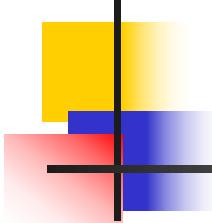


Data Mining: Concepts and Techniques

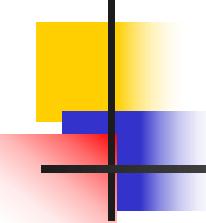
— Slides for Textbook —
— Chapter 2 —

©Jiawei Han and Micheline Kamber
Intelligent Database Systems Research Lab
School of Computing Science
Simon Fraser University, Canada
<http://www.cs.sfu.ca>



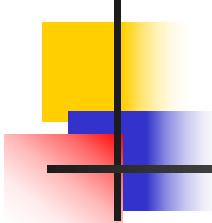
Chapter 2: Data Warehousing and OLAP Technology for Data Mining

- n What is a data warehouse?
- n A multi-dimensional data model
- n Data warehouse architecture
- n Data warehouse implementation
- n Further development of data cube technology
- n From data warehousing to data mining



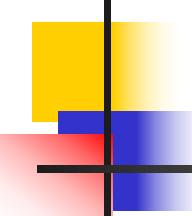
What is Data Warehouse?

- n Defined in many different ways, but not rigorously.
 - n A decision support database that is maintained **separately** from the organization's operational database
 - n Support **information processing** by providing a solid platform of consolidated, historical data for analysis.
- n “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
- n Data warehousing:
 - n The process of constructing and using data warehouses



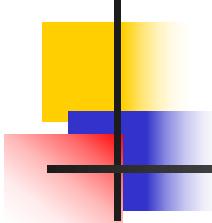
Data Warehouse—Subject-Oriented

- „ 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**.



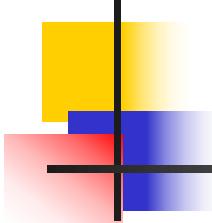
Data Warehouse—Integrated

- n Constructed by integrating multiple, heterogeneous data sources
 - n relational databases, flat files, on-line transaction records
- n Data cleaning and data integration techniques are applied.
 - n Ensure consistency in naming conventions, encoding structures, attribute measures, etc. among different data sources
 - n E.g., Hotel price: currency, tax, breakfast covered, etc.
 - n When data is moved to the warehouse, it is converted.



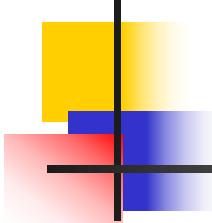
Data Warehouse—Time Variant

- 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”.



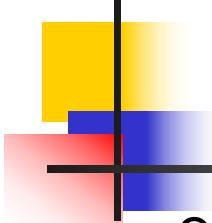
Data Warehouse—Non-Volatile

- 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

- n Traditional heterogeneous DB integration:
 - n Build **wrappers/mediators** on top of heterogeneous databases
 - n **Query driven** approach
 - n 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
 - n Complex information filtering, compete for resources
- n Data warehouse: **update-driven**, high performance
 - n Information from heterogeneous sources is integrated in advance and stored in warehouses for direct query and analysis

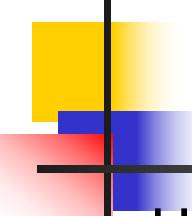


Data Warehouse vs. Operational DBMS

- 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

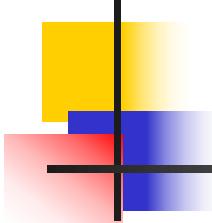
OLTP vs. OLAP

	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



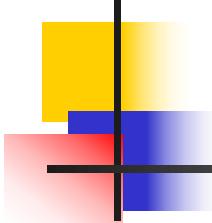
Why Separate Data Warehouse?

- n High performance for both systems
 - n DBMS— tuned for OLTP: access methods, indexing, concurrency control, recovery
 - n Warehouse—tuned for OLAP: complex OLAP queries, multidimensional view, consolidation.
- n Different functions and different data:
 - n missing data: Decision support requires historical data which operational DBs do not typically maintain
 - n data consolidation: DS requires consolidation (aggregation, summarization) of data from heterogeneous sources
 - n data quality: different sources typically use inconsistent data representations, codes and formats which have to be reconciled



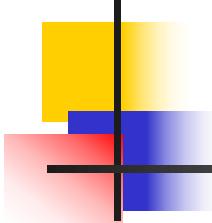
Chapter 2: Data Warehousing and OLAP Technology for Data Mining

- n What is a data warehouse?
- n A multi-dimensional data model
- n Data warehouse architecture
- n Data warehouse implementation
- n Further development of data cube technology
- n From data warehousing to data mining



A multi-dimensional data model

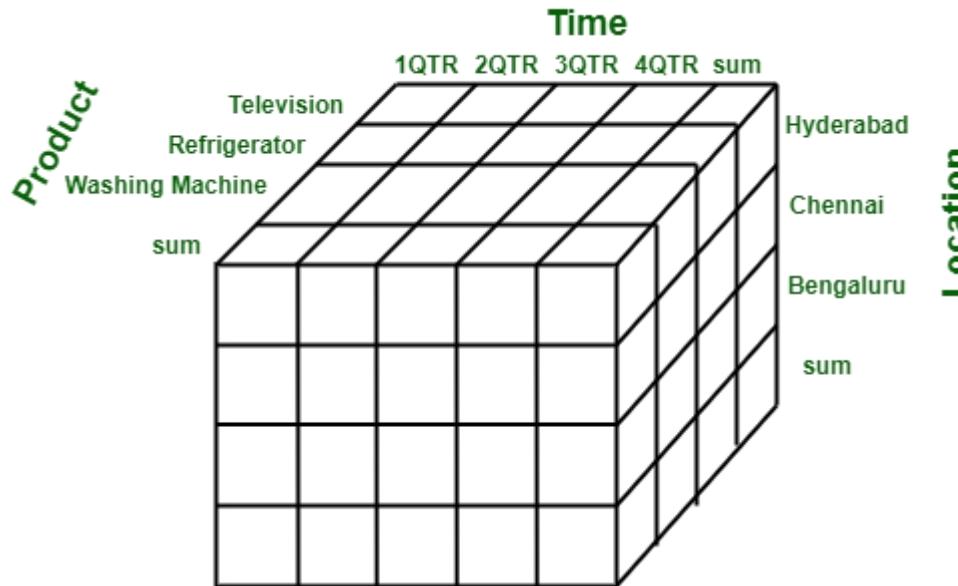
- A multidimensional model views data in the form of a data-cube. A data cube enables data to be modeled and viewed in multiple dimensions. It is defined by dimensions and facts.
- The dimensions are the perspectives or entities concerning which an organization keeps records. For example, a shop may create a sales data warehouse to keep records of the store's sales for the dimension time, item, and location. These dimensions allow the user to keep track of things, for example, monthly sales of items and the locations at which the items were sold. Each dimension has a table related to it, called a dimensional table, which describes the dimension further. For example, a dimensional table for an item may contain the attributes item_name, brand, and type.
- A multidimensional data model is organized around a central theme, for example, sales. This theme is represented by a fact table. Facts are numerical measures. The fact table contains the names of the facts or measures of the related dimensional tables.



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**.

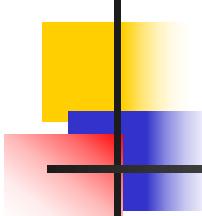
A multi-dimensional data model



Working on a Multidimensional Data Model

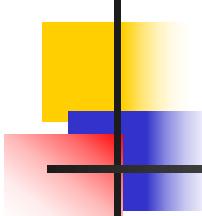
On the basis of the pre-decided steps, the Multidimensional Data Model works.

The following stages should be followed by every project for building a Multi Dimensional Data Model :



Working on a Multidimensional Data Model

- n Stage 1 : Assembling data from the client : In first stage, a Multi Dimensional Data Model collects correct data from the client. Mostly, software professionals provide simplicity to the client about the range of data which can be gained with the selected technology and collect the complete data in detail.
- n Stage 2 : Grouping different segments of the system : In the second stage, the Multi Dimensional Data Model recognizes and classifies all the data to the respective section they belong to and also builds it problem-free to apply step by step.
- n Stage 3 : Noticing the different proportions : In the third stage, it is the basis on which the design of the system is based. In this stage, the main factors are recognized according to the user's point of view. These factors are also known as "Dimensions".
- n Stage 4 : Preparing the actual-time factors and their respective qualities : In the fourth stage, the factors which are recognized in the previous step are used further for identifying the related qualities. These qualities are also known as "attributes" in the database.



Working on a Multidimensional Data Model

- n Stage 5 : Finding the actuality of factors which are listed previously and their qualities : In the fifth stage, A Multi Dimensional Data Model separates and differentiates the actuality from the factors which are collected by it. These actually play a significant role in the arrangement of a Multi Dimensional Data Model.
- n Stage 6 : Building the Schema to place the data, with respect to the information collected from the steps above : In the sixth stage, on the basis of the data which was collected previously, a Schema is built.

Example

- n Let us take the example of the data of a factory which sells products per quarter in Bangalore. The data is represented in the table given below :

Location = "Bangalore"				
Time (quarter)	Type of item			
	Jam	Bread	Sugar	Milk
Q1	350	389	35	50
Q2	260	528	50	90
Q3	483	256	20	60
Q4	436	396	15	40

n 2D factory data

- n In the above given presentation, the factory's sales for Bangalore are, for the time dimension, which is organized into quarters and the dimension of items, which is sorted according to the kind of item which is sold. The facts here are represented in rupees (in thousands).

Example

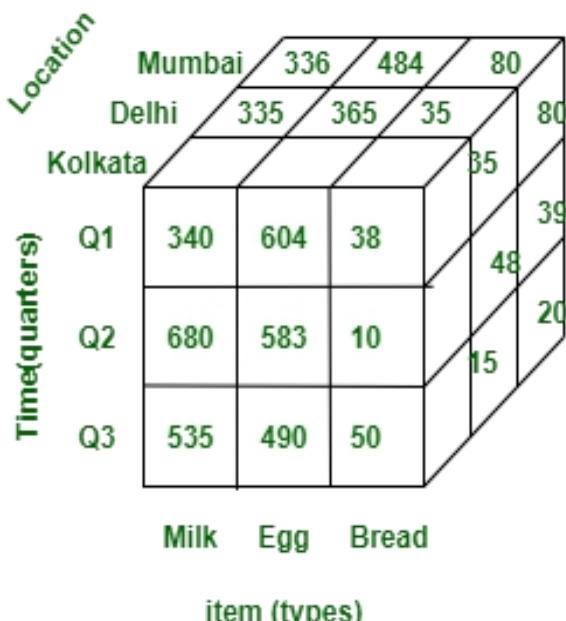
- Now, if we desire to view the data of the sales in a three-dimensional table, then it is represented in the diagram given below. Here the data of the sales is represented as a two dimensional table. Let us consider the data according to item, time and location (like Kolkata, Delhi, Mumbai). Here is the table :

Time	Location="Kolkata"			Location="Delhi"			Location="Mumbai"		
	item			item			item		
	Milk	Egg	Bread	Milk	Egg	Bread	Milk	Egg	Bread
Q1	340	604	38	335	365	35	336	484	80
Q2	680	583	10	684	490	48	595	594	39
Q3	535	490	50	389	385	15	366	385	20

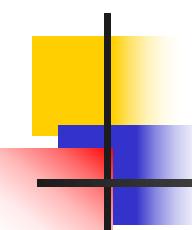
3D data representation as 2D

Example

This data can be represented in the form of three dimensions conceptually, which is shown in the image below :

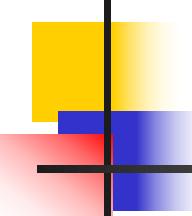


Time	Location="Kolkata"			Location="Delhi"			Location="Mumbai"		
	item			item			item		
	Milk	Egg	Bread	Milk	Egg	Bread	Milk	Egg	Bread
Q1	340	604	38	335	365	35	336	484	80
Q2	680	583	10	684	490	48	595	594	39
Q3	535	490	50	389	385	15	366	385	20



Multi Dimensional Data Model

- **Advantages of Multi Dimensional Data Model**
- The following are the advantages of a multi-dimensional data model :
 - A multi-dimensional data model is easy to handle.
 - It is easy to maintain.
 - Its performance is better than that of normal databases (e.g. relational databases).
 - The representation of data is better than traditional databases. That is because the multi-dimensional databases are multi-viewed and carry different types of factors.
 - It is workable on complex systems and applications, contrary to the simple one-dimensional database systems.
 - The compatibility in this type of database is an upliftment for projects having lower bandwidth for maintenance staff.



