# LECTURE 11

## DATA WAREHOUSE AND OLAP

# WHAT IS A 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.

- 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

# DATA WAREHOUSE—INTEGRATED

- Constructed by integrating multiple, heterogeneous data sources
  - relational databases, flat files, on-line transaction records
- Data cleaning and data integration techniques are applied.
  - Ensure consistency in naming conventions, encoding structures, attribute measures, etc. among different data sources
    - 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"
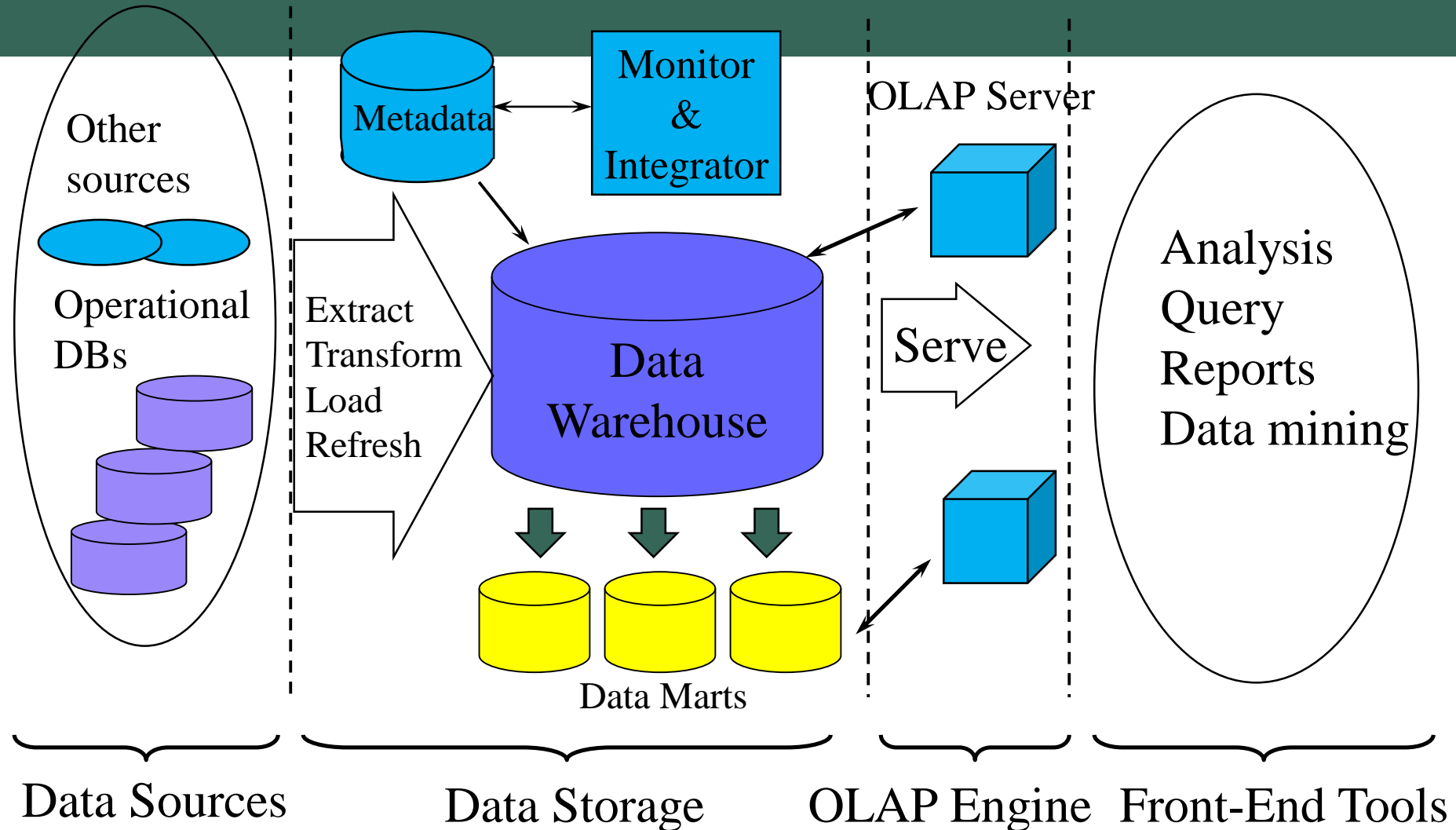
# DATA WAREHOUSE—NONVOLATILE

- 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*

# WHY A SEPARATE DATA WAREHOUSE?

- High performance for both systems

  - DBMS— tuned for OLTP: access methods, indexing, concurrency control, recovery

  - Warehouse—tuned for OLAP: complex OLAP queries, multidimensional view, consolidation

- Different functions and different data:

  - missing data: Decision support requires historical data which operational DBs do not typically maintain

  - data consolidation:  DS requires consolidation (aggregation, summarization) of data from heterogeneous sources

  - data quality: different sources typically use inconsistent data representations, codes and formats which have to be reconciled

- Note: There are more and more systems which perform OLAP analysis directly on relational databases

