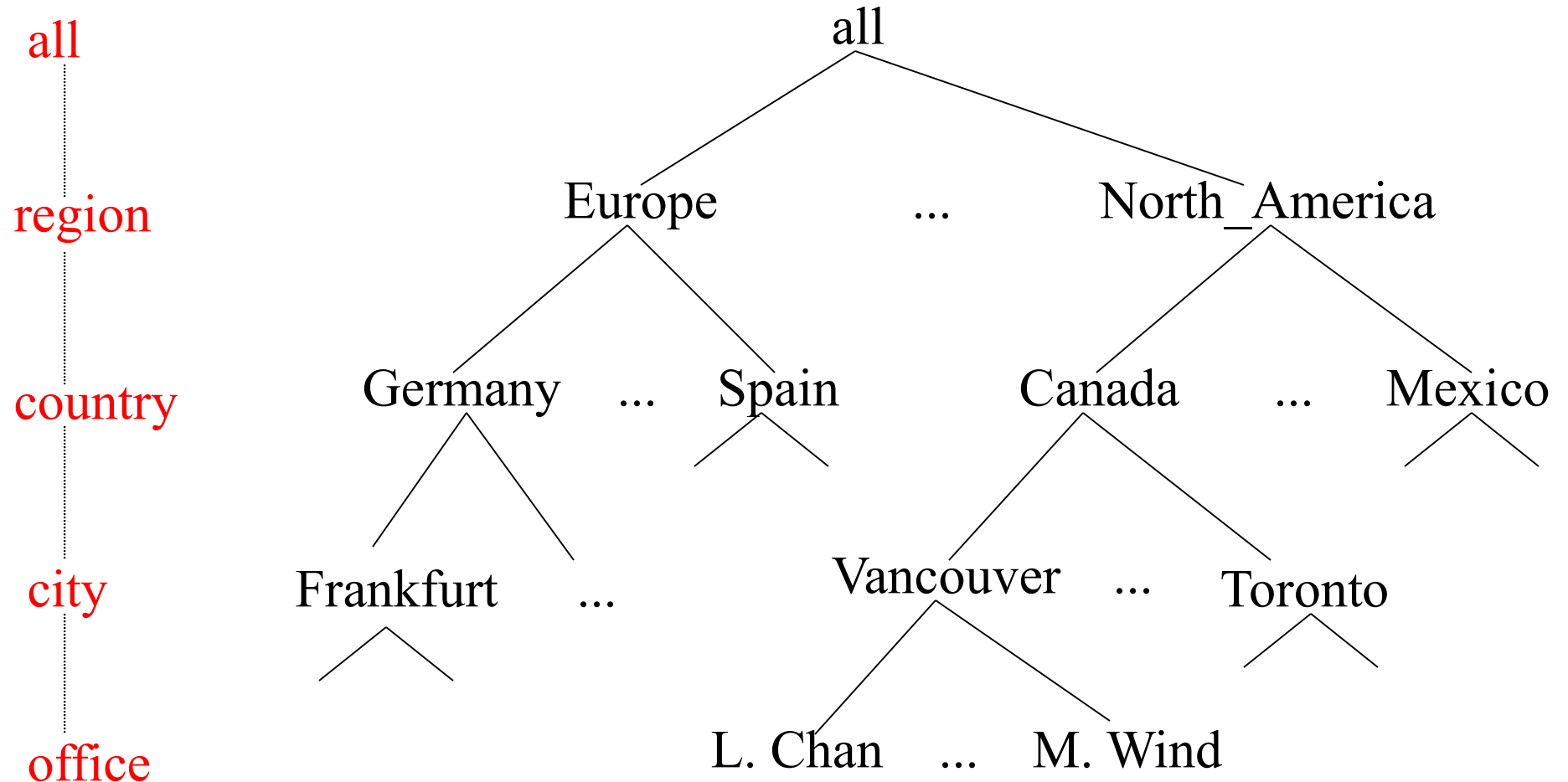
```
define cube sales [time, item, branch, location]:  
    dollars_sold = sum(sales_in_dollars), avg_sales =  
        avg(sales_in_dollars), units_sold = count(*)  
define dimension time as (time_key, day, day_of_week, month, quarter,  
    year)  
define dimension item as (item_key, item_name, brand, 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
```


Measures of Data Cube: Three Categories

- **Distributive**: if the result derived by applying the function to n aggregate values is the same as that derived by applying the function on all the data without partitioning
 - E.g., `count()`, `sum()`, `min()`, `max()`
- **Algebraic**: if it can be computed by an algebraic function with M arguments (where M is a bounded integer), each of which is obtained by applying a distributive aggregate function
 - E.g., `avg()`, `min_N()`, `standard_deviation()`
- **Holistic**: if there is no constant bound on the storage size needed to describe a subaggregate.
 - E.g., `median()`, `mode()`, `rank()`

A Concept Hierarchy: Dimension (location)



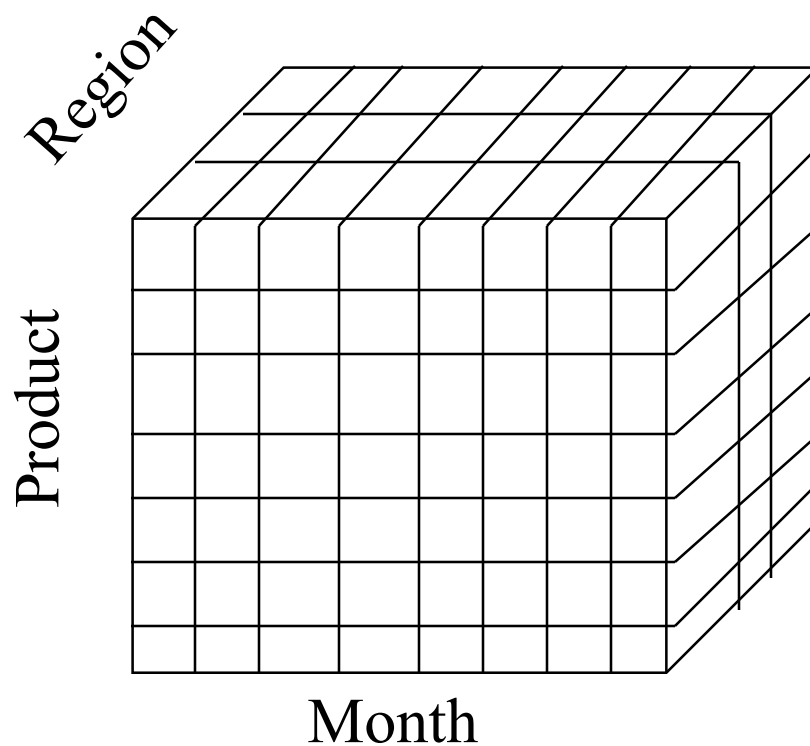
View of Warehouses and Hierarchies

Specification of hierarchies

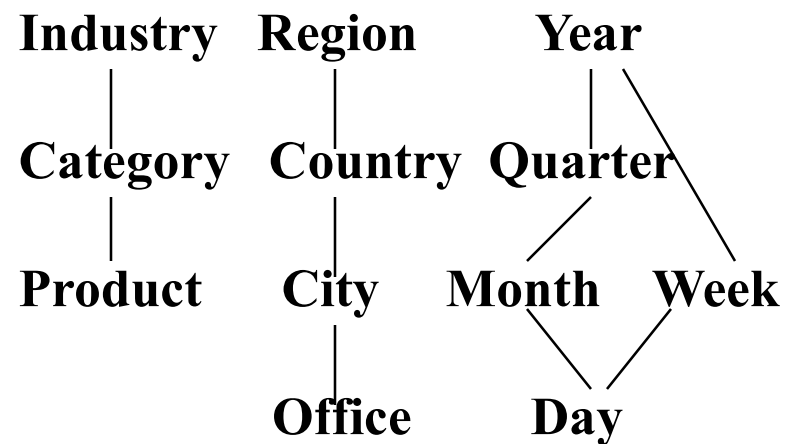
- Schema hierarchy
day < {month < quarter; week} < year
- Set_grouping hierarchy
{1..10} < inexpensive

