

A · P · U
ASIA PACIFIC UNIVERSITY
OF TECHNOLOGY & INNOVATION

Data Management

CT051-3-M

Topic 6 – Data Warehouse

Learning Outcomes

By the end of this lecture, YOU should be able to:

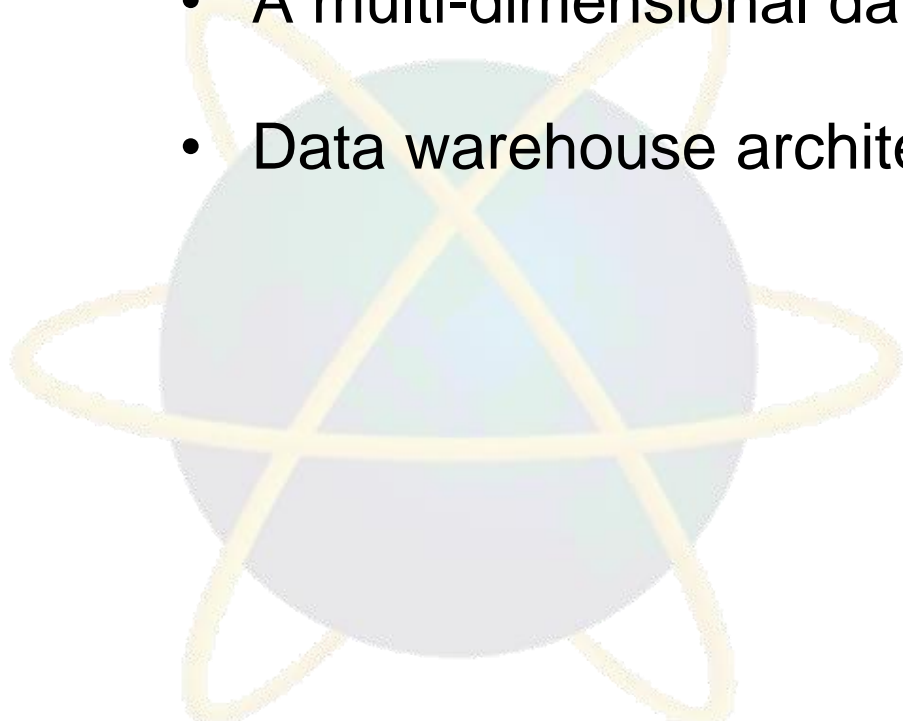
- Explain data warehouse and OLAP concepts
- Explain the data warehouse architecture and schemes
- Introduce the concept of OLAP

Key Terms you must be able to use

- If you have mastered this topic, you should be able to use the following terms correctly in your assignments and exams:
 - OLTP
 - OLAP
 - Data Cube
 - Star Scheme
 - Snowflake schema
 - Fact constellations
 - Information processing
 - Analytical processing

Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture



What is Data Warehouse?

- Defined in many different ways, but not strictly.
 - 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, sale**.
- 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—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*.

How are organizations using the information from data warehouses ?

- Many organizations use this information to support business decision making activities:
- Increasing customer focus, which includes the analysis of customer buying patterns (such as buying preference, buying time, budget cycles, and appetites for spending).
- Repositioning products and managing product portfolios by comparing the performance of sales by quarter, by year, and by geographic regions in order to fine-tune production strategies.

OLTP vs. OLAP

	OLTP	OLAP
users	clerk, IT professional	knowledge worker
function	day to day operations	decision support
DB design	ER diagram, application-oriented	Star/snowflake, subject-oriented
data	current, up-to-date detailed, flat relational isolated	historical, summarized, multidimensional integrated, consolidated
usage	repetitive	ad-hoc
unit of work	short, simple transaction	complex query
# records accessed	tens	millions
#users	thousands	hundreds
DB size	100MB-GB	100GB-TB

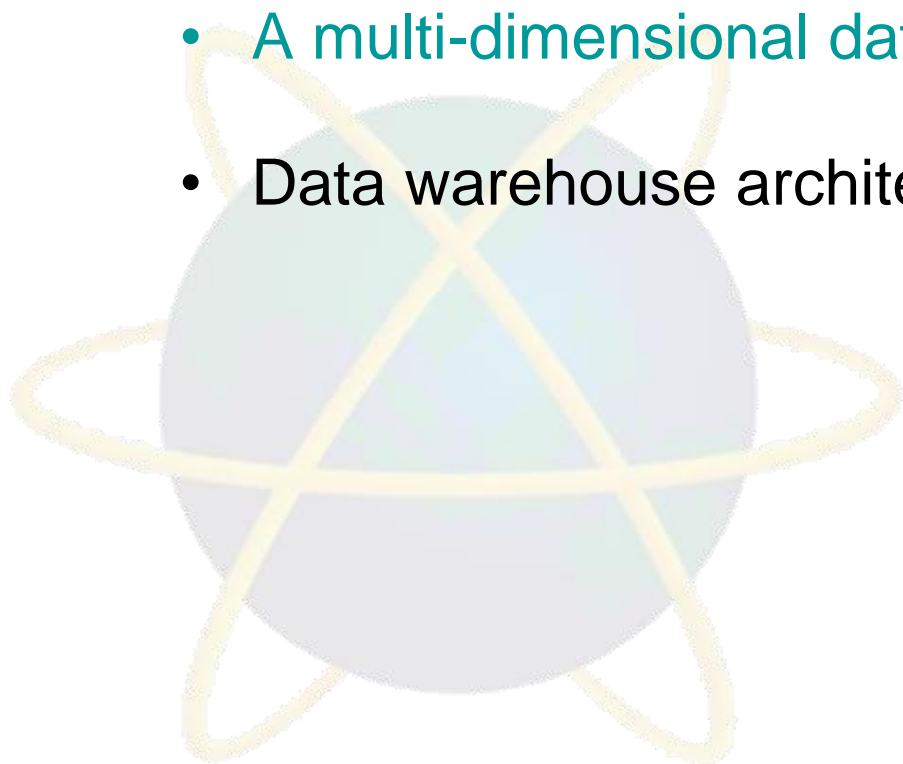


Why Separate Data Warehouse?

- 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.
- 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

Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture



A 2-D view of sales data for AllElectronics

location = "Vancouver"

item (type)

time (quarter)	<i>home entertainment</i>	<i>computer</i>	<i>phone</i>	<i>security</i>
Q1	605	825	14	400
Q2	680	952	31	512
Q3	812	1023	30	501
Q4	927	1038	38	580

3-D view of sales data for AllElectronics