# DATA WAREHOUSE: A MULTI-TIERED ARCHITECTURE



Other sources

Metadata

Monitor & Integrator

OLAP Server

Operational DBs

Extract Transform Load Refresh

Data Warehouse

Serve

Analysis Query Reports Data mining

Data Marts

Data Sources

Data Storage

OLAP Engine

Front-End Tools

# 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

9

# EXTRACTION, TRANSFORMATION, AND LOADING (ETL)

- **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 indicies and partitions
- **Refresh**
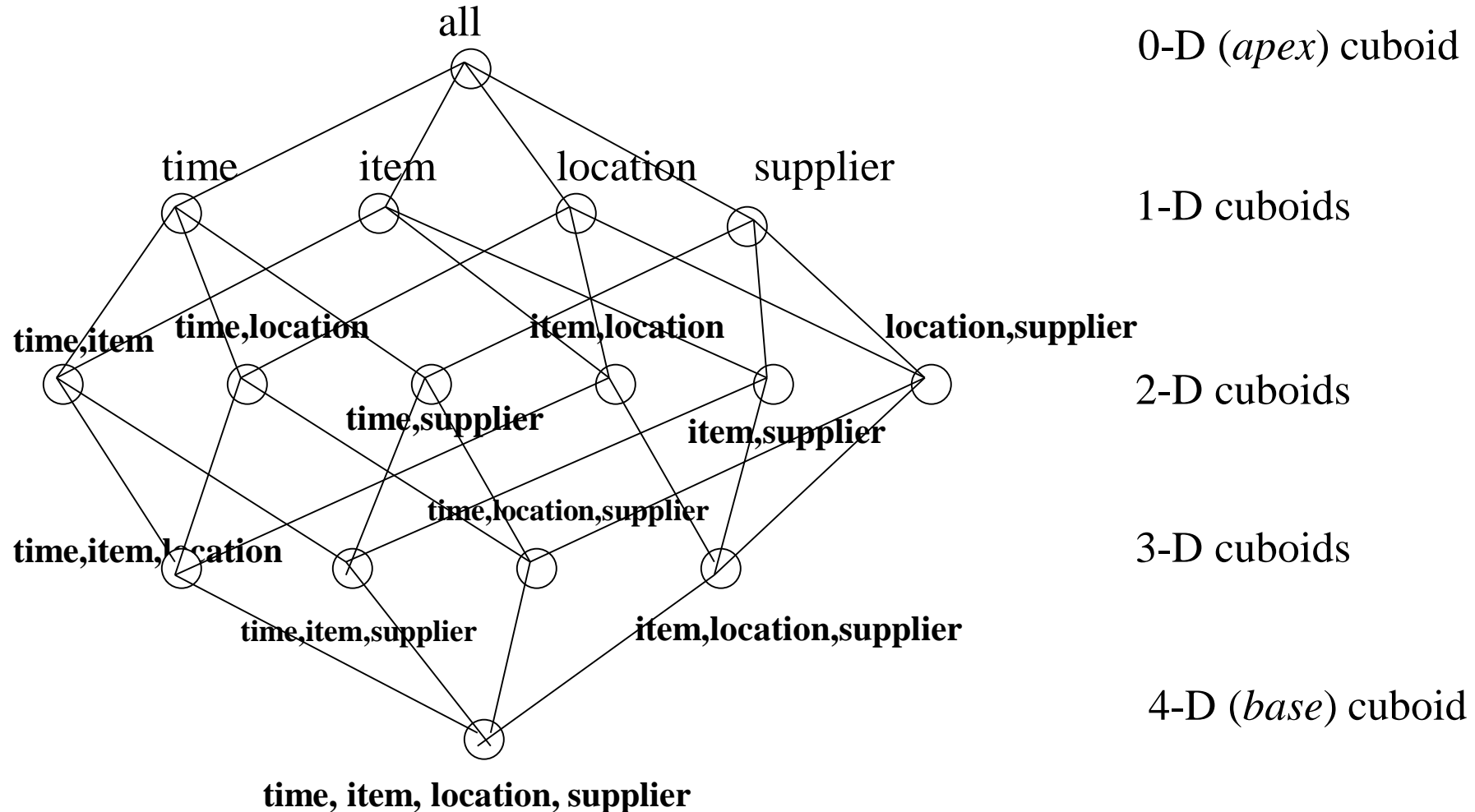  - propagate the updates from the data sources to the warehouse

# METADATA REPOSITORY

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

- Description of the structure of the data warehouse

  - schema, view, dimensions, hierarchies, derived data defn, data mart locations and contents

- Operational meta-data

  - 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)