Multidimensional Data

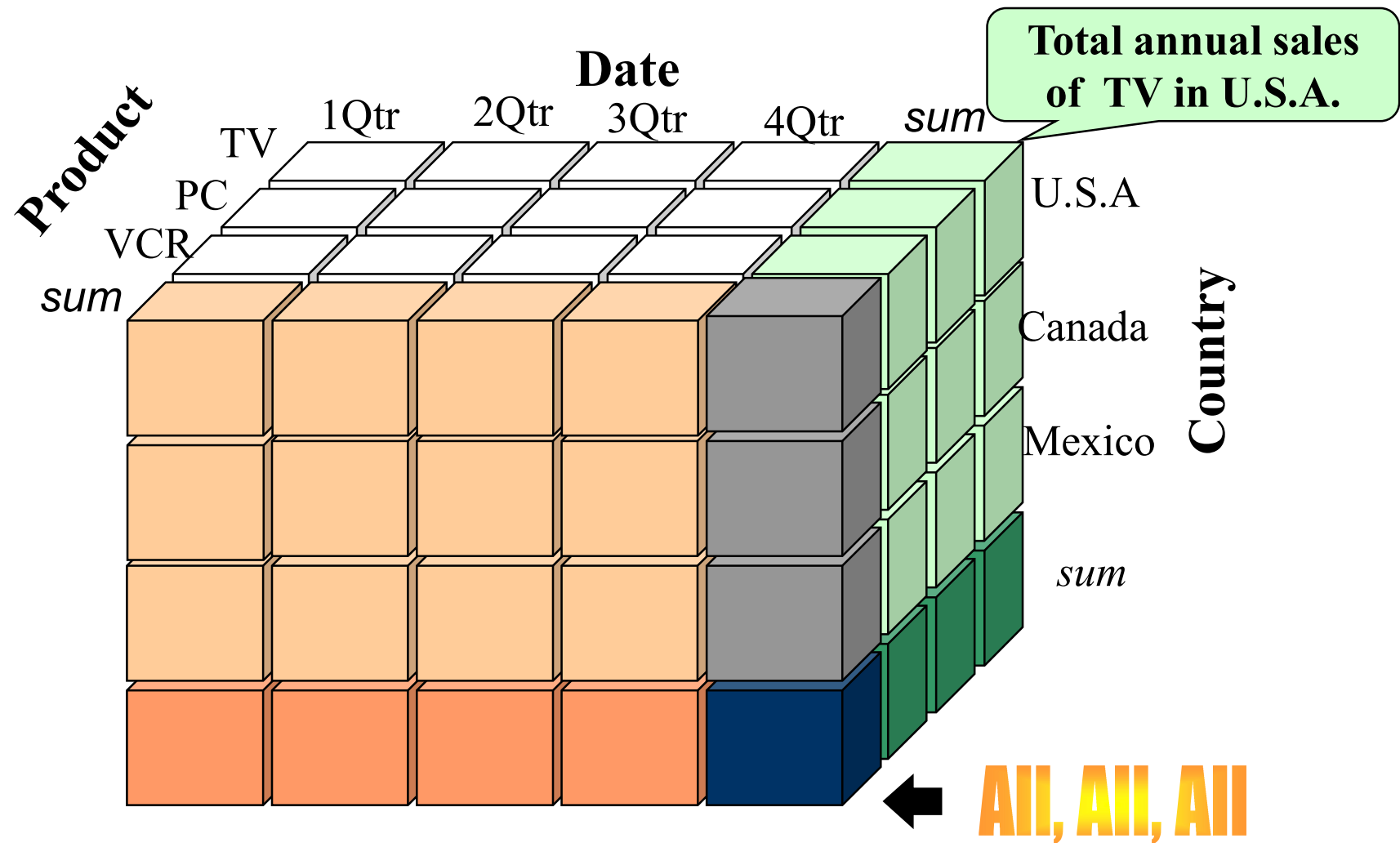
- Sales volume as a function of product, month, and region



Dimensions: Product, Location, Time
Hierarchical summarization paths



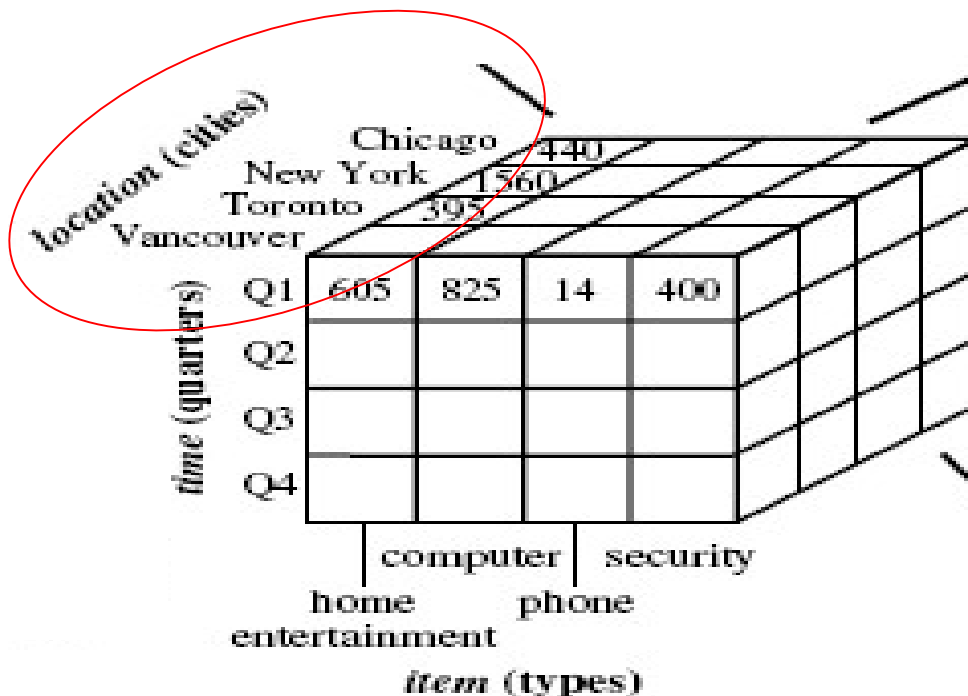
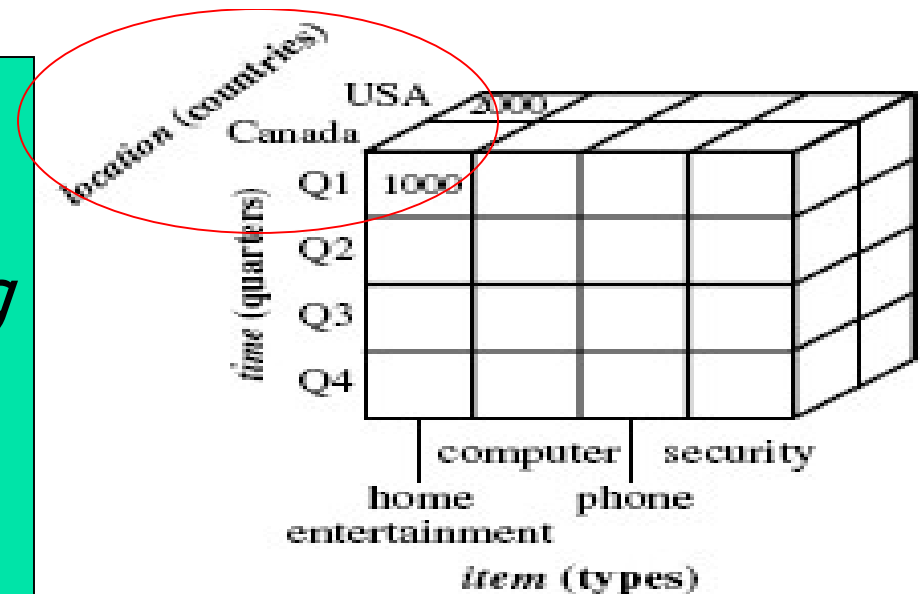
A Sample Data Cube