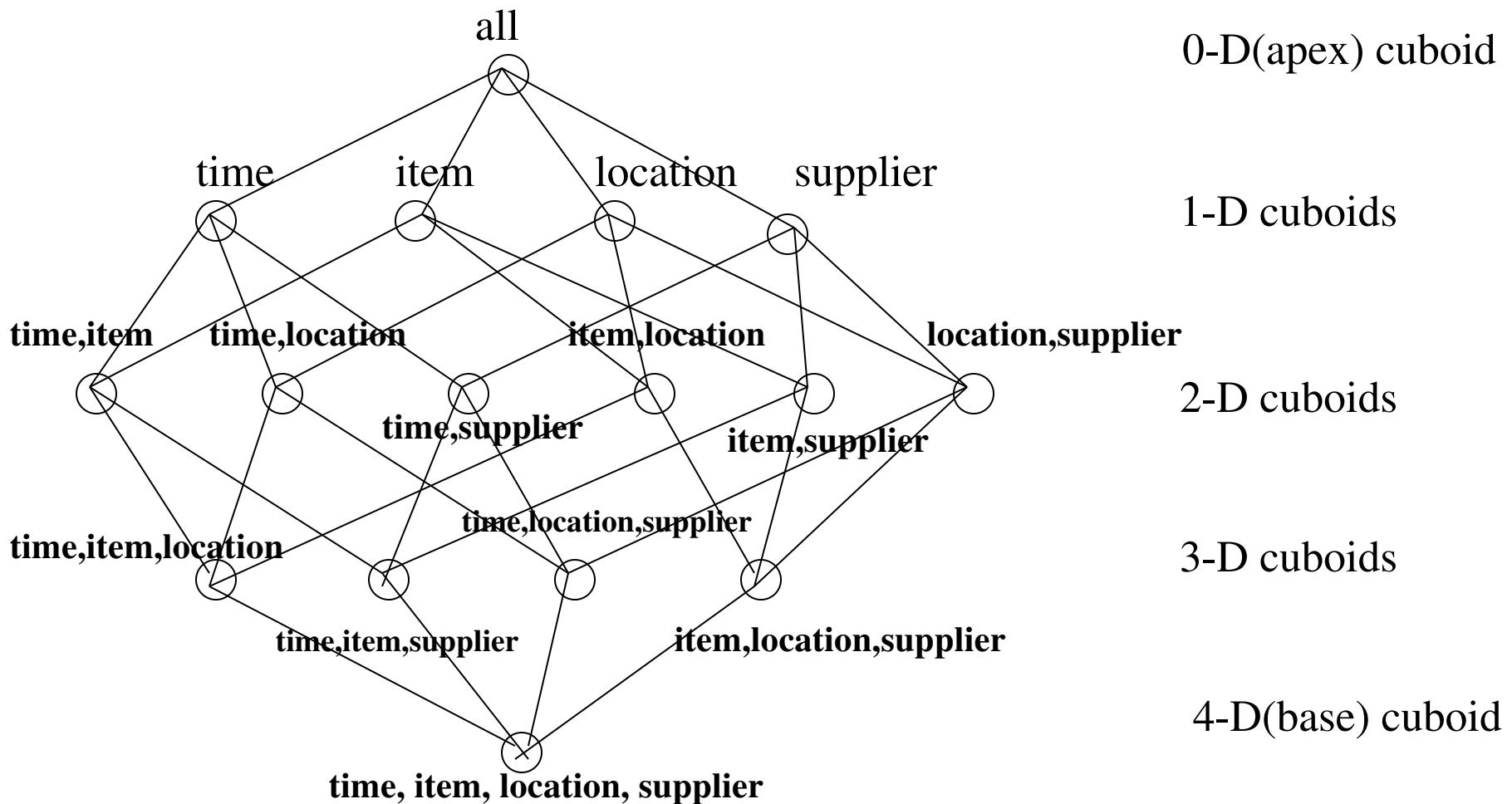
Multi Dimensional Data Model

Disadvantages of Multi Dimensional Data Model

The following are the disadvantages of a Multi Dimensional Data Model :

- The multi-dimensional Data Model is slightly complicated in nature and it requires professionals to recognize and examine the data in the database.
- During the work of a Multi-Dimensional Data Model, when the system caches, there is a great effect on the working of the system.
- It is complicated in nature due to which the databases are generally dynamic in design.
- The path to achieving the end product is complicated most of the time.

Cube: A Lattice of Cuboids

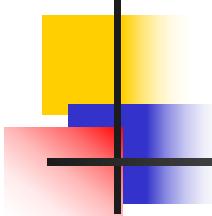


Data Warehousing - Schemas

Schema is a logical description of the entire database. It includes the name and description of records of all record types including all associated data-items and aggregates. Much like a database, a data warehouse also requires to maintain a schema. A database uses relational model, while a data warehouse uses Star, Snowflake, and Fact Constellation schema.

n **Star schema:**

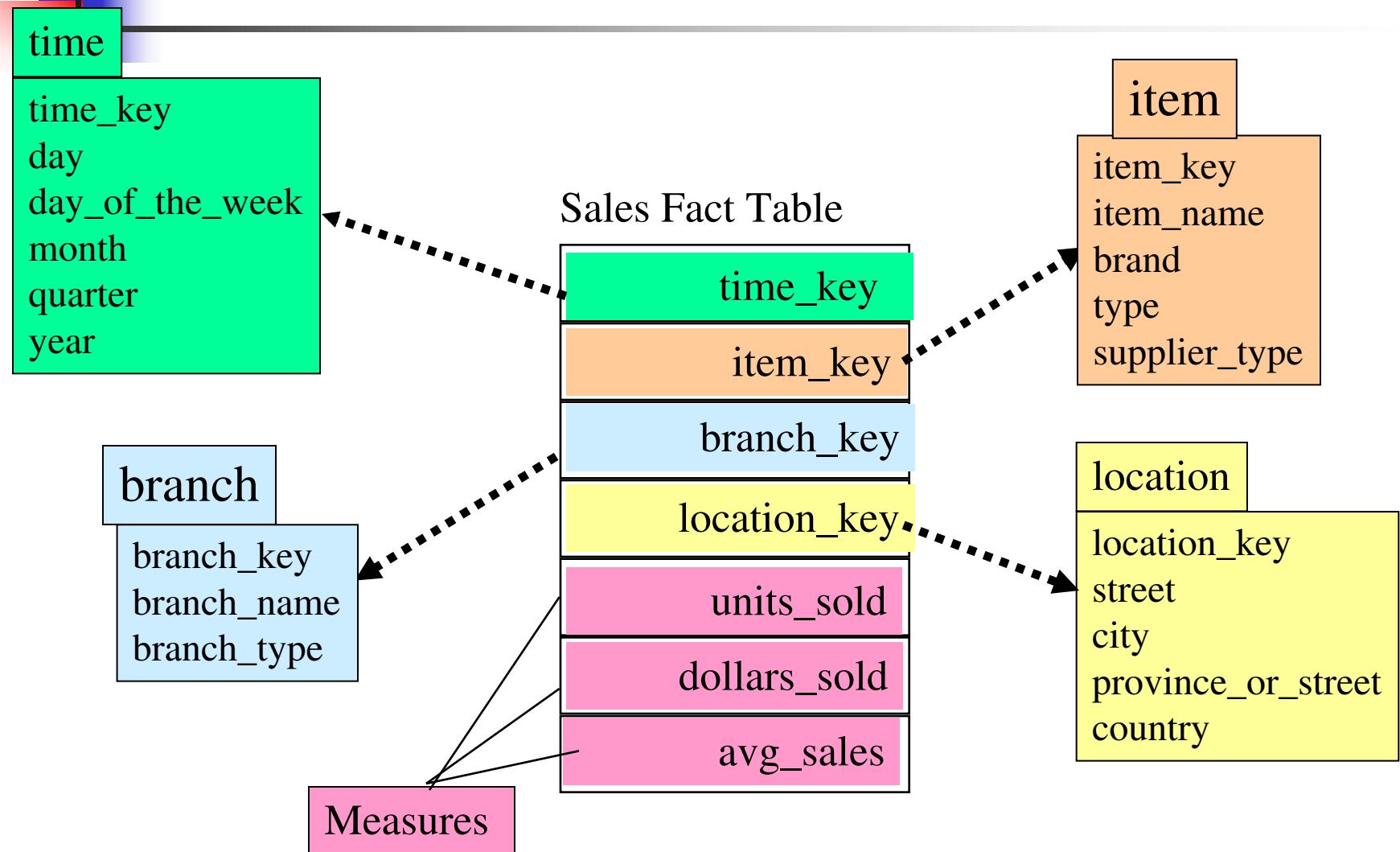
- A fact table in the middle connected to a set of dimension tables
- Each dimension in a star schema is represented with only one-dimension table
- This dimension table contains the set of attributes.
- The following diagram shows the sales data of a company with respect to the four dimensions, namely time, item, branch, and location.

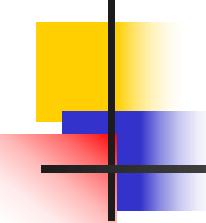


Conceptual Modeling of Data Warehouses

- n Modeling data warehouses: dimensions & measures
 - n **Star schema**: A fact table in the middle connected to a set of dimension tables
 - n **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
 - n **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





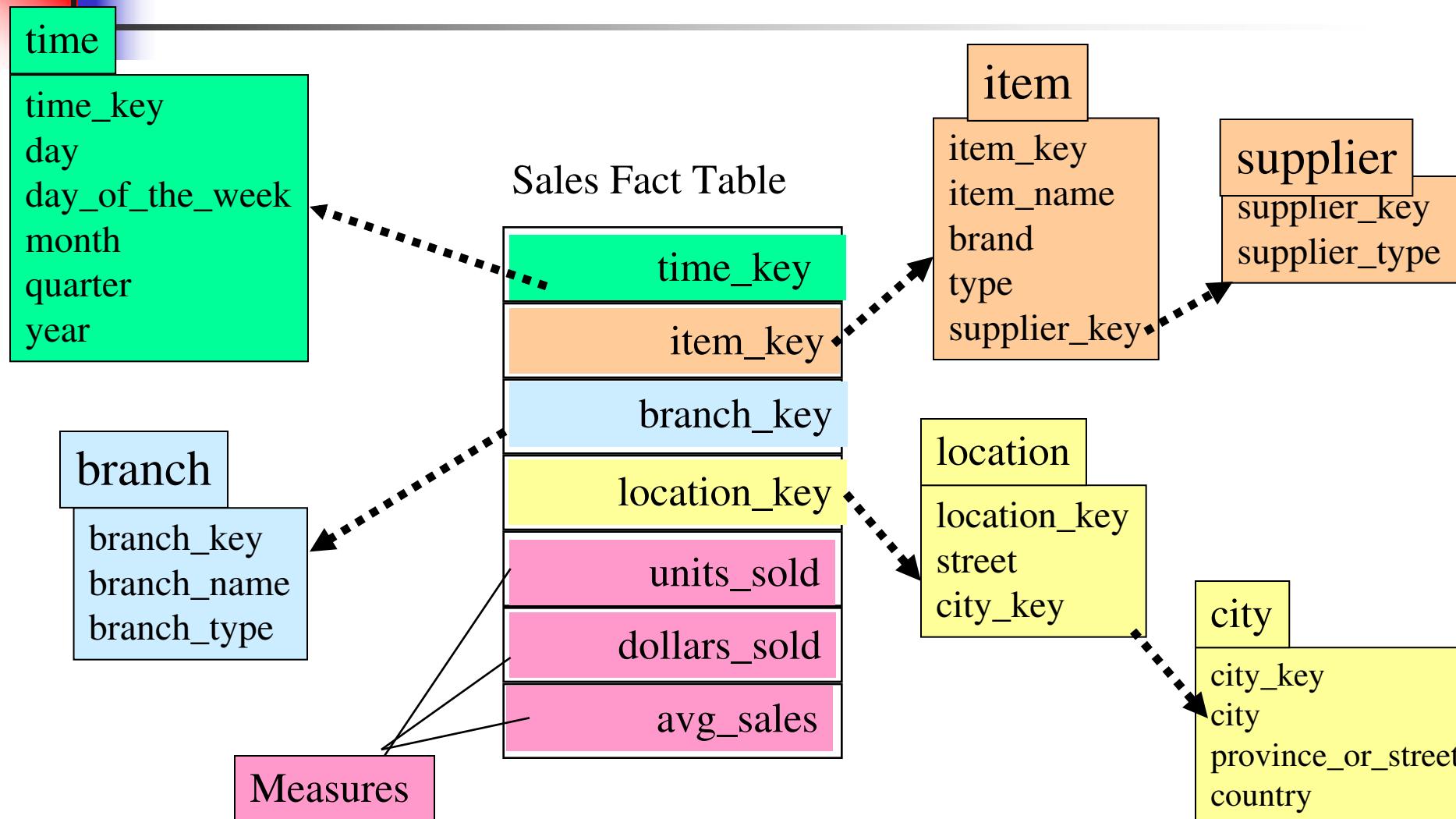
Star Schema

- There is a fact table at the center. It contains the keys to each of four dimensions.
- The fact table also contains the attributes, namely dollars sold and units sold.
- Note – Each dimension has only one dimension table and each table holds a set of attributes. For example, the location dimension table contains the attribute set {location_key, street, city, province_or_state, country}. This constraint may cause data redundancy. For example, "Vancouver" and "Victoria" both the cities are in the Canadian province of British Columbia. The entries for such cities may cause data redundancy along the attributes province_or_state and country.