- The algorithms used for summarization

- The mapping from operational environment to the data warehouse

- Data related to system performance

  - warehouse schema, view and derived data definitions

- Business data

  - business terms and definitions, ownership of data, charging policies

11

- 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.
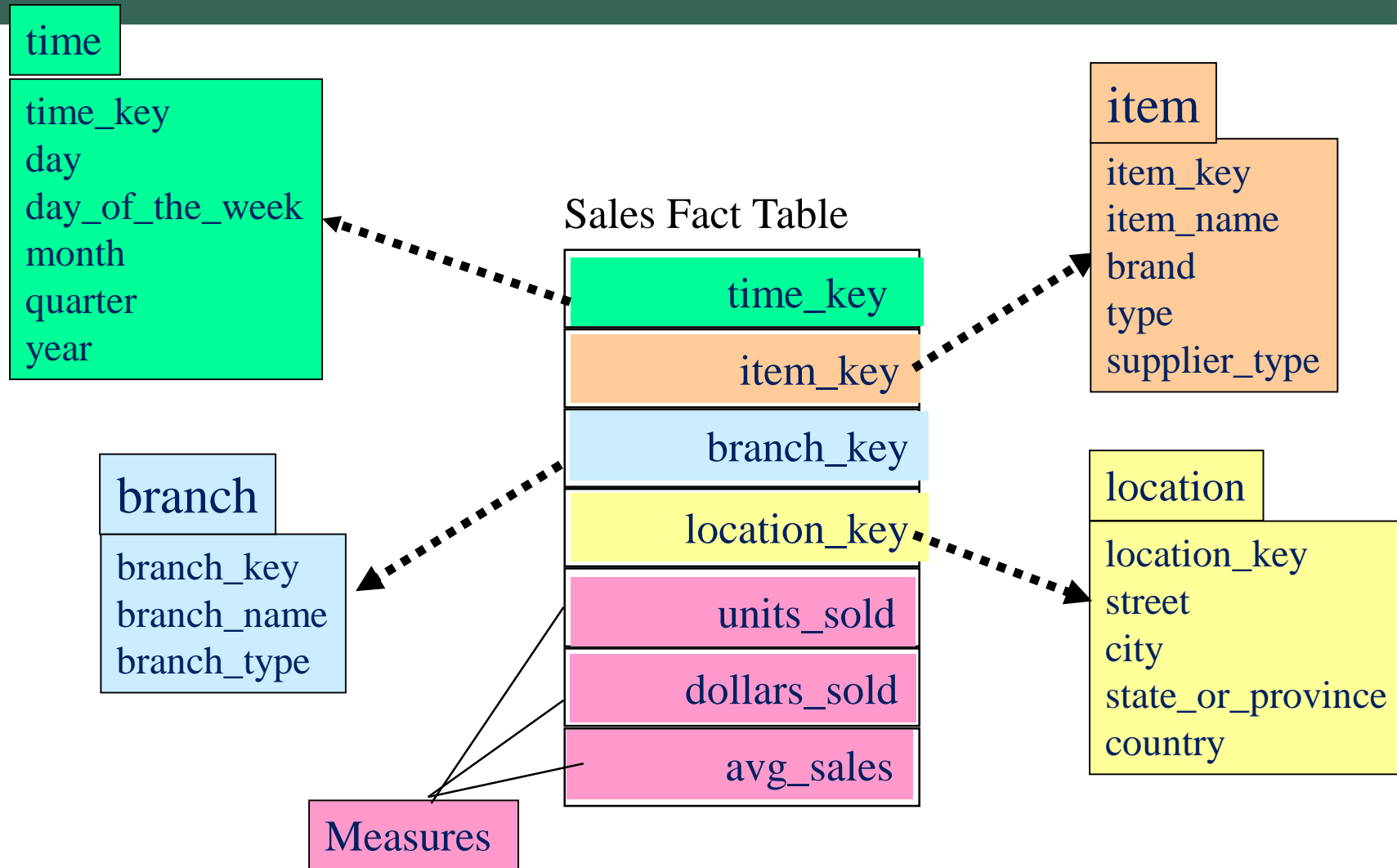
12

# CUBE: A LATTICE OF CUBOIDS



all — 0-D (*apex*) cuboid

time, item, location, supplier — 1-D cuboids

time,item | time,location | item,location | location,supplier | time,supplier | item,supplier — 2-D cuboids

time,item,location | time,location,supplier | time,item,supplier | item,location,supplier — 3-D cuboids

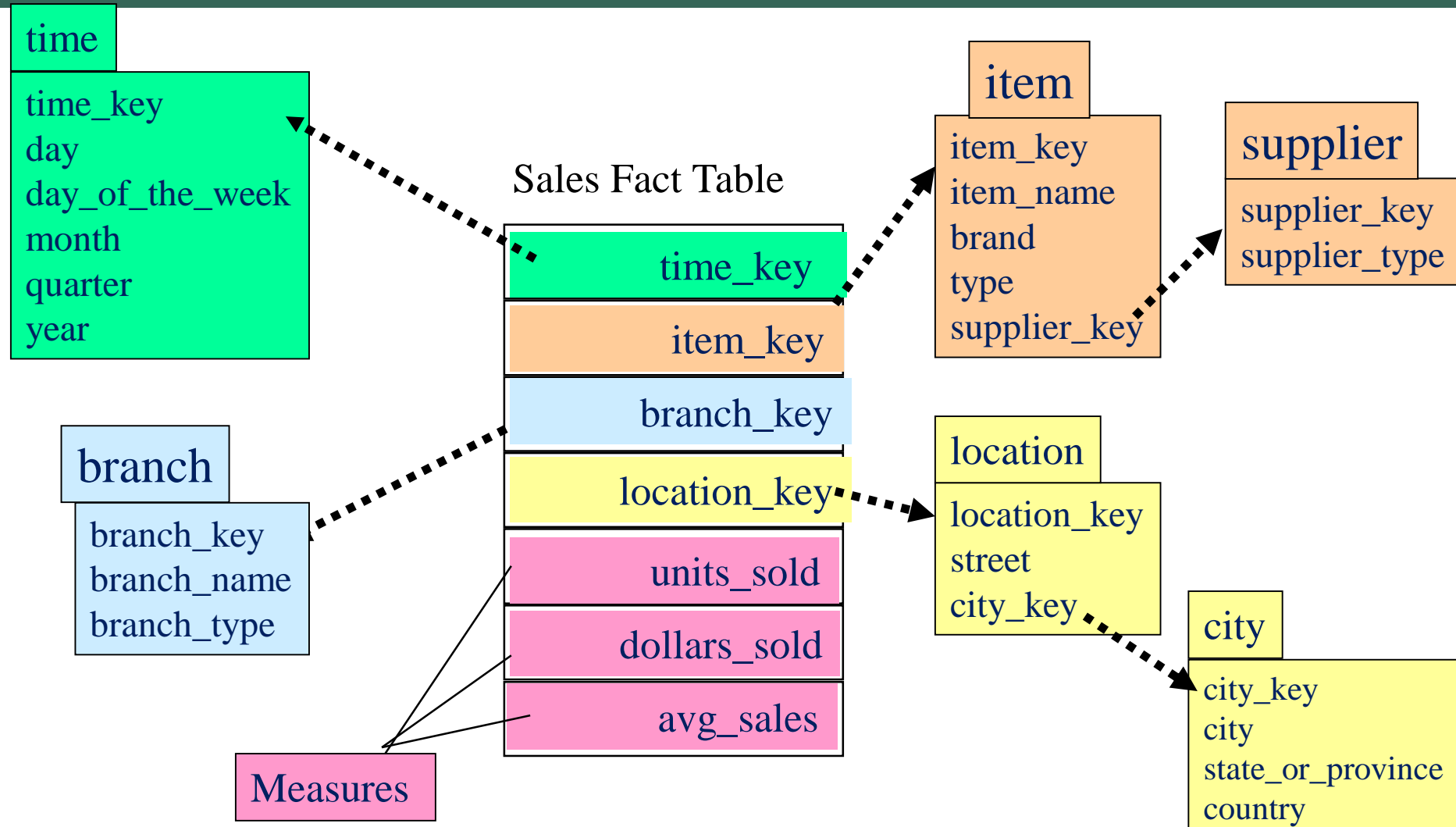time, item, location, supplier — 4-D (*base*) cuboid