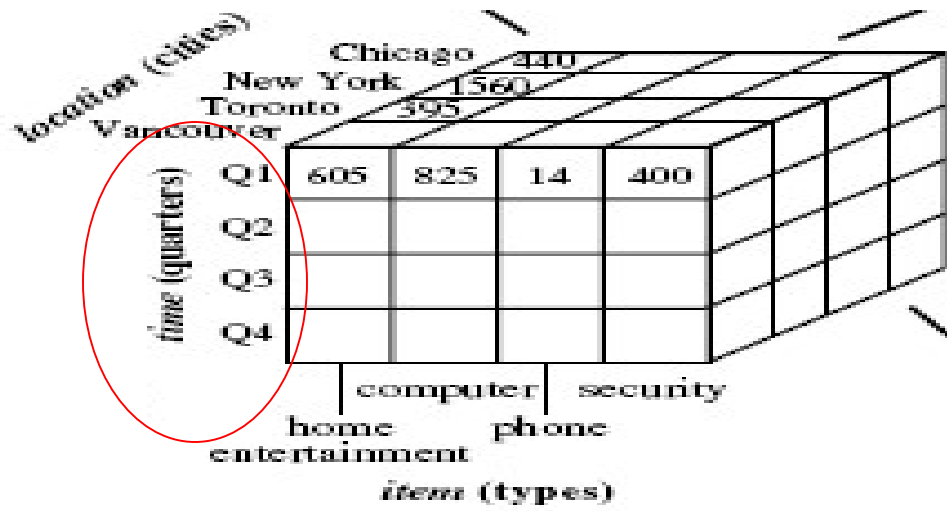
Typical OLAP Operations

- **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)

The roll-up operation performs aggregation on a data cube, either by *climbing up a concept hierarchy* for a dimension or by *dimension reduction*.

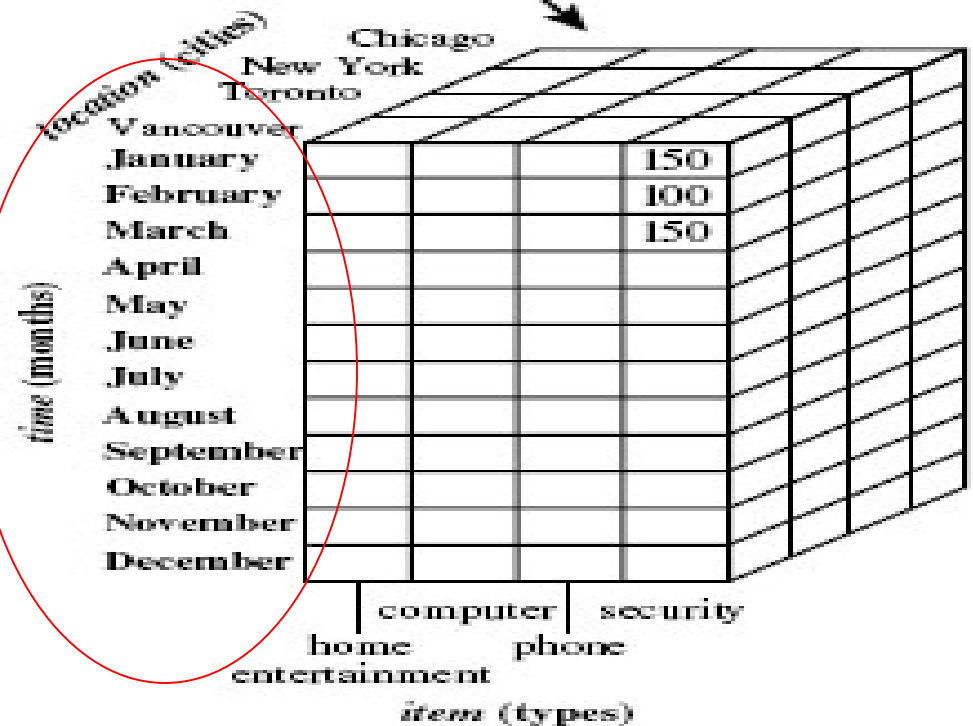


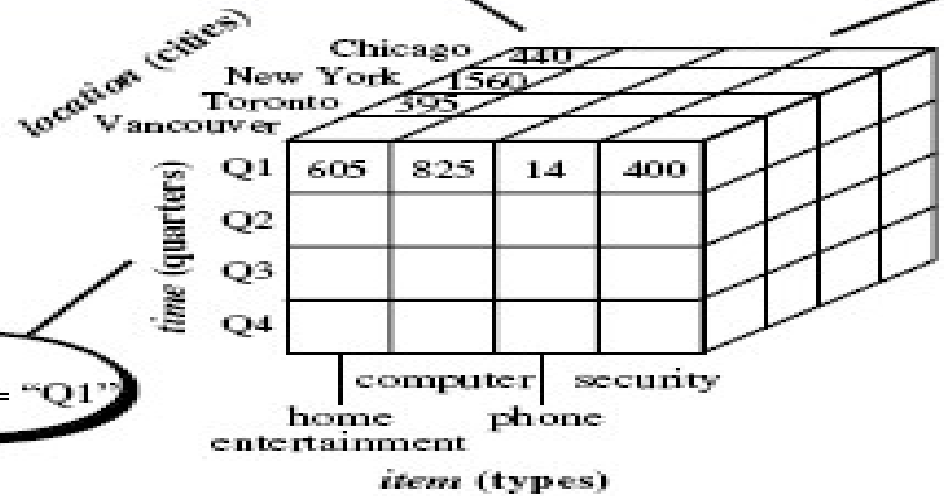
roll-up
on location
(from cities
to countries)



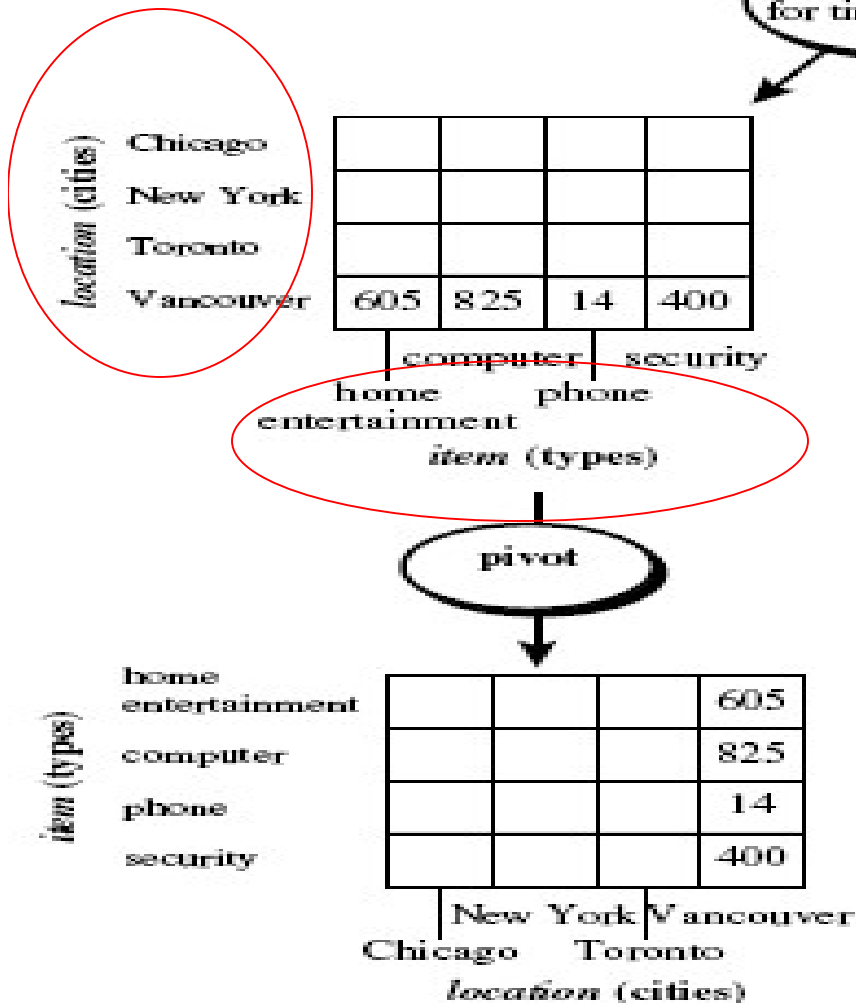
drill-down
on time
(from quarters
to months)