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
- Some dimension tables in the Snowflake schema are normalized.
- The normalization splits up the data into additional tables
- Unlike Star schema, the dimensions table in a snowflake schema are normalized. For example, the item dimension table in star schema is normalized and split into two dimension tables, namely item and supplier table.

Example of Snowflake Schema



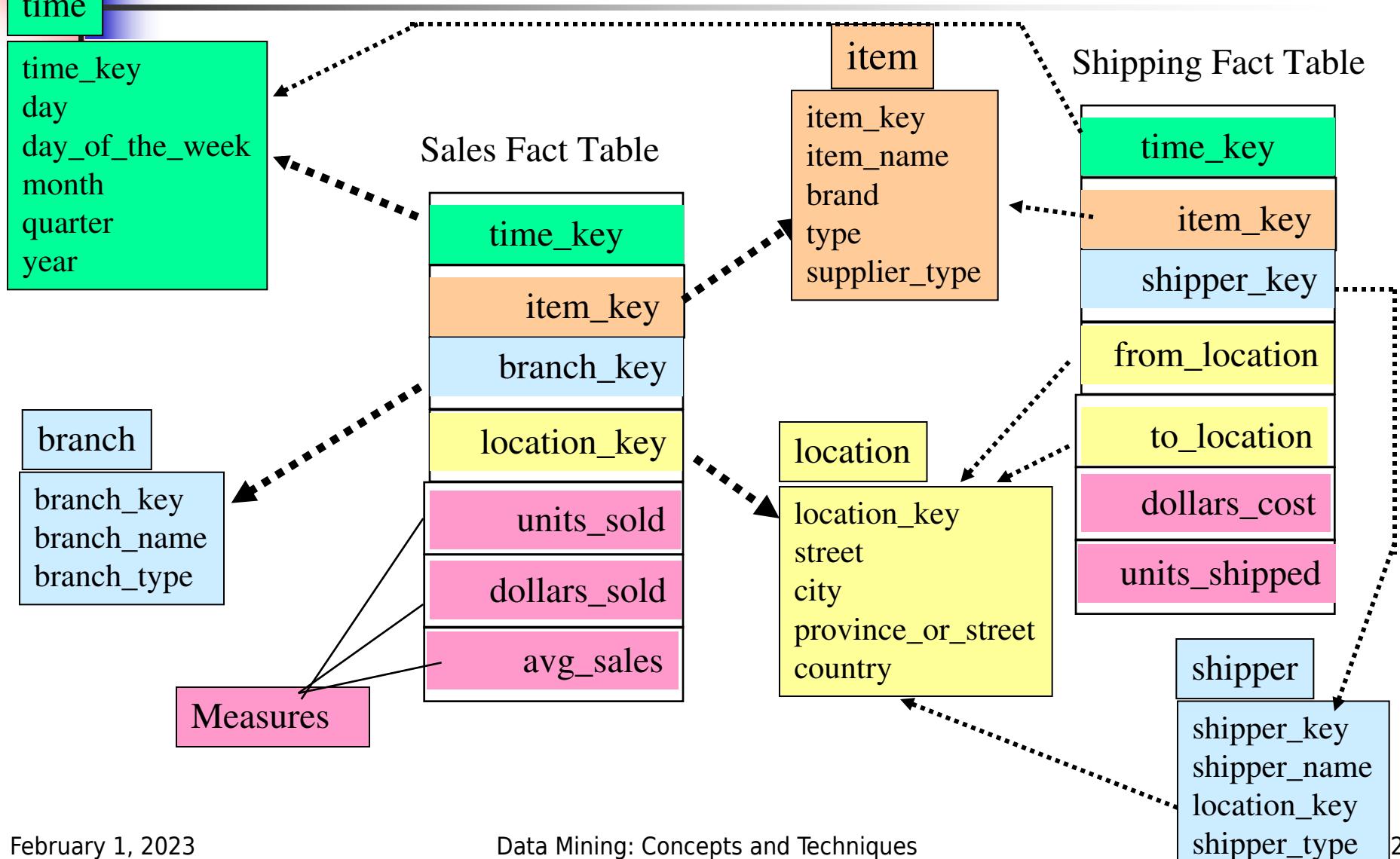
Snowflake Schema

- Now the item dimension table contains the attributes item_key, item_name, type, brand, and supplier-key.
 - The supplier key is linked to the supplier dimension table. The supplier dimension table contains the attributes supplier_key and supplier_type.
- n Note – Due to normalization in the Snowflake schema, the redundancy is reduced and therefore, it becomes easy to maintain and the save storage space.

Fact Constellation

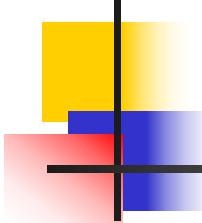
- Multiple fact tables share dimension tables, viewed as a collection of stars, therefore called **galaxy schema** or fact constellation
- A fact constellation has multiple fact tables. It is also known as galaxy schema.
- The following diagram shows two fact tables, namely sales and shipping.

Example of Fact Constellation



Fact Constellation

- The sales fact table is same as that in the star schema.
- The shipping fact table has the five dimensions, namely item_key, time_key, shipper_key, from_location, to_location.
- The shipping fact table also contains two measures, namely dollars sold and units sold.
- It is also possible to share dimension tables between fact tables. For example, time, item, and location dimension tables are shared between the sales and shipping fact table.



A Data Mining Query Language, DMQL: Language Primitives

- „ Cube Definition (Fact Table)

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

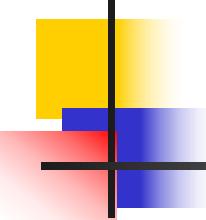
- „ Dimension Definition (Dimension Table)

```
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 a Star Schema in DMQL

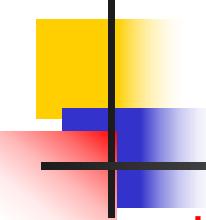
```
define cube sales_star [time, item, branch, location]:  
    dollars_sold = sum(sales_in_dollars), avg_sales =  
        avg(sales_in_dollars), units_sold = count(*)  
define dimension time as (time_key, day, day_of_week,  
    month, quarter, year)  
define dimension item as (item_key, item_name,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

- n **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.
 - n E.g., count(), sum(), min(), max().
- n **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.
 - n E.g., avg(), min_N(), standard_deviation().
- n **holistic**: if there is no constant bound on the storage size needed to describe a subaggregate.

A Concept Hierarchy: Dimension (location)