13

# 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
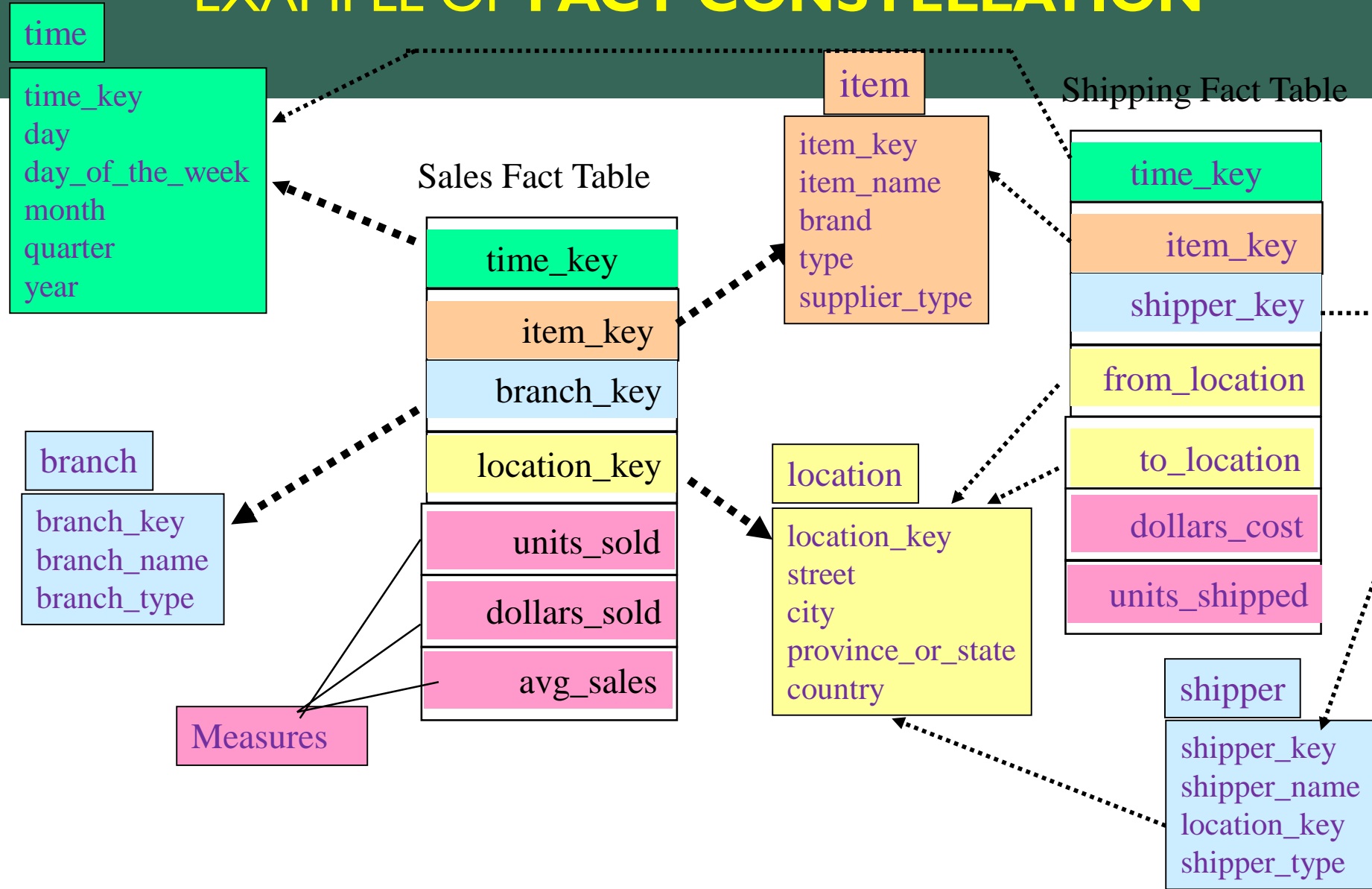
14

# EXAMPLE OF **STAR SCHEMA**

**time**
- time_key
- day
- day_of_the_week
- month
- quarter
- year

**item**
- item_key
- item_name
- brand
- type
- supplier_type

Sales Fact Table

- time_key
- item_key
- branch_key
- location_key
- units_sold
- dollars_sold
- avg_sales