Drill-down is the reverse of roll-up. It navigates from less detailed data to more detailed data. Drill-down can be realized by either stepping down a concept hierarchy for a dimension or introducing additional dimensions.



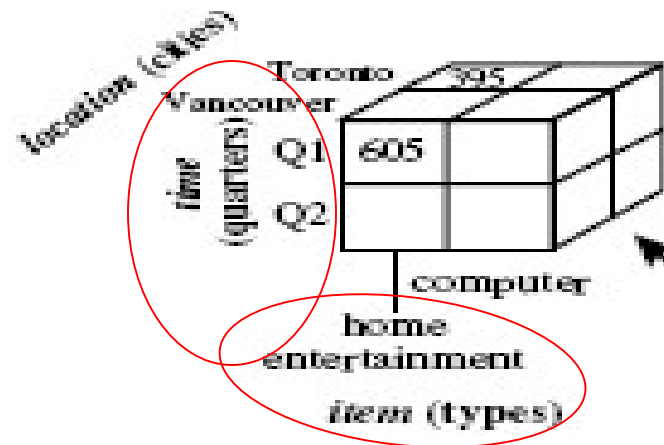


slice
for time = "Q1"

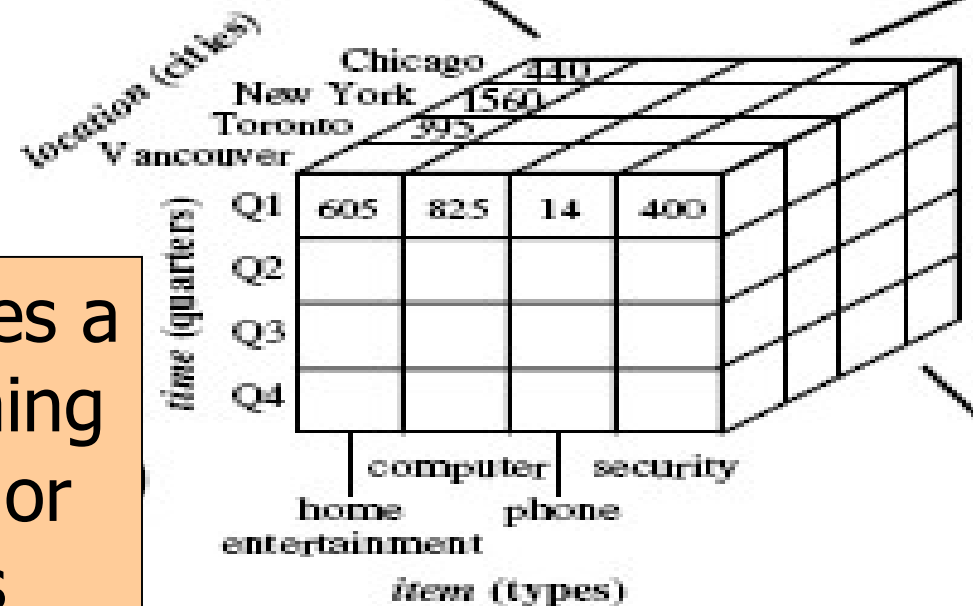


slice operation performs a selection on one dimension of the given cube, resulting in a subcube.

Pivot (also called *rotate*) is a visualization operation that rotates the data axes in view in order to provide an alternative presentation of the data.

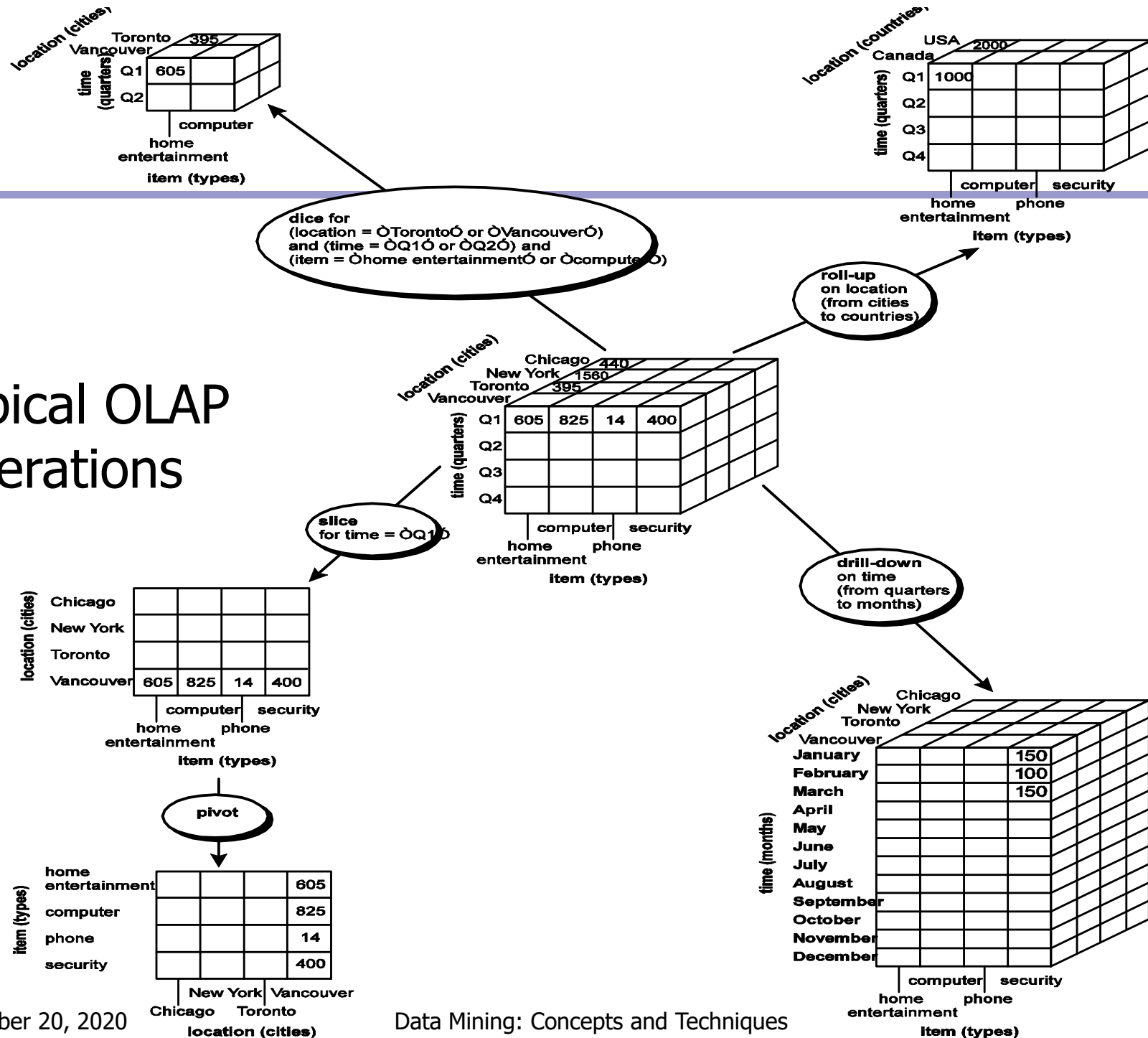


dice for
 (location = "Toronto" or "Vancouver")
 and (time = "Q1" or "Q2") and
 (item = "home entertainment" or "computer")



dice operation defines a subcube by performing a selection on two or more dimensions