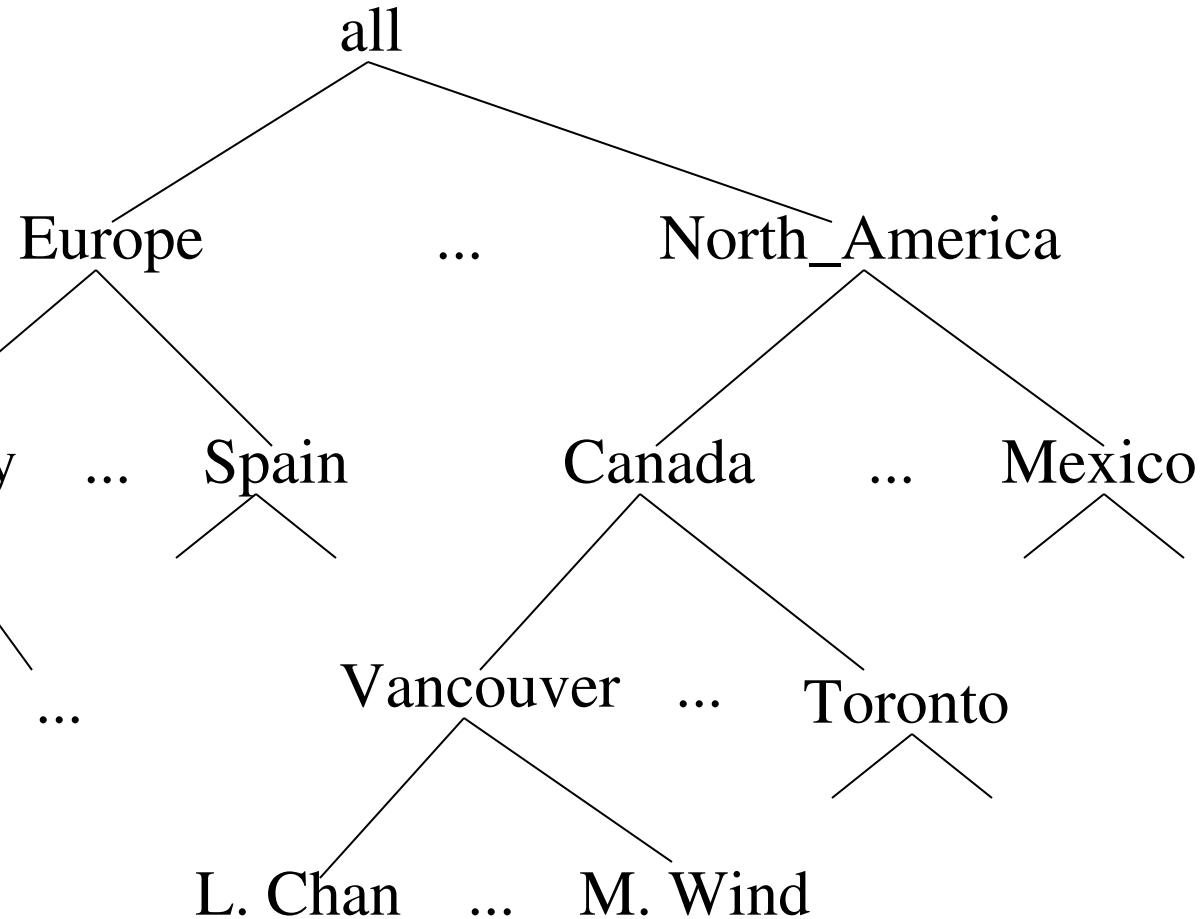
all

region

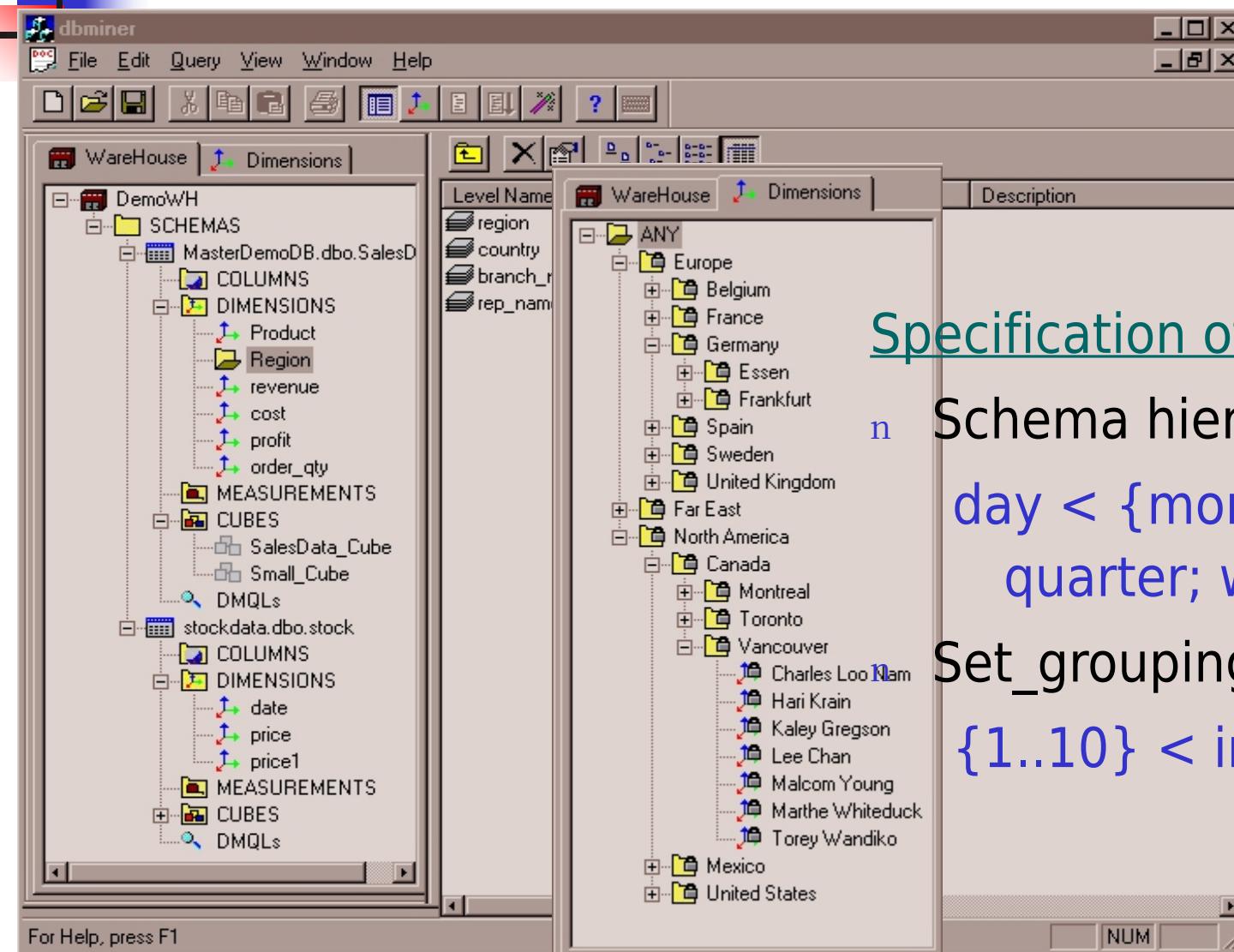
country

city

office



View of Warehouses and Hierarchies



Specification of hierarchies

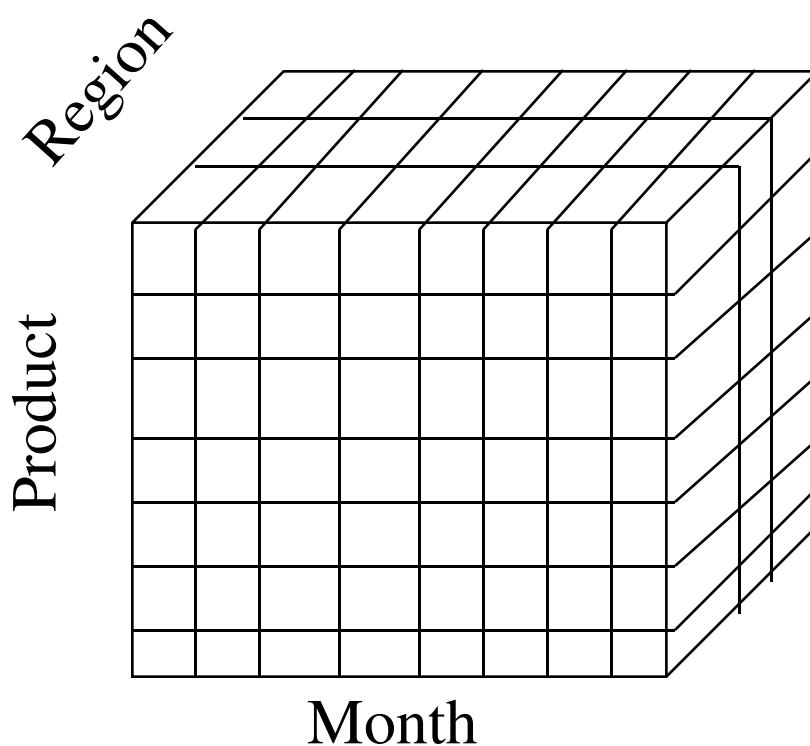
n Schema hierarchy
day < {month <
quarter; week} < year

Set_grouping hierarchy

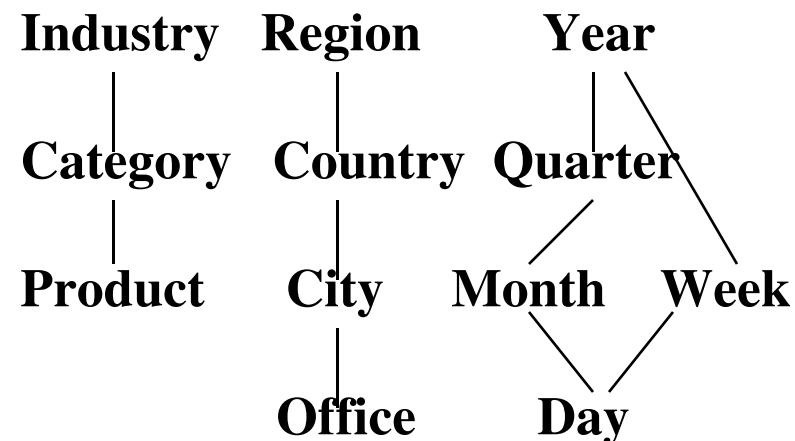
{1..10} < inexpensive

Multidimensional Data

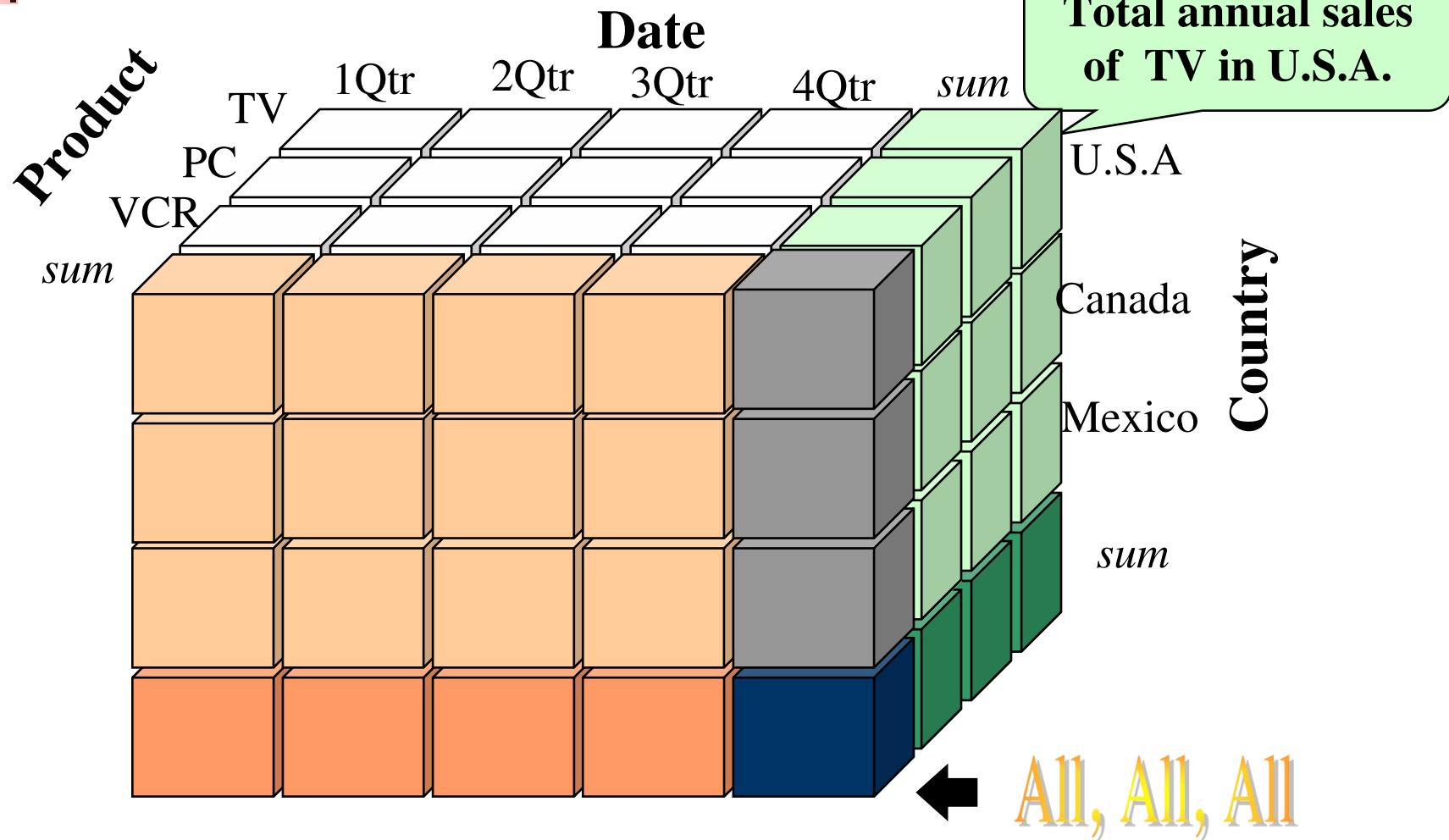
- n Sales volume as a function of product, month, and region



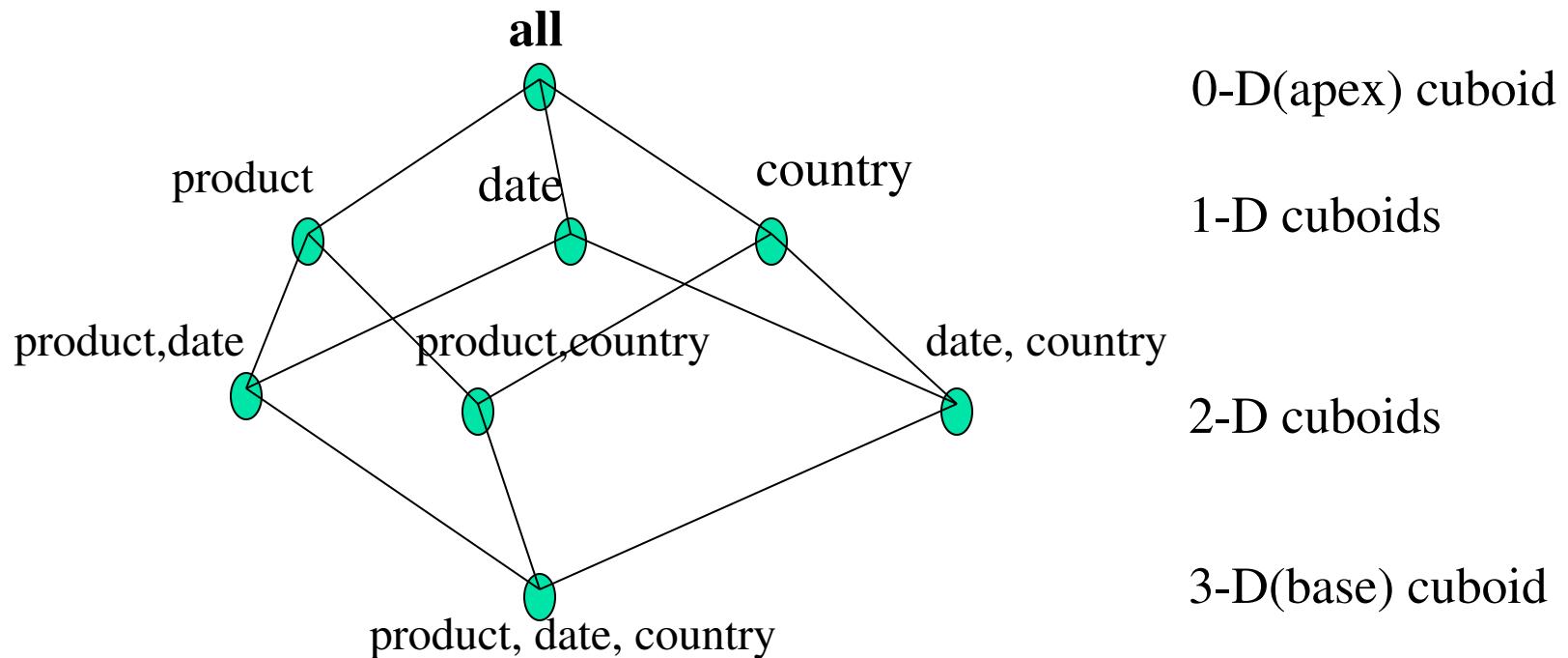
Dimensions: Product, Location, Time
Hierarchical summarization paths