**branch**
- branch_key
- branch_name
- branch_type

**location**
- location_key
- street
- city
- state_or_province
- country

Measures

15

# EXAMPLE OF **SNOWFLAKE SCHEMA**

# EXAMPLE OF **FACT CONSTELLATION**

**time**

time_key
day
day_of_the_week
month
quarter
year

Sales Fact Table

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

**branch**

branch_key
branch_name
branch_type

Measures

**item**

item_key
item_name
brand
type
supplier_type

**location**

location_key
street
city
province_or_state
country

Shipping Fact Table

| time_key |
| item_key |
| shipper_key |
| from_location |
| to_location |
| dollars_cost |
| units_shipped |

**shipper**

shipper_key
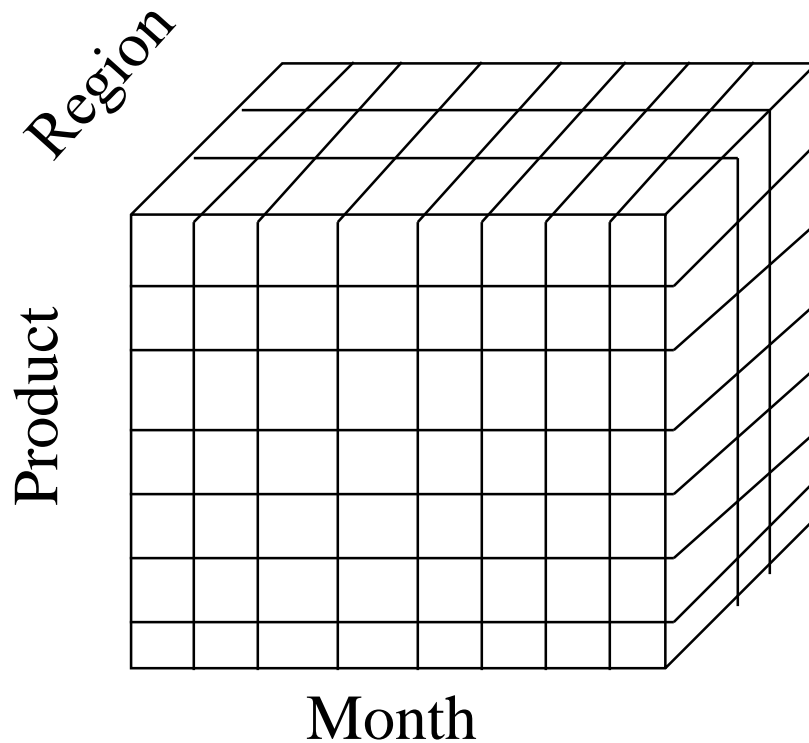shipper_name
location_key
shipper_type
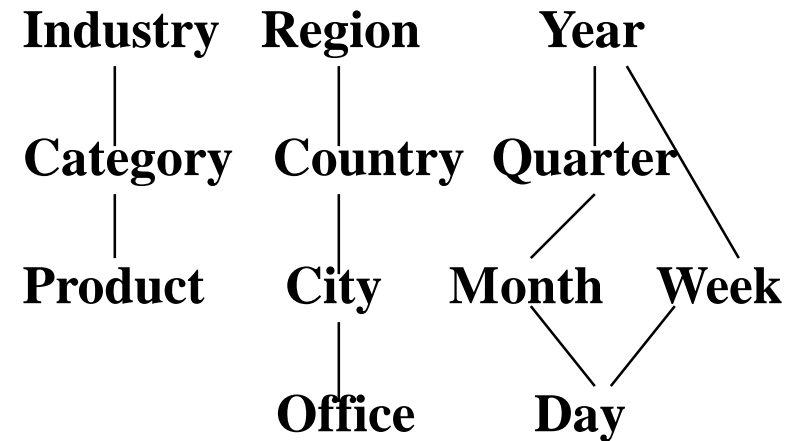
17