Typical OLAP Operations



Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- From data warehousing to data mining

Design of Data Warehouse: A Business Analysis Framework

Four views regarding the design of a data warehouse

- **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

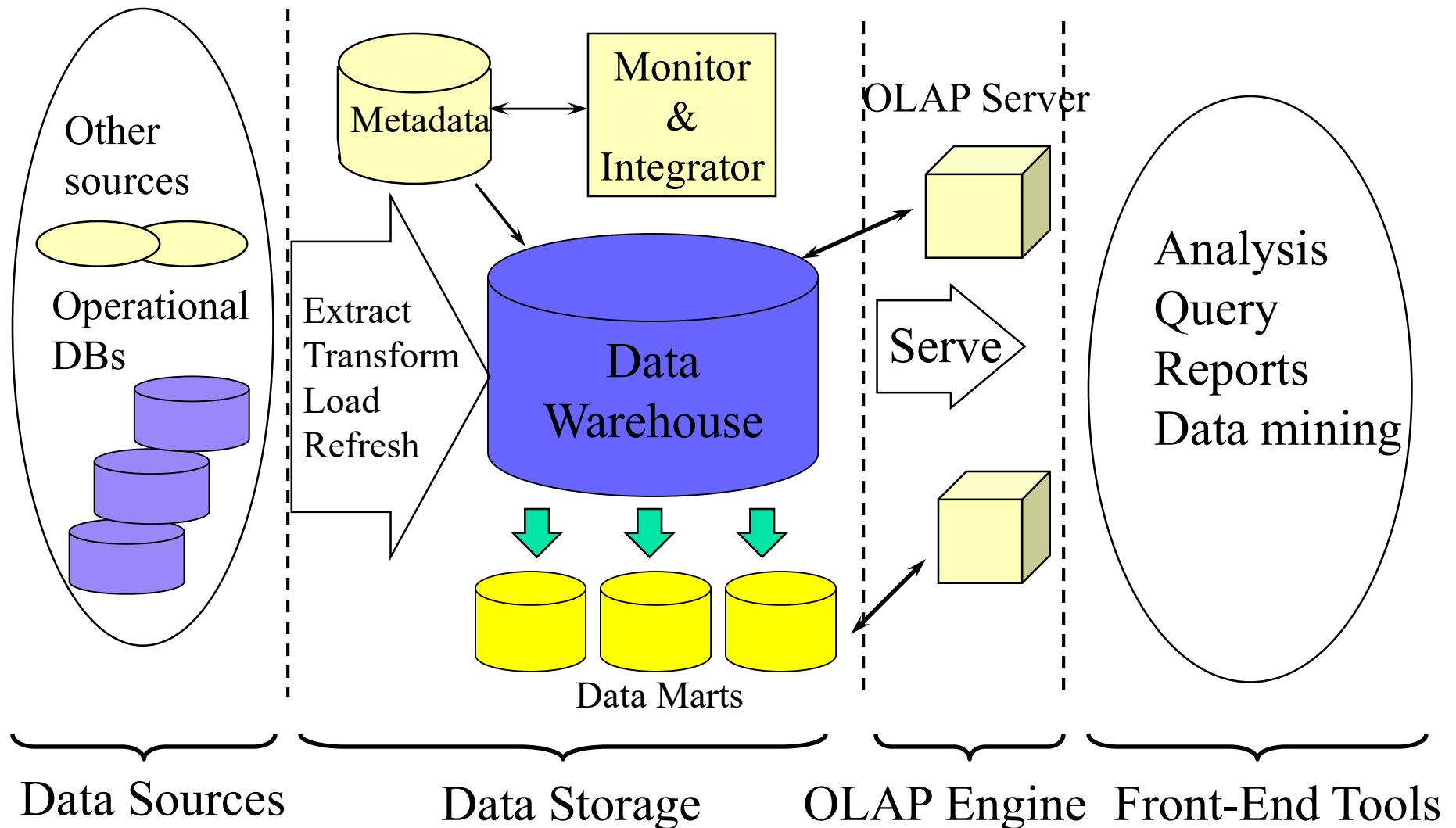
Approaches of DW Design Process

- Top-down, bottom-up approaches or a combination of both
 - Top-down: Starts with overall design and planning (mature)
 - Bottom-up: Starts with experiments and prototypes (rapid)
- From software engineering point of view
 - Waterfall: structured and systematic analysis at each step before proceeding to the next
 - Spiral: rapid generation of increasingly functional systems, short turn around time, quick turn around

Steps for DW design process

- Typical data warehouse design process
 1. Choose a **business process** to model, e.g., orders, invoices, etc.
 2. Choose the **grain (*atomic level of data*)** of the business process
 3. Choose the **dimensions** that will apply to each fact table record
 4. Choose the **measure** that will populate each fact table record

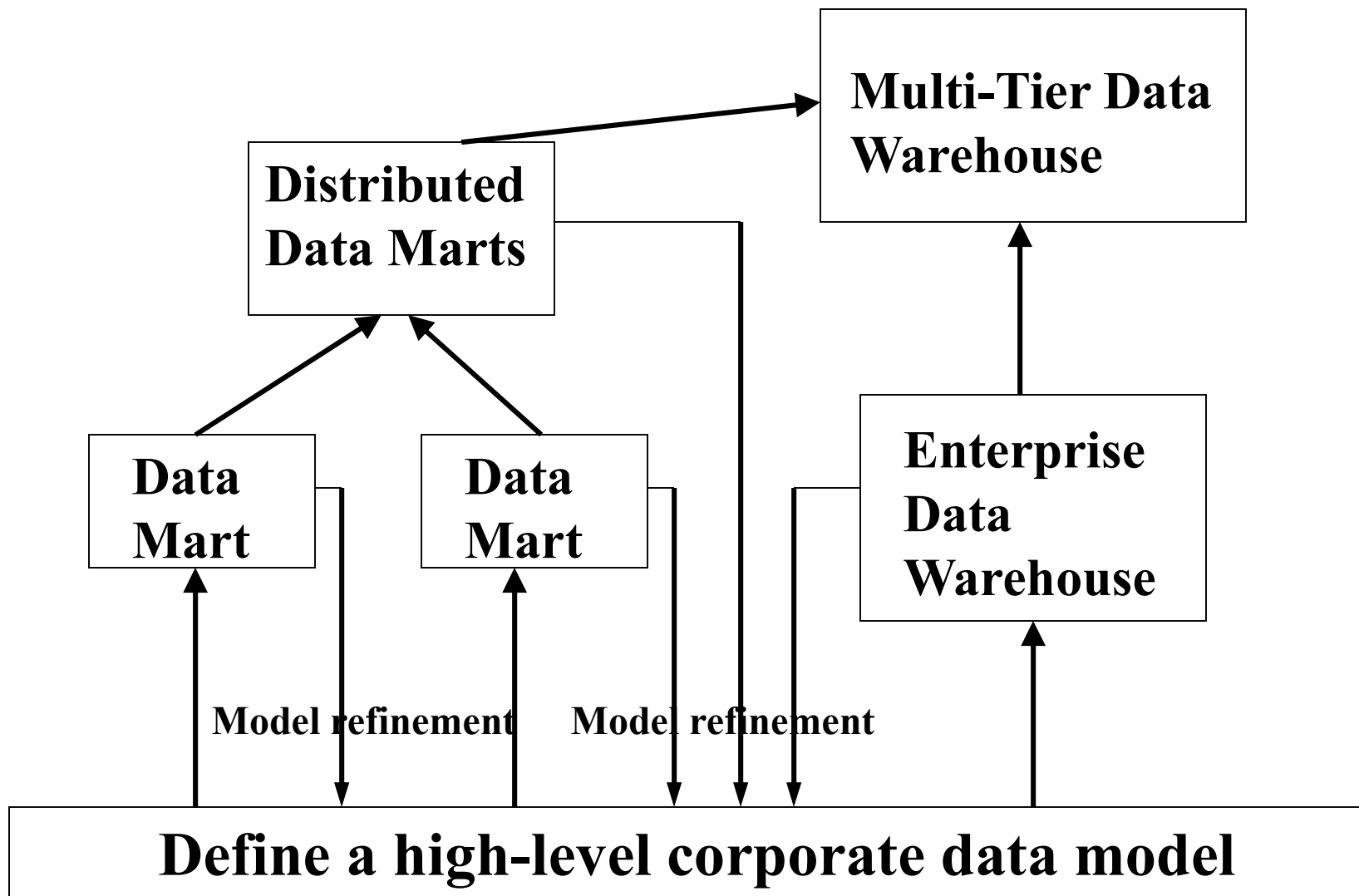
Data Warehouse: A Multi-Tiered Architecture