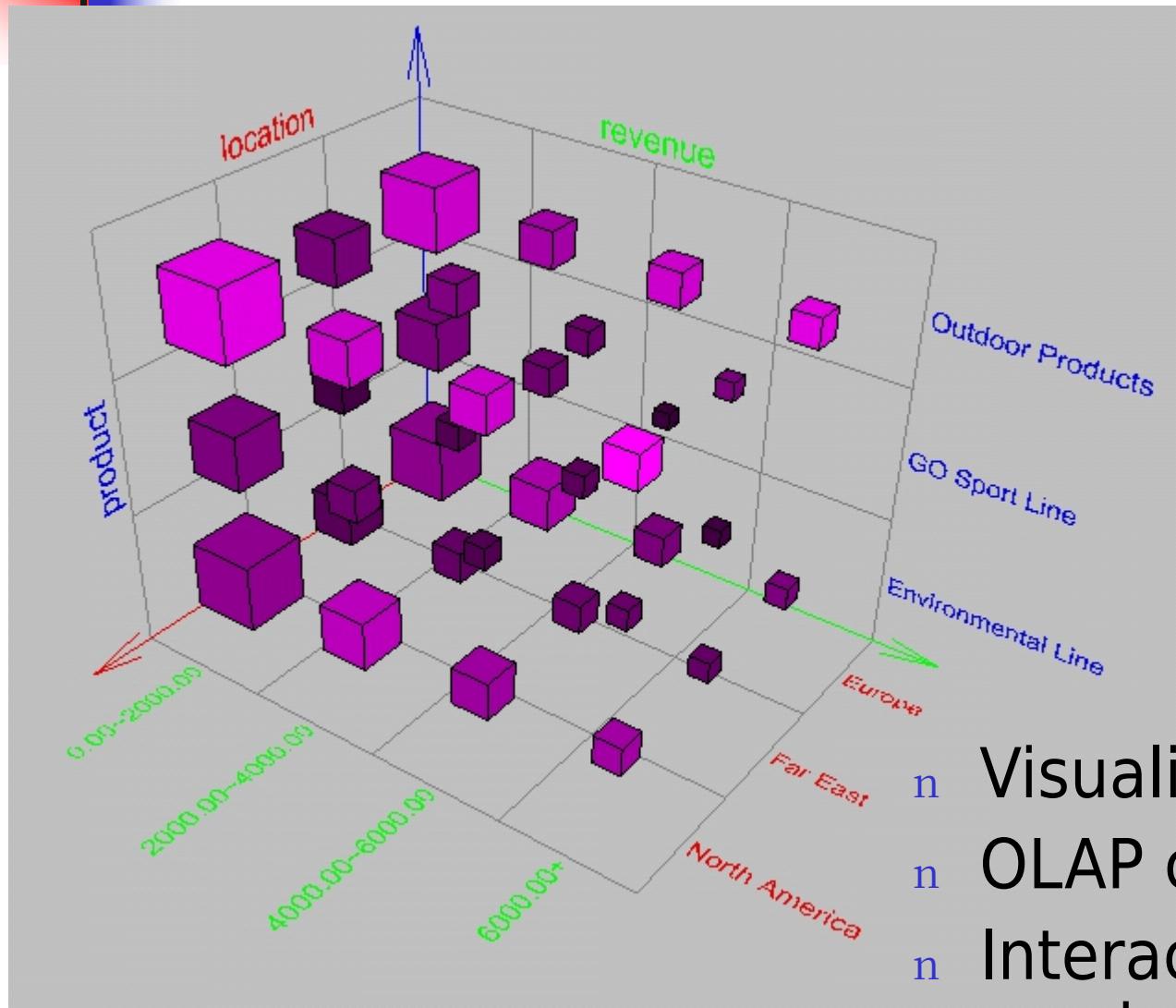
A Sample Data Cube



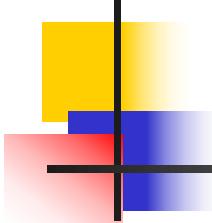
Cuboids Corresponding to the Cube



Browsing a Data Cube



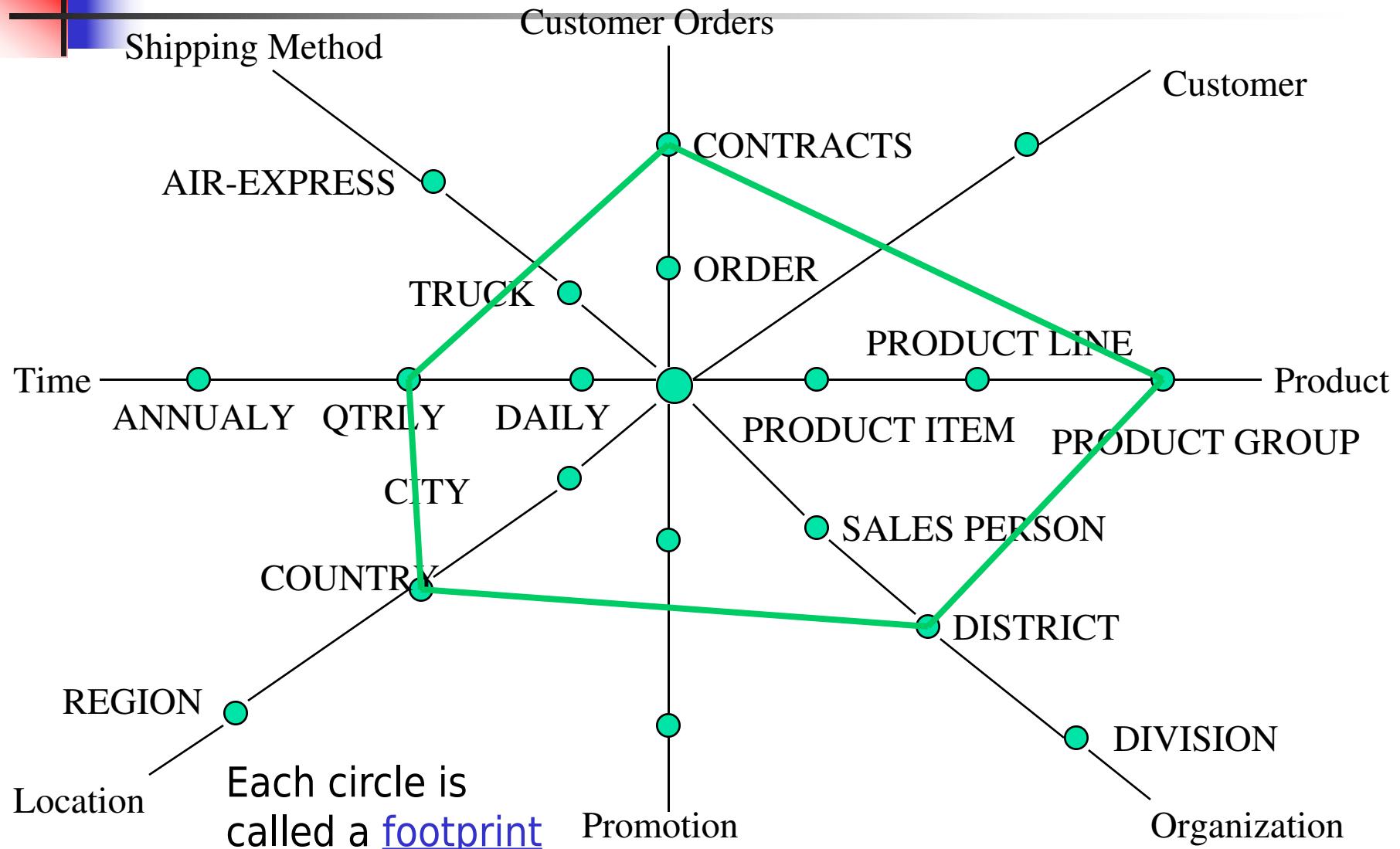
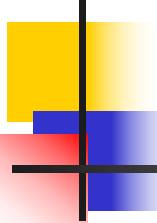
- n Visualization
- n OLAP capabilities
- n Interactive manipulation

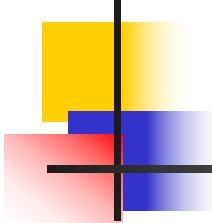


Typical OLAP Operations

- n Roll up (drill-up): summarize data
 - n *by climbing up hierarchy or by dimension reduction*
- n Drill down (roll down): reverse of roll-up
 - n *from higher level summary to lower level summary or detailed data, or introducing new dimensions*
- n Slice and dice:
 - n *project and select*
- n Pivot (rotate):
 - n *reorient the cube, visualization, 3D to series of 2D planes.*
- n Other operations
 - n *drill across: involving (across) more than one fact table*
 - n *drill through: through the bottom level of the cube to its back-end relational tables (using SQL)*

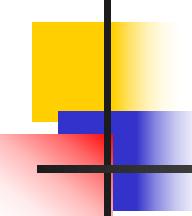
A Star-Net Query Model





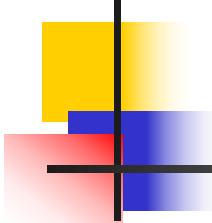
Chapter 2: Data Warehousing and OLAP Technology for Data Mining

- n What is a data warehouse?
- n A multi-dimensional data model
- n **Data warehouse architecture**
- n Data warehouse implementation
- n Further development of data cube technology
- n From data warehousing to data mining



Design of a Data Warehouse: A Business Analysis Framework

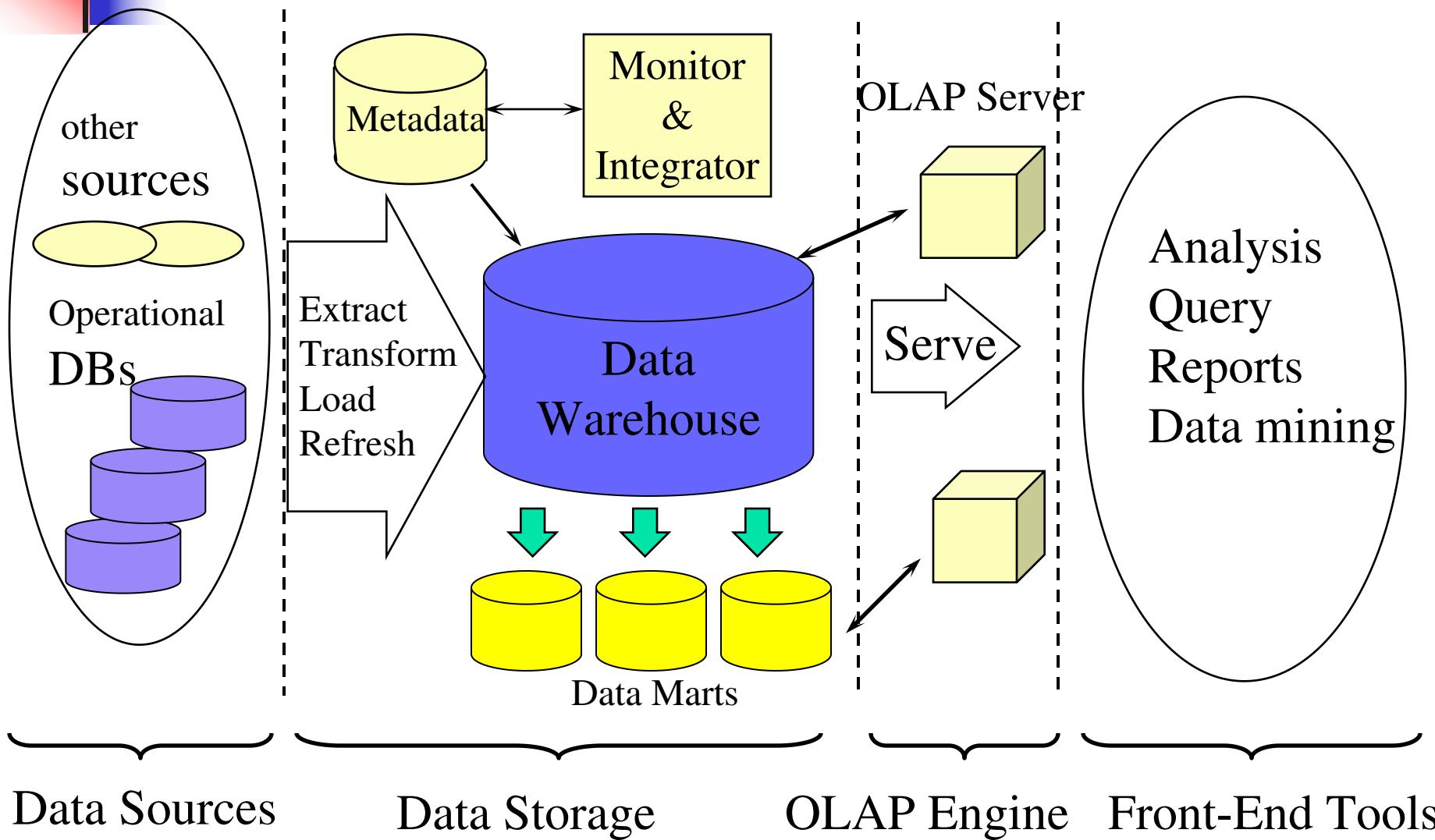
- n Four views regarding the design of a data warehouse
 - n **Top-down view**
 - n allows selection of the relevant information necessary for the data warehouse
 - n **Data source view**
 - n exposes the information being captured, stored, and managed by operational systems
 - n **Data warehouse view**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