# DATA CUBE MEASURES: 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()

# MULTIDIMENSIONAL DATA
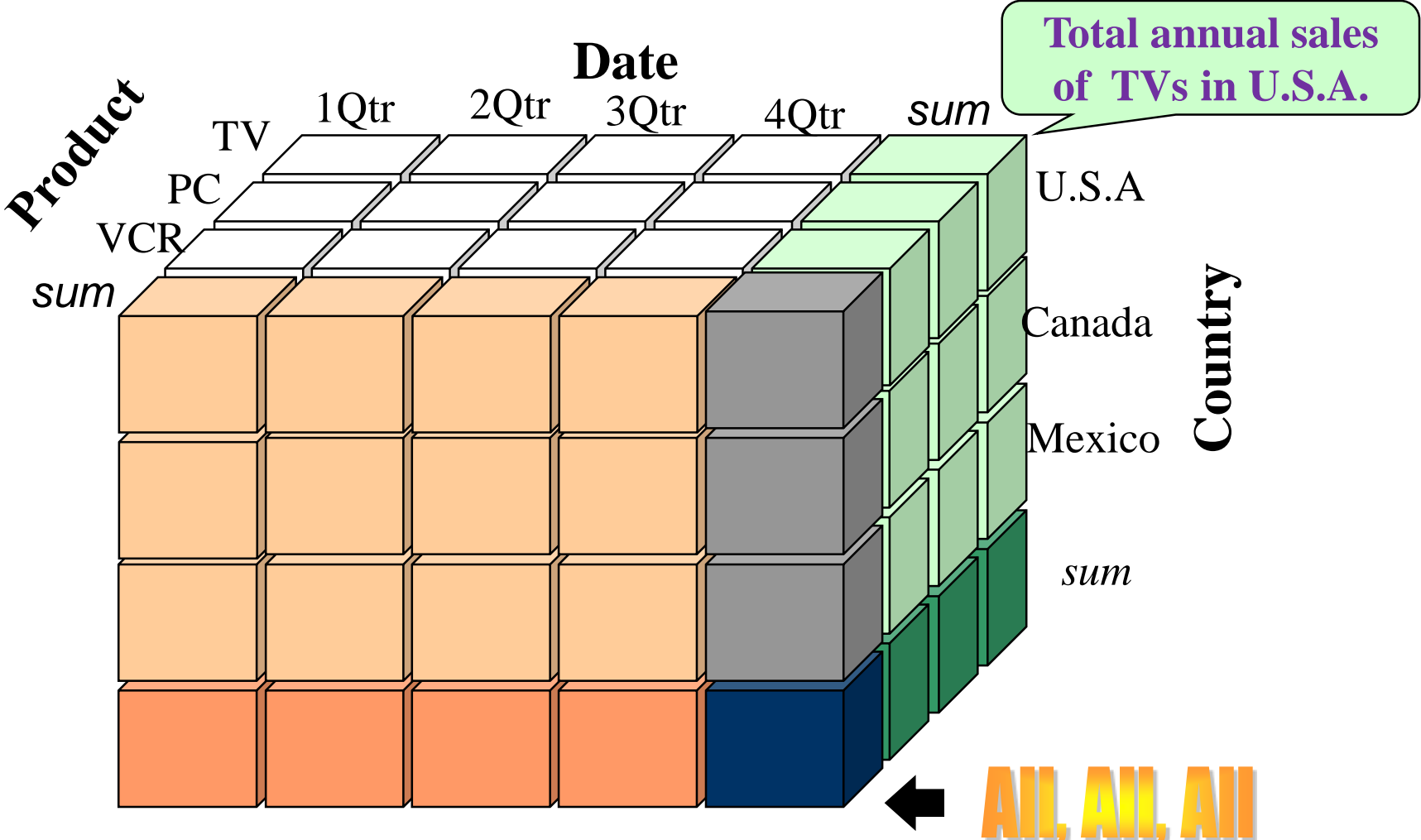
- Sales volume as a function of product, month, and region
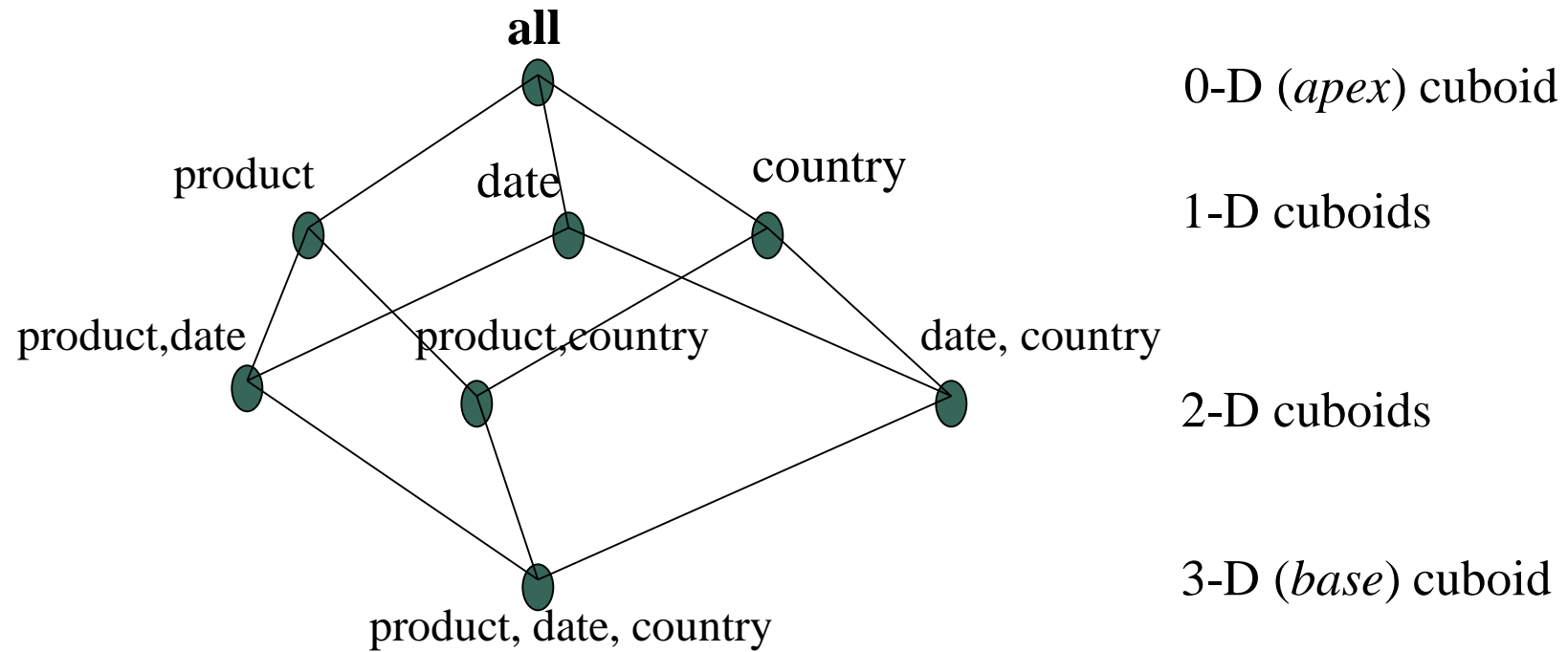
**Dimensions: *Product, Location, Time***
**Hierarchical summarization paths**