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



Data Warehouse Back-End Tools and Utilities

- **Data extraction**
 - get data from multiple, heterogeneous, and external sources
- **Data cleaning**
 - detect errors in the data and rectify them when possible
- **Data transformation**
 - convert data from legacy or host format to warehouse format
- **Load**
 - sort, summarize, consolidate, compute views, check integrity, and build indices and partitions
- **Refresh**
 - propagate the updates from the data sources to the warehouse

Metadata Repository

Meta data is the data defining warehouse objects. It stores:

1. Description of the structure of the DW
2. Operational meta-data
 - data lineage
 - data currency
 - monitoring info
3. The mapping from operational environment to the DW
4. The algorithms used for summarization
5. Data related to system performance
6. Business metadata

Metadata Repository

- Description of the structure of the DW
 - schema, view, dimensions, hierarchies, derived data defn, data mart locations and contents
- Operational meta-data
 - data lineage (history of migrated data and transformation path), data currency (active, archived, or purged), monitoring info. (warehouse usage statistics, error reports, audit trails)
- The mapping from operational environment to the DW
- The algorithms used for summarization
- Data related to system performance
 - warehouse schema, view and derived data definitions
- Business metadata
 - business terms and definitions, ownership of data, charging policies

OLAP Server Architectures

- Relational OLAP (ROLAP)

- Use relational or extended-relational DBMS to store and manage warehouse data and OLAP middle ware
- Include optimization of DBMS backend, implementation of aggregation navigation logic, and additional tools and services
- Greater scalability

- Multidimensional OLAP (MOLAP)

- Sparse array-based multidimensional storage engine
- Fast indexing to pre-computed summarized data

- Hybrid OLAP (HOLAP) (e.g., Microsoft SQLServer)

- Flexibility, e.g., low level: relational, high-level: array

- Specialized SQL servers (e.g., Redbricks)

- Specialized support for SQL queries over star/snowflake schemas

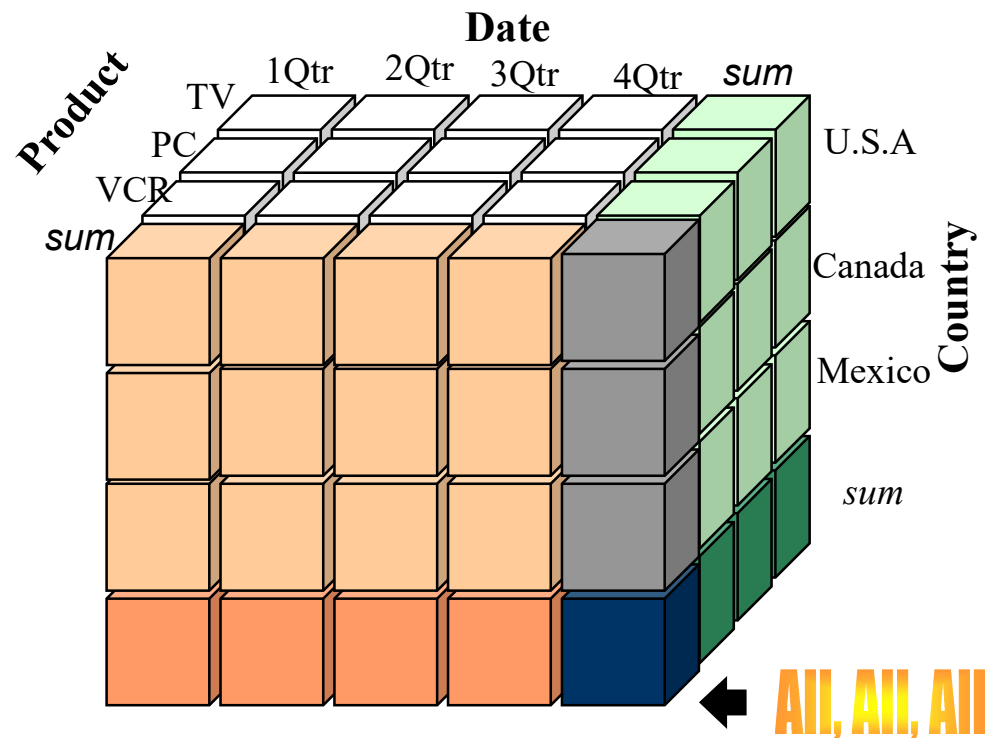
Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- From data warehousing to data mining

Efficient Data Cube Computation

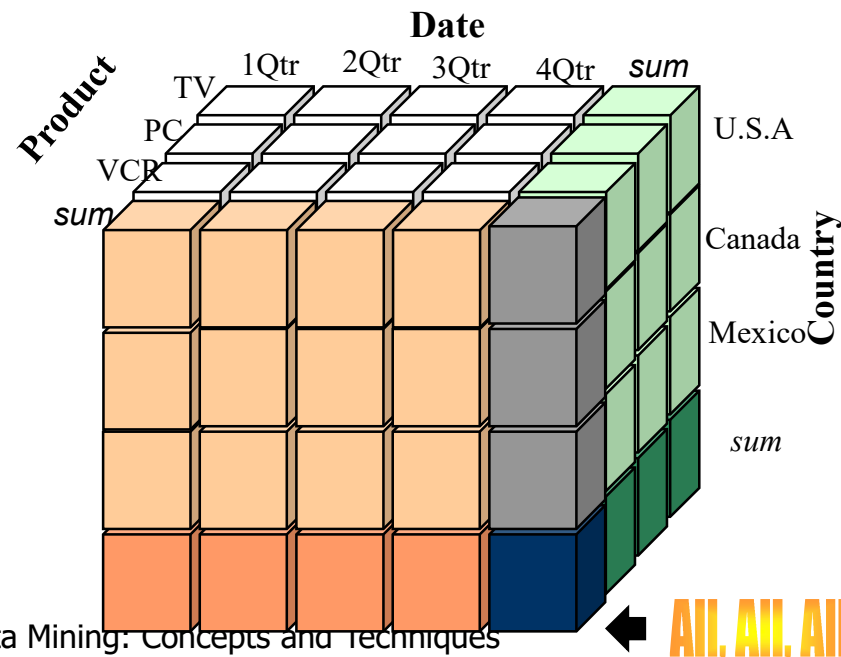
- Data cube can be viewed as a lattice of cuboids
 - The bottom-most cuboid is the base cuboid
 - The top-most cuboid (apex) contains only one cell
 - How many cuboids in an n-dimensional cube with L levels?