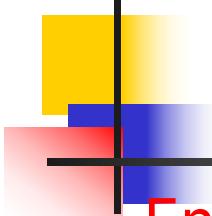


Data Warehouse Design Process

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

Multi-Tiered Architecture





Three Data Warehouse Models

n Enterprise warehouse

- n collects all of the information about subjects spanning the entire organization

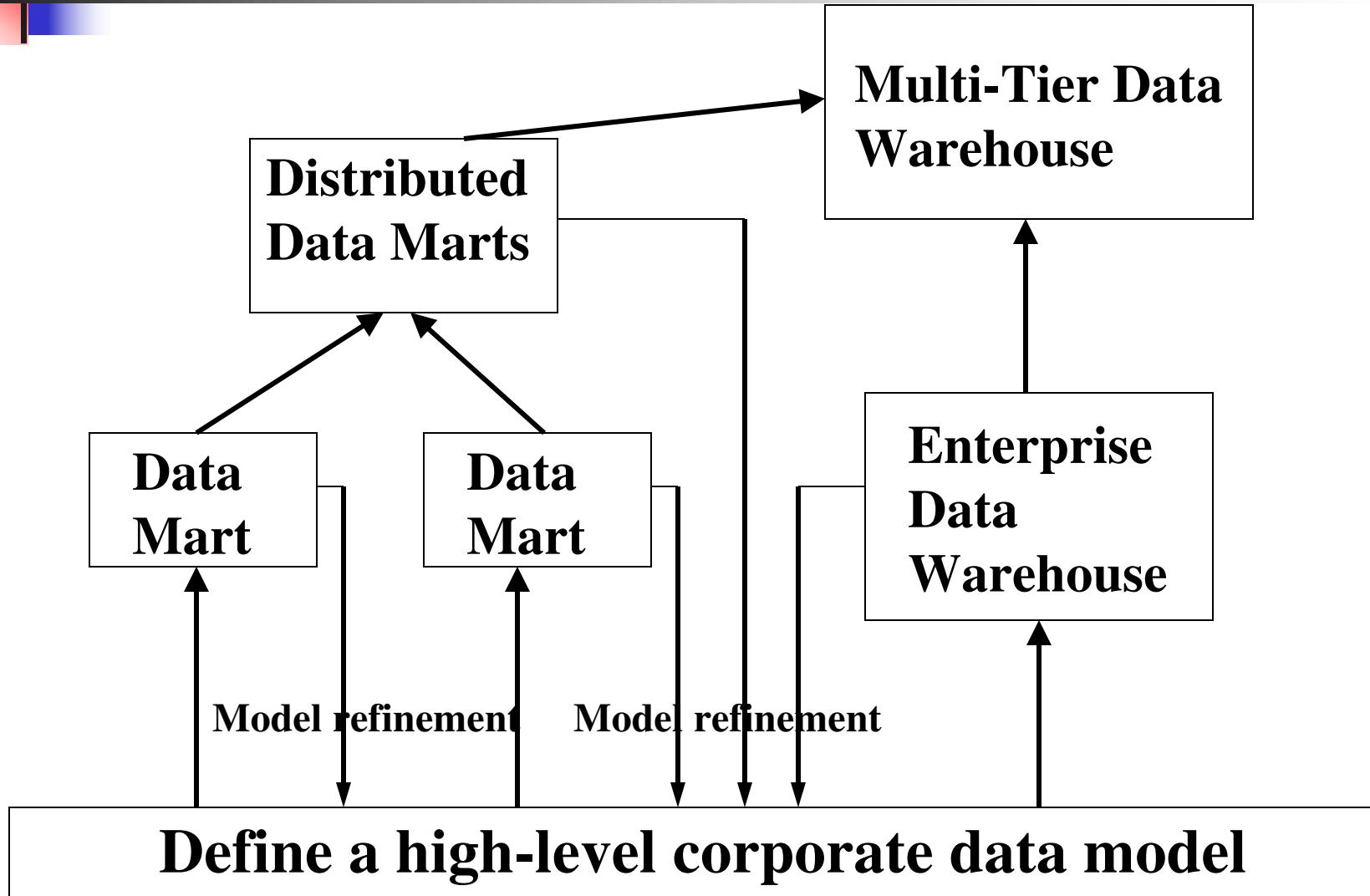
n Data Mart

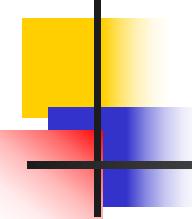
- n a subset of corporate-wide data that is of value to a specific groups of users. Its scope is confined to specific, selected groups, such as marketing data mart
 - n Independent vs. dependent (directly from warehouse) data mart

n Virtual warehouse

- n A set of views over operational databases
- n Only some of the possible summary views may be materialized

Data Warehouse Development: A Recommended Approach





OLAP Server Architectures

n Relational OLAP (ROLAP)

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

n Multidimensional OLAP (MOLAP)

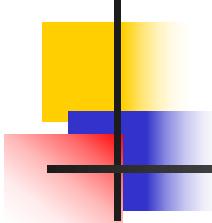
- n Array-based multidimensional storage engine (sparse matrix techniques)
- n fast indexing to pre-computed summarized data

n Hybrid OLAP (HOLAP)

- n User flexibility, e.g., low level: relational, high-level: array

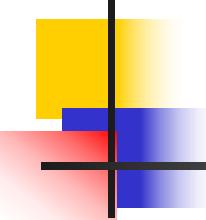
n Specialized SQL servers

specialized support for SQL queries over star/snowflake schemas



Chapter 2: Data Warehousing and OLAP Technology for Data Mining

- n What is a data warehouse?
- n A multi-dimensional data model
- n Data warehouse architecture
- n **Data warehouse implementation**
- n Further development of data cube technology
- n From data warehousing to data mining



Efficient Data Cube Computation

- n Data cube can be viewed as a lattice of cuboids
 - n The bottom-most cuboid is the base cuboid
 - n The top-most cuboid (apex) contains only one cell
 - n How many cuboids in an n-dimensional cube with L levels?
$$T = \prod_{i=1}^n (L_i + 1)$$
- n Materialization of data cube
 - n Materialize every (cuboid) (full materialization), none (no materialization), or some (partial materialization)
 - n Selection of which cuboids to materialize
 - n Based on size, sharing, access frequency, etc.

Cube Operation

- „ Cube definition and computation in DMQL

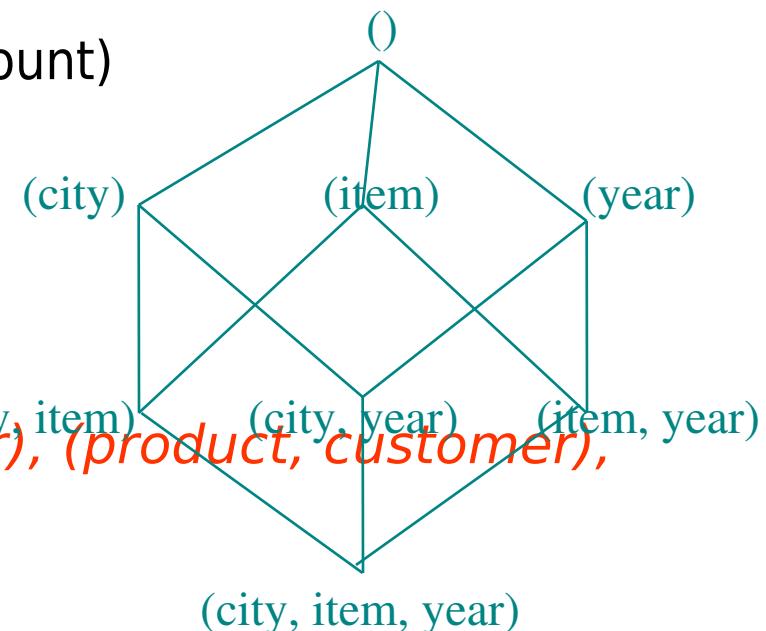
```
define cube sales[item, city, year]: sum(sales_in_dollars)  
compute cube sales
```