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

# CUBOIDS CORRESPONDING TO THE CUBE



**all**

product    date    country

product,date    product,country    date, country

product, date, country

0-D (*apex*) cuboid

1-D cuboids

2-D cuboids

3-D (*base*) cuboid
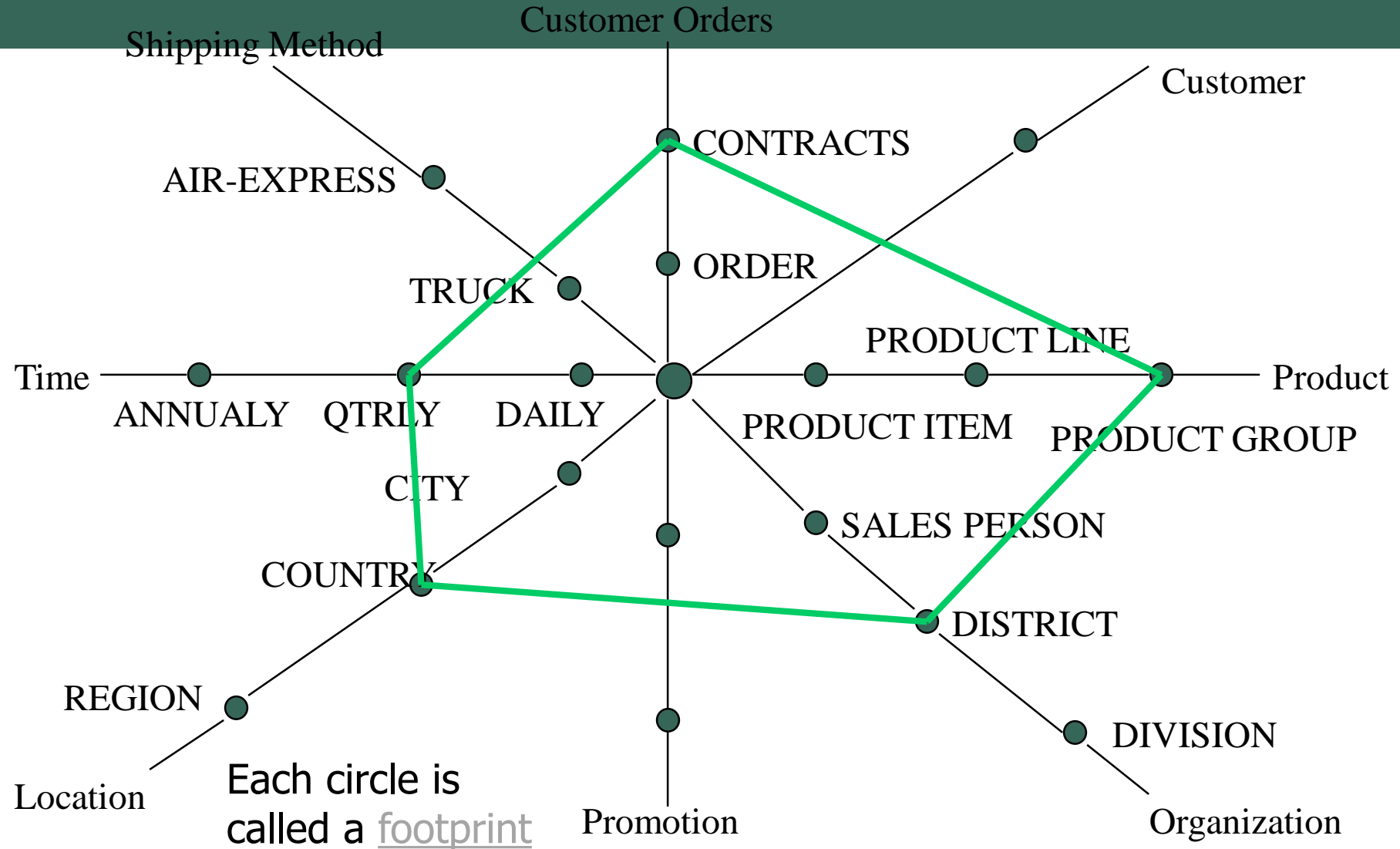
# TYPICAL OLAP OPERATIONS