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



Efficient Data Cube Computation Contd..)

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



Cube Operation

- Cube definition and computation in DMQL

define cube sales[item, city, year]: sum(sales_in_dollars)

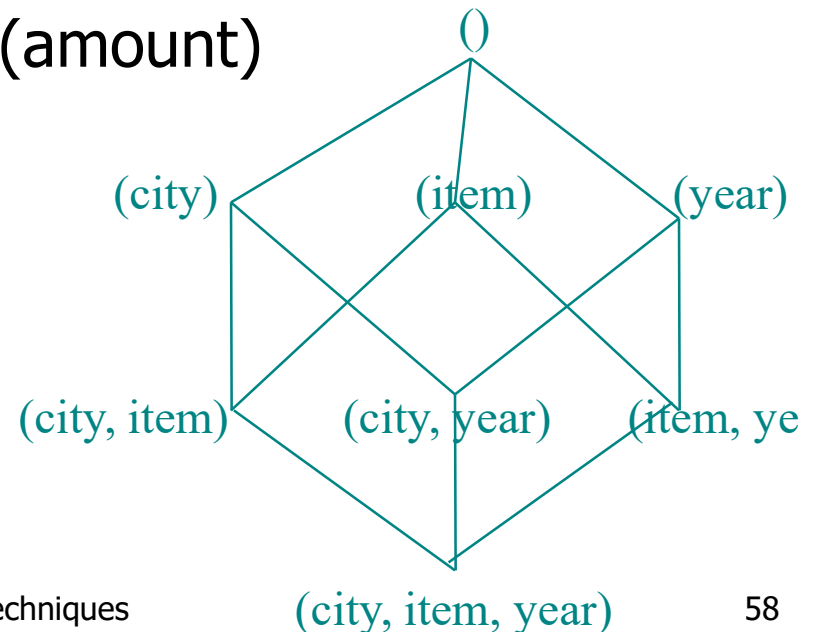
compute cube sales

- Transform it into a SQL-like language (with a new operator **cube by**, introduced by Gray et al.'96)

SELECT item, city, year, SUM (amount)

FROM SALES

CUBE BY item, city, year



Indexing OLAP Data: Bitmap Index

- Index on a particular column
- Each value in the column has a bit vector: **bit-op is fast**
- The length of the bit vector: # of records in the base table
- The i -th bit is set if the i -th row of the base table has the value for the indexed column
- not suitable for high cardinality domains

Base table

Cust	Region	Type
C1	Asia	Retail
C2	Europe	Dealer
C3	Asia	Dealer
C4	America	Retail
C5	Europe	Dealer

Index on Region

RecID	Asia	Europe	America
1	1	0	0
2	0	1	0
3	1	0	0
4	0	0	1
5	0	1	0

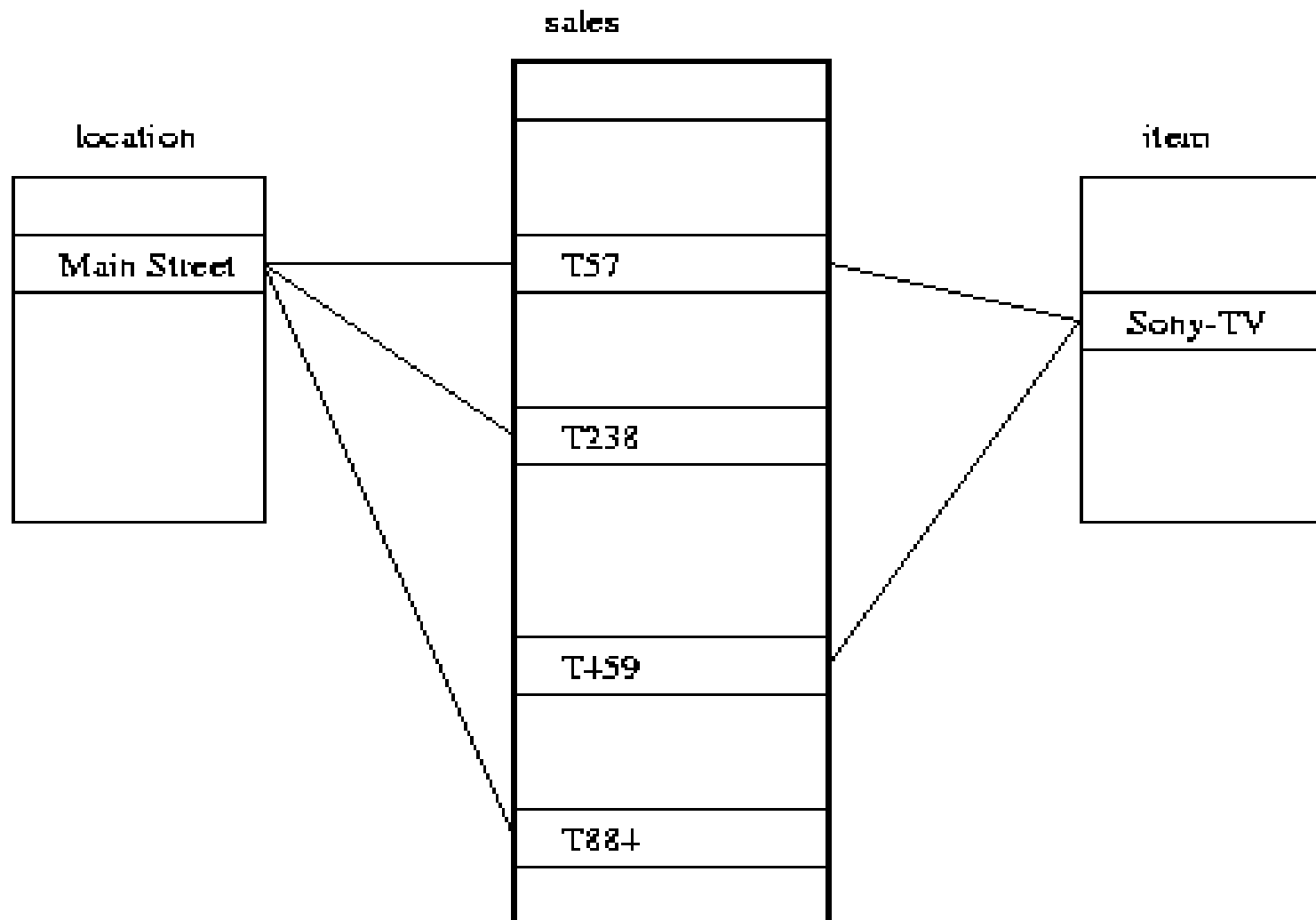
Index on Type

RecID	Retail	Dealer
1	1	0
2	0	1
3	0	1
4	1	0
5	0	1

Indexing OLAP Data: Join Indices

- Join index: $JI(R\text{-id}, S\text{-id})$ where $R(R\text{-id}, \dots) \triangleright \triangleleft S(S\text{-id}, \dots)$
- Traditional indices map the values to a list of record ids
 - It materializes relational join in JI file and speeds up relational join
- In data warehouses, join index relates the values of the dimensions of a star schema to rows in the fact table.
 - E.g. fact table: *Sales* and two dimensions *city* and *product*
 - A join index on *city* maintains for each distinct city a list of R-IDs of the tuples recording the Sales in the city
 - Join indices can span multiple dimensions