- „ 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
```

- „ Need compute the following Group-Bys

*(date, product, customer),
(date,product),(date, customer), (product, customer),
(date), (product), (customer)
()*



Cube Computation: ROLAP-Based Method

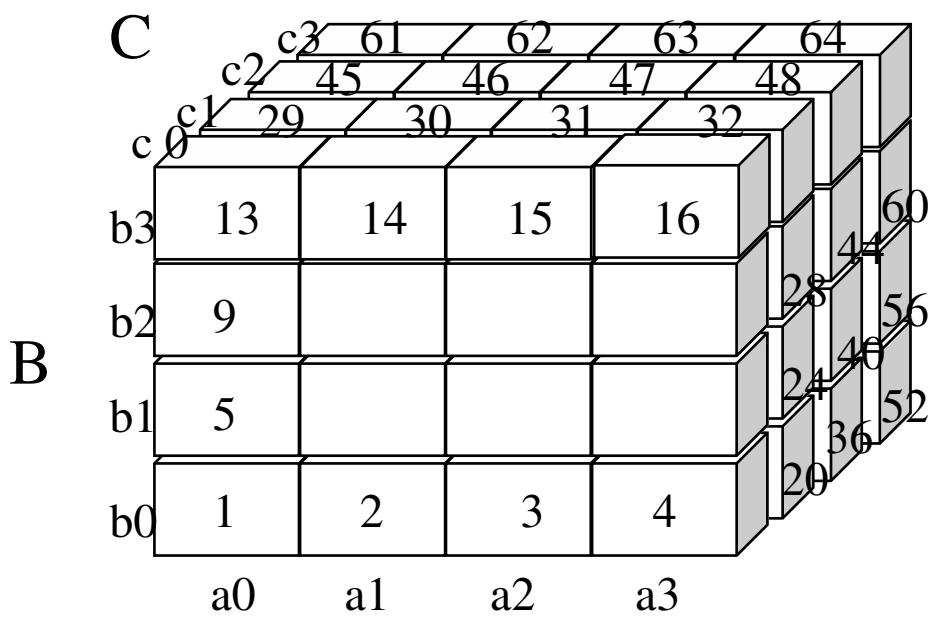
- n Efficient cube computation methods
 - n ROLAP-based cubing algorithms (Agarwal et al'96)
 - n Array-based cubing algorithm (Zhao et al'97)
 - n Bottom-up computation method (Bayer & Ramakrishnan'99)
- n ROLAP-based cubing algorithms
 - n Sorting, hashing, and grouping operations are applied to the dimension attributes in order to reorder and cluster related tuples
 - n Grouping is performed on some subaggregates as a “partial grouping step”
 - n Aggregates may be computed from previously computed aggregates, rather than from the base fact table

Cube Computation: ROLAP-Based Method (2)

- n This is not in the textbook but in a research paper
- n Hash/sort based methods (Agarwal et. al. VLDB'96)
 - n **Smallest-parent:** computing a cuboid from the smallest cuboid previously computed cuboid.
 - n **Cache-results:** caching results of a cuboid from which other cuboids are computed to reduce disk I/Os
 - n **Amortize-scans:** computing as many as possible cuboids at the same time to amortize disk reads
 - n **Share Sorts:** sharing sorting costs cross multiple cuboids when sort-based method is used
 - n **Share-partitions:** sharing the partitioning cost cross multiple cuboids when hash-based algorithms are used

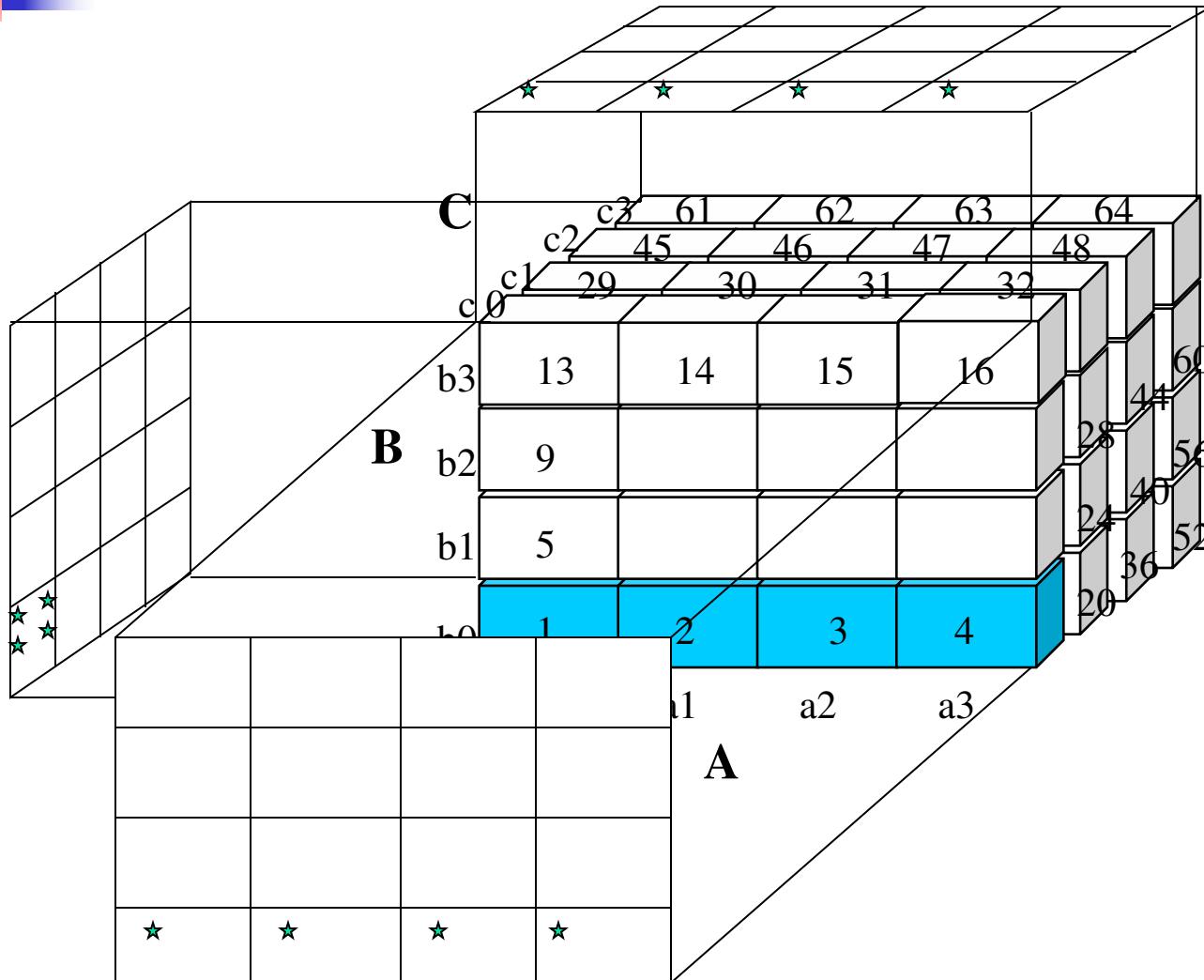
Multi-way Array Aggregation for Cube Computation

- Partition arrays into chunks (a small subcube which fits in memory).
- Compressed sparse array addressing: (chunk_id, offset)
- Compute aggregates in “multiway” by visiting cube cells in the order which minimizes the # of times to visit each cell, and reduces memory access and storage cost.

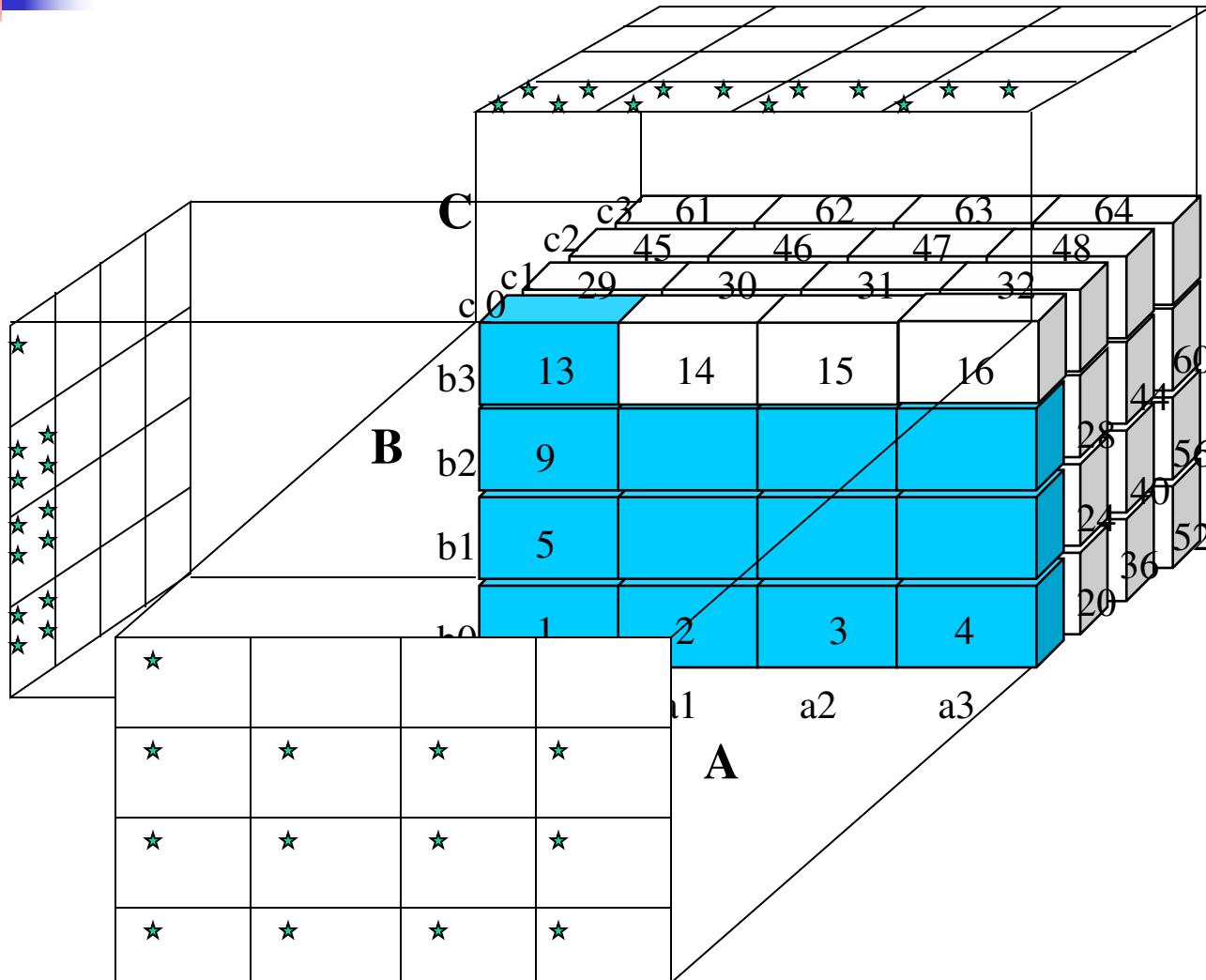


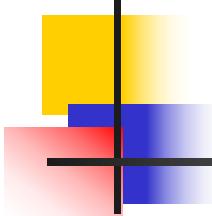
What is the best traversing order to do multi-way aggregation?

Multi-way Array Aggregation for Cube Computation



Multi-way Array Aggregation for Cube Computation





Multi-Way Array Aggregation for Cube Computation (Cont.)

- Method: the planes should be sorted and computed according to their size in ascending order.
 - See the details of Example 2.12 (pp. 75-78)
 - Idea: keep the smallest plane in the main memory, fetch and compute only one chunk at a time for the largest plane
- Limitation of the method: computing well only for a small number of dimensions
 - If there are a large number of dimensions, “bottom-up computation” and iceberg cube computation methods can be explored

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

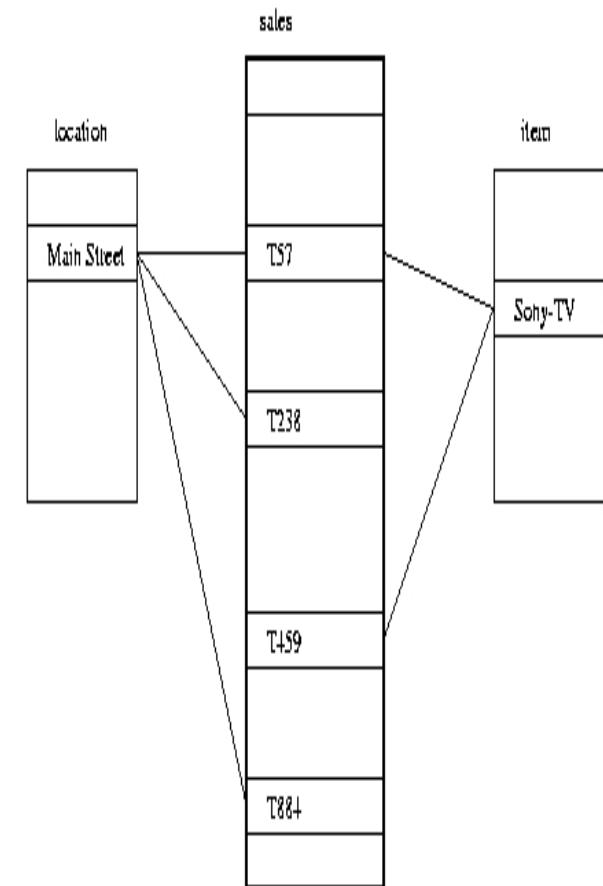
RecID	Asia	Europe	America
1	1	0	0
2	0	1	0
3	1	0	0
4	0	0	1
5	0	1	0

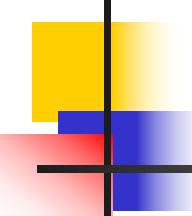
Index on Type

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

Indexing OLAP Data: Join Indices

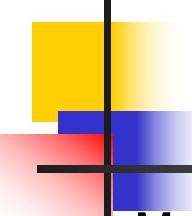
- n Join index: $JI(R\text{-id}, S\text{-id}) \text{ where } R(R\text{-id}, \dots) \bowtie S(S\text{-id}, \dots)$
- n Traditional indices map the values to a list of record ids
 - n It materializes relational join in JI file and speeds up relational join — a rather costly operation
- n In data warehouses,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





Efficient Processing OLAP Queries

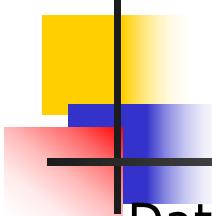
- „ Determine which operations should be performed on the available cuboids:
 - „ transform drill, roll, etc. into corresponding SQL and/or OLAP operations, e.g, dice = selection + projection
- „ Determine to which materialized cuboid(s) the relevant operations should be applied.
- „ Exploring indexing structures and compressed vs. dense array structures in MOLAP



Metadata Repository

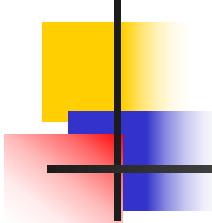
n Meta data is the data defining warehouse objects. It has the following kinds

- n Description of the structure of the warehouse
 - n schema, view, dimensions, hierarchies, derived data defn, data mart locations and contents
- n Operational meta-data
 - n 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)
- n The algorithms used for summarization
- n The mapping from operational environment to the data warehouse
- n Data related to system performance
 - n warehouse schema, view and derived data definitions
- n Business data



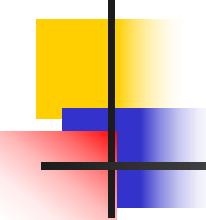
Data Warehouse Back-End Tools and Utilities

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



Chapter 2: Data Warehousing and OLAP Technology for Data Mining

- n What is a data warehouse?
- n A multi-dimensional data model
- n Data warehouse architecture
- n Data warehouse implementation
- n **Further development of data cube technology**
- n From data warehousing to data mining



Discovery-Driven Exploration of Data Cubes

- Hypothesis-driven: exploration by user, huge search space
- Discovery-driven (Sarawagi et al.'98)
 - pre-compute measures indicating exceptions, guide user in the data analysis, at all levels of aggregation
 - Exception: significantly different from the value anticipated, based on a statistical model
 - Visual cues such as background color are used to reflect the degree of exception of each cell
 - Computation of exception indicator (modeling fitting and computing SelfExp, InExp, and PathExp values) can be overlapped with cube construction

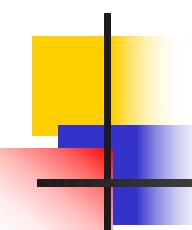
Examples: Discovery-Driven Data Cubes

item	all
region	all

Sum of sales		month											
		Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
Total		1%	-1%	0%	1%	3%	-1	-9%	-1%	2%	-4%	3%	

Avg sales		month											
item		Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
Sony b/w printer		9%	-8%	2%	-5%	14%	-4%	0%	41%	-13%	-15%	-11%	
Sony color printer		0%	0%	3%	2%	4%	-10%	-13%	0%	4%	-6%	4%	
HP b/w printer		-2%	1%	2%	3%	8%	0%	-12%	-9%	3%	-3%	6%	
HP color printer		0%	0%	-2%	1%	0%	-1%	-7%	-2%	1%	-5%	1%	