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

23

# A STAR-NET QUERY MODEL



Customer Orders

Shipping Method

Customer

CONTRACTS

AIR-EXPRESS

ORDER

PRODUCT LINE

TRUCK

Time

Product

ANNUALY   QTRLY   DAILY

PRODUCT ITEM   PRODUCT GROUP

CITY

SALES PERSON

COUNTRY

DISTRICT

REGION

DIVISION

Location

Each circle is called a footprint
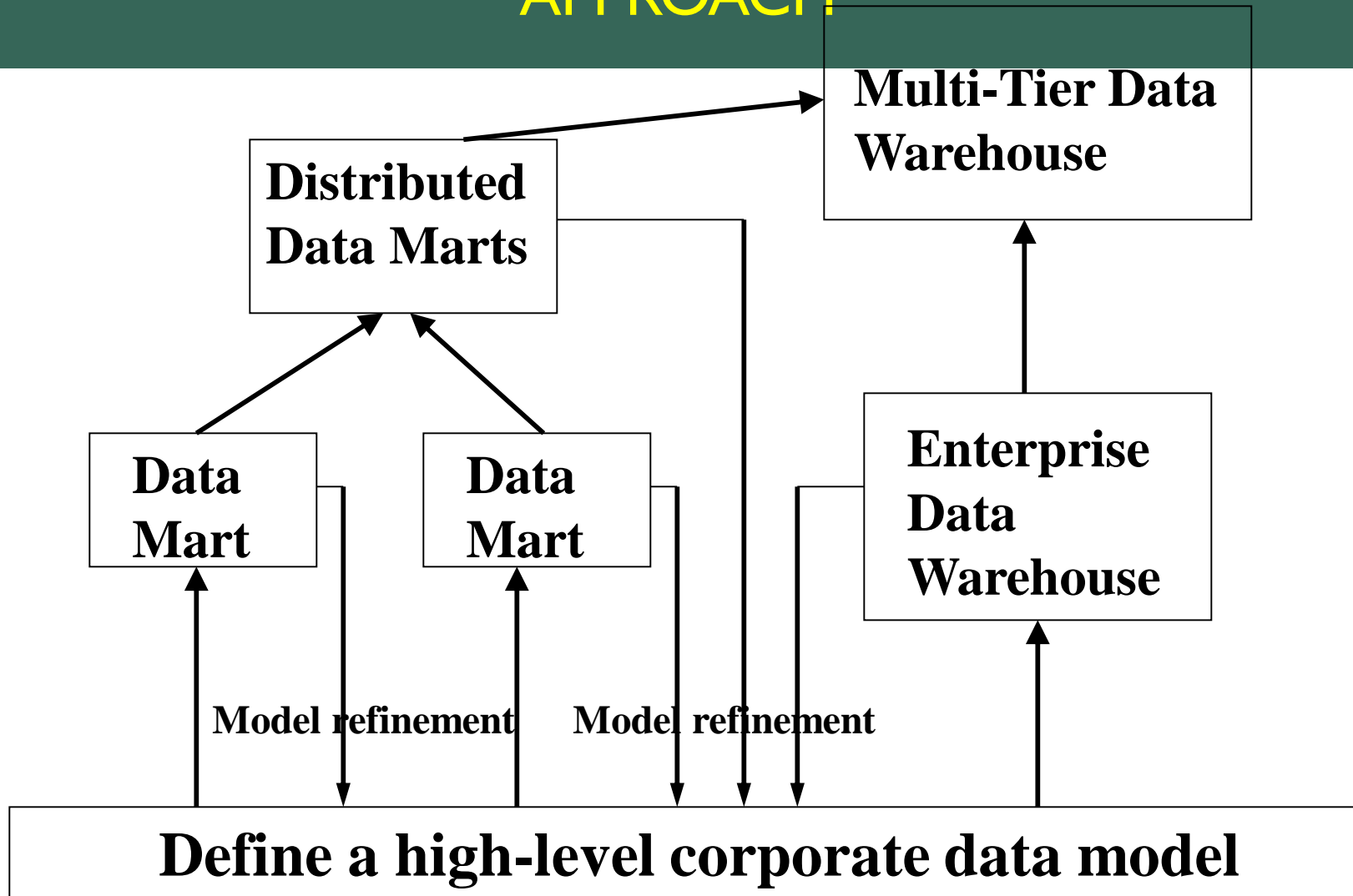
Promotion

Organization

24

# 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

25

# DATA WAREHOUSE DESIGN PROCESS

- **Top-down, bottom-up approaches or a combination** of both
  - <u>Top-down</u>: Starts with overall design and planning (mature)
  - <u>Bottom-up</u>: Starts with experiments and prototypes (rapid)
- **From software engineering point of view**
  - <u>Waterfall</u>: structured and systematic analysis at each step before proceeding to the next
  - <u>Spiral</u>: 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

26

# 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

# 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?

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

- Materialization of data cube

  - Materialize <u>every</u> (cuboid) (**full materialization**), <u>none </u>(**no materialization**), or <u>some (**partial materialization)**</u>

  - Selection of which cuboids to materialize

    - Based on size, sharing, access frequency, etc.
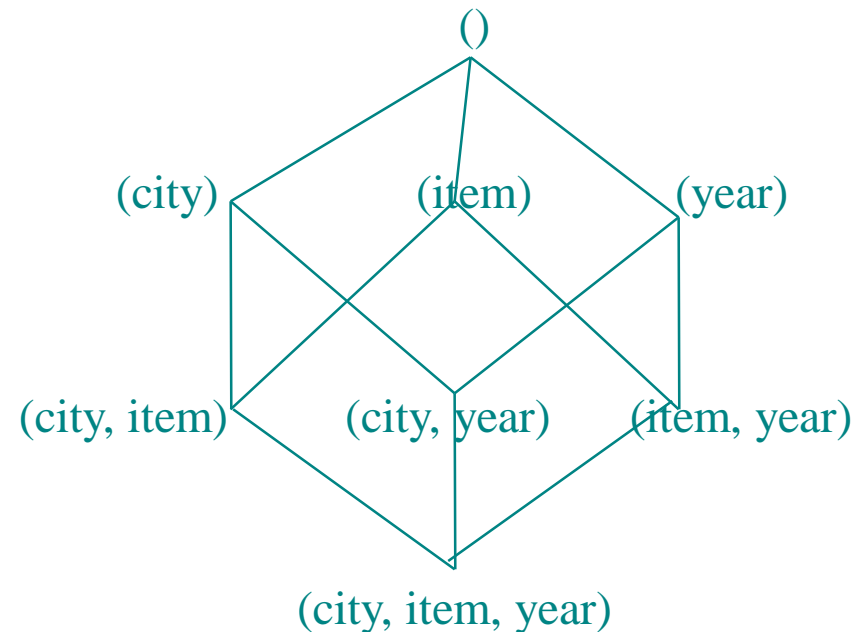
# THE "COMPUTE CUBE" OPERATOR

- 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)*

  *()*



()

(city)  (item)  (year)

(city, item) (city, year) (item, year)

(city, item, year)

30

# 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

- A recent bit compression technique, Word-Aligned Hybrid (WAH), makes it work for high cardinality domain as well [Wu, et al. TODS'06]

## 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 |

31

# 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

32