Indexing OLAP Data: Join Indices



Efficient Processing OLAP Queries

- **Determine which operations should be performed on the available cuboids**
 - Transform drill, roll, etc. into corresponding SQL and/or OLAP operations, e.g., dice = selection + projection
- **Determine which materialized cuboid(s) should be selected for OLAP op. -- Which should be selected to process the query?**
 - 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 = 2004
- **Explore indexing structures and compressed vs. dense array structs in MOLAP**

Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- From data warehousing to data mining

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

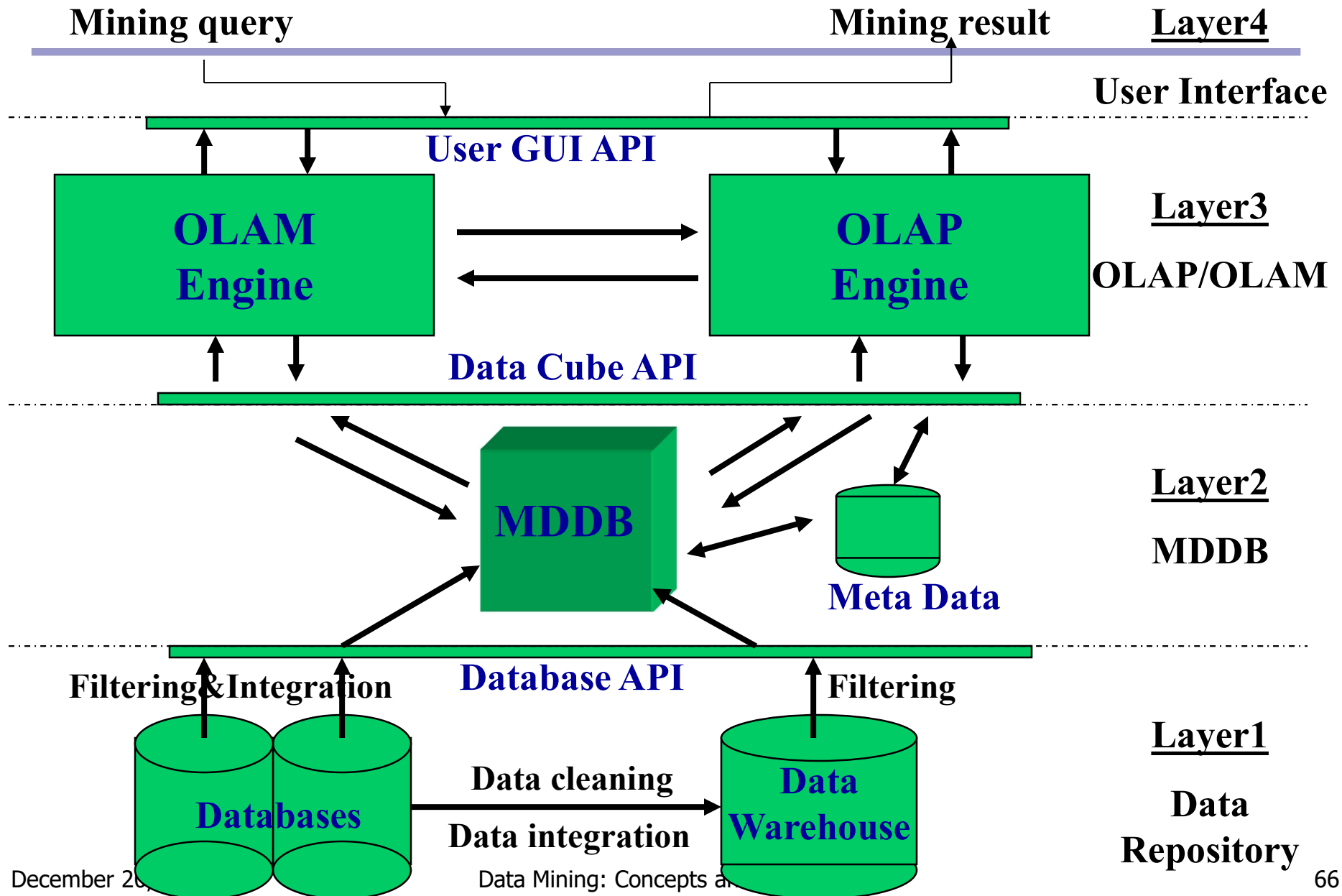
From On-Line Analytical Processing (OLAP) to On Line Analytical Mining (OLAM)

- (OLAM=OLAP+DM & mining kn in multi-D DBs.)

Why online analytical mining?

- High quality of data in DW
 - DW contains integrated, consistent, cleaned data
- 1. Available information processing structure surrounding DW
 - ODBC, OLEDB, Web accessing, service facilities, reporting and OLAP tools
- 2. OLAP-based exploratory data analysis
 - Mining with drilling, dicing, pivoting, etc.
- 3. On-line selection of data mining functions
 - Integration and swapping of multiple mining functions, algorithms, and tasks

An OLAM System Architecture



Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- From data warehousing to data mining
- Summary

Summary: Data Warehouse and OLAP Technology

- Why data warehousing?
- A **multi-dimensional model** of a data warehouse
 - Star schema, snowflake schema, fact constellations
 - A data cube consists of dimensions & measures
- **OLAP** operations: drilling, rolling, slicing, dicing and pivoting
- Data warehouse architecture
- OLAP servers: ROLAP, MOLAP, HOLAP
- Efficient computation of data cubes
 - Partial vs. full vs. no materialization
 - Indexing OALP data: Bitmap index and join index
 - OLAP query processing
- From OLAP to OLAM (on-line analytical mining)

Exercises

1. Briefly compare the following concepts with example to explain:

- (a) Snowflake schema, fact constellation, starlet query model
- (b) Data cleaning, data transformation, refresh
- (c) Discovery-driven cube, multifeature cube, virtual warehouse

2. Suppose that a data warehouse consists of the three dimensions *time*, *doctor*, and *patient*, and the two measures *count* and *charge*, where *charge* is the fee that a doctor charges a patient for a visit.

- a) Enumerate three classes of schemas that are popularly used for modeling data warehouses.
- b) Draw a schema diagram for the above data warehouse using one of the schema classes
- c) Starting with the base cuboid [day, doctor, patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2010?
- d) To obtain the same list, write an SQL query assuming the data are stored in a relational database with the schema fee (day, month, year, doctor, hospital, patient, count, charge).

Exercises

3. Suppose that a data warehouse for Big University consists of the four dimensions student, course, semester, and instructor, and two measures count and avg grade. At the lowest conceptual level (e.g., for a given student, course, semester, and instructor combination), the avg grade measure stores the actual course grade of the student. At higher conceptual levels, avg grade stores the average grade for the given combination.
- (a) Draw a snowflake schema diagram for the data warehouse.
 - (b) Starting with the base cuboid [student, course, semester, instructor], what specific OLAP operations (e.g., roll-up from semester to year) should you perform in order to list the average grade of CS courses for each Big University student.
 - (c) If each dimension has five levels (including all), such as “student < major < status < university < all”, how many cuboids will this cube contain (including the base and apex cuboids)?

.

Exercises

4. Suppose that a data warehouse consists of the four dimensions date, spectator, location, and game, and the two measures count and charge, where charge is the fare that a spectator pays when watching a game on a given date. Spectators may be students, adults, or seniors, with each category having its own charge rate.
- (a) Draw a star schema diagram for the data warehouse.
 - (b) Starting with the base cuboid [date, spectator, location, game], what specific OLAP operations should you perform in order to list the total charge paid by student spectators at GM Place in 2010?
 - (c) Bitmap indexing is useful in data warehousing. Taking this cube as an example, briefly discuss advantages and problems of using a bitmap index structure.