IBM home computer		1%	-2%	-1%	-1%	3%	3%	-10%	4%	1%	-4%	-1%	
IBM laptop computer		0%	0%	-1%	3%	4%	2%	-10%	-2%	0%	-9%	3%	
Toshiba home computer		-2%	-5%	1%	1%	-1%	1%	5%	-3%	-5%	-1%	-1%	
Toshiba laptop computer		1%	0%	3%	0%	-2%	-2%	-5%	3%	2%	-1%	0%	
Logitech mouse		3%	-2%	-1%	0%	4%	6%	-11%	2%	1%	-4%	0%	
Ergo-way mouse		0%	0%	2%	3%	1%	-2%	-2%	-5%	0%	-5%	8%	

item	IBM home computer											
Avg sales	month											
region	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
North		-1%	-3%	-1%	0%	3%	4%	-7%	1%	0%	-3%	-3%
South		-1%	1%	-9%	6%	-1%	-39%	9%	-34%	4%	1%	7%
East		-1%	-2%	2%	-3%	1%	18%	-2%	11%	-3%	-2%	-1%
West		4%	0%	-1%	-3%	5%	1%	-18%	8%	5%	-8%	1%

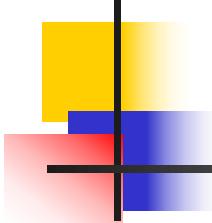


Complex Aggregation at Multiple Granularities: Multi-Feature Cubes

- n Multi-feature cubes (Ross, et al. 1998): Compute complex queries involving multiple dependent aggregates at multiple granularities
- n Ex. Grouping by all subsets of {item, region, month}, find the maximum price in 1997 for each group, and the total sales among all maximum price tuples

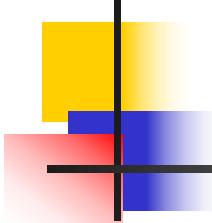
```
select item, region, month, max(price), sum(R.sales)  
from purchases  
where year = 1997  
cube by item, region, month: R  
such that R.price = max(price)
```

- n Continuing the last example, among the max price tuples, find the min and max shelf life, and find the fraction of the total sales due to tuple that have min shelf life within the set of all max price tuples



Chapter 2: Data Warehousing and OLAP Technology for Data Mining

- n What is a data warehouse?
- n A multi-dimensional data model
- n Data warehouse architecture
- n Data warehouse implementation
- n Further development of data cube technology
- n From data warehousing to data mining



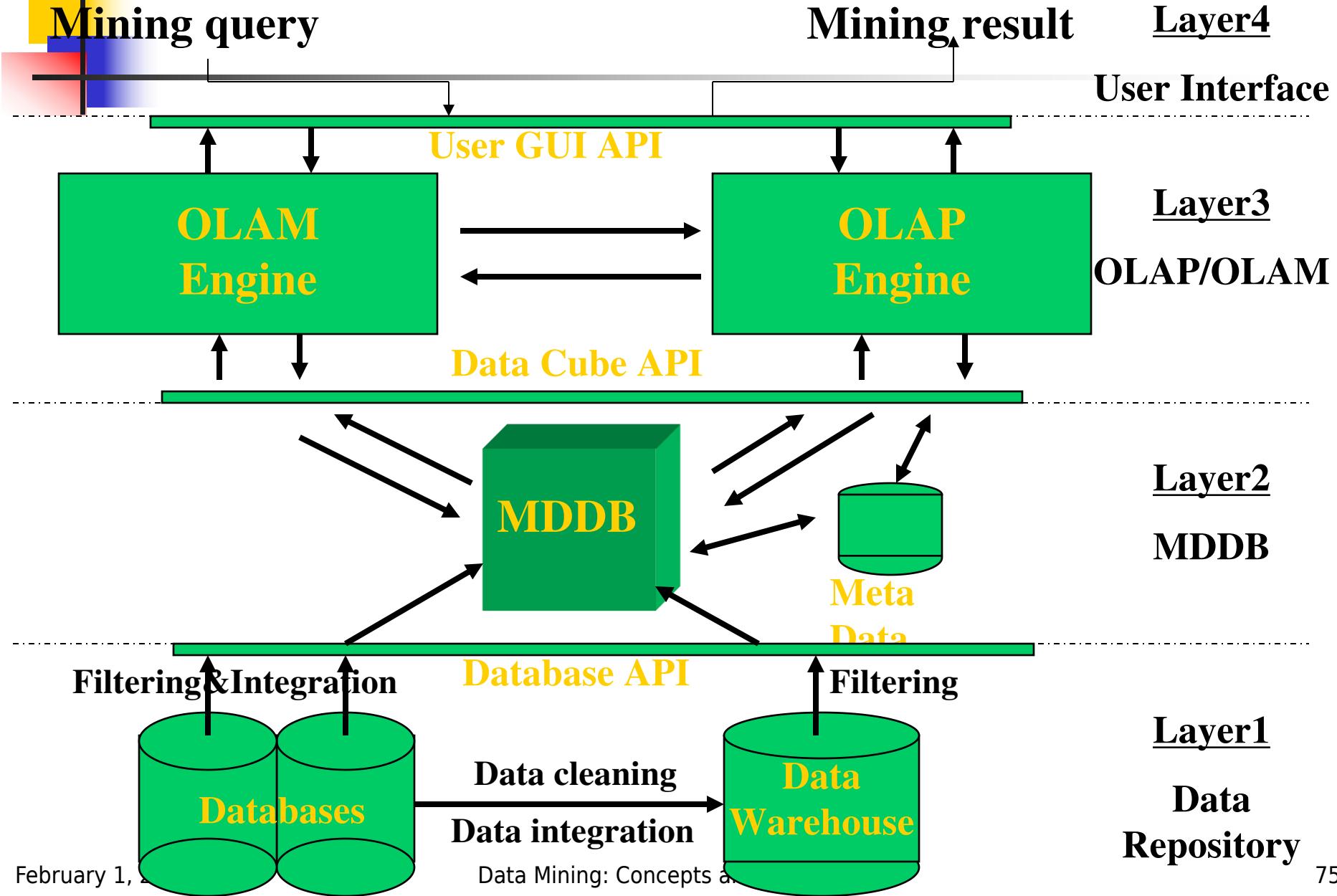
Data Warehouse Usage

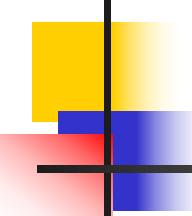
- n Three kinds of data warehouse applications
 - n **Information processing**
 - n supports querying, basic statistical analysis, and reporting using crosstabs, tables, charts and graphs
 - n **Analytical processing**
 - n multidimensional analysis of data warehouse data
 - n supports basic OLAP operations, slice-dice, drilling, pivoting
 - n **Data mining**
 - n knowledge discovery from hidden patterns
 - n supports associations, constructing analytical models, performing classification and prediction, and presenting the mining results using visualization tools.

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

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

An OLAM Architecture



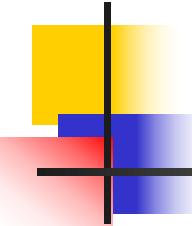


Summary

- n **Data warehouse**
 - n A subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process
- n A **multi-dimensional model** of a data warehouse
 - n Star schema, snowflake schema, fact constellations
 - n A data cube consists of dimensions & measures
- n **OLAP** operations: drilling, rolling, slicing, dicing and pivoting
- n OLAP servers: ROLAP, MOLAP, HOLAP
- n Efficient computation of data cubes
 - n Partial vs. full vs. no materialization
 - n Multiway array aggregation
 - n Bitmap index and join index implementations
- n Further development of data cube technology
 - n Discovery-drive and multi-feature cubes
 - n From OLAP to OLAM (on-line analytical mining)

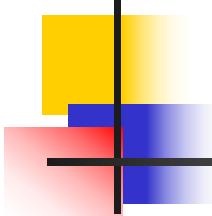
References (I)

- n S. Agarwal, R. Agrawal, P. M. Deshpande, A. Gupta, J. F. Naughton, R. Ramakrishnan, and S. Sarawagi. On the computation of multidimensional aggregates. In Proc. 1996 Int. Conf. Very Large Data Bases, 506-521, Bombay, India, Sept. 1996.
- n D. Agrawal, A. E. Abbadi, A. Singh, and T. Yurek. Efficient view maintenance in data warehouses. In Proc. 1997 ACM-SIGMOD Int. Conf. Management of Data, 417-427, Tucson, Arizona, May 1997.
- n R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. In Proc. 1998 ACM-SIGMOD Int. Conf. Management of Data, 94-105, Seattle, Washington, June 1998.
- n R. Agrawal, A. Gupta, and S. Sarawagi. Modeling multidimensional databases. In Proc. 1997 Int. Conf. Data Engineering, 232-243, Birmingham, England, April 1997.
- n K. Beyer and R. Ramakrishnan. Bottom-Up Computation of Sparse and Iceberg CUBEs. In Proc. 1999 ACM-SIGMOD Int. Conf. Management of Data (SIGMOD'99), 359-370, Philadelphia, PA, June 1999.
- n S. Chaudhuri and U. Dayal. An overview of data warehousing and OLAP technology. ACM SIGMOD Record, 26:65-74, 1997.
- n OLAP council. MD API specification version 2.0. In <http://www.olapcouncil.org/research/api.htm>, 1998.
- n J. Gray, S. Chaudhuri, A. Bosworth, A. Layman, D. Reichart, M. Venkatrao, F. Pellow, and H. Pirahesh. Data cube: A relational aggregation operator generalizing group-by, cross-tab and sub-totals. Data Mining and Knowledge Discovery, 1:29-54, 1997.



References (II)

- V. Harinarayan, A. Rajaraman, and J. D. Ullman. Implementing data cubes efficiently. In Proc. 1996 ACM-SIGMOD Int. Conf. Management of Data, pages 205-216, Montreal, Canada, June 1996.
- Microsoft. OLEDB for OLAP programmer's reference version 1.0. In <http://www.microsoft.com/data/oledb/olap>, 1998.
- K. Ross and D. Srivastava. Fast computation of sparse datacubes. In Proc. 1997 Int. Conf. Very Large Data Bases, 116-125, Athens, Greece, Aug. 1997.
- K. A. Ross, D. Srivastava, and D. Chatziantoniou. Complex aggregation at multiple granularities. In Proc. Int. Conf. of Extending Database Technology (EDBT'98), 263-277, Valencia, Spain, March 1998.
- S. Sarawagi, R. Agrawal, and N. Megiddo. Discovery-driven exploration of OLAP data cubes. In Proc. Int. Conf. of Extending Database Technology (EDBT'98), pages 168-182, Valencia, Spain, March 1998.
- E. Thomsen. OLAP Solutions: Building Multidimensional Information Systems. John Wiley & Sons, 1997.
- Y. Zhao, P. M. Deshpande, and J. F. Naughton. An array-based algorithm for simultaneous multidimensional aggregates. In Proc. 1997 ACM-SIGMOD Int. Conf. Management of Data, 159-170, Tucson, Arizona, May 1997.



<http://www.cs.sfu.ca/~han>



Thank you !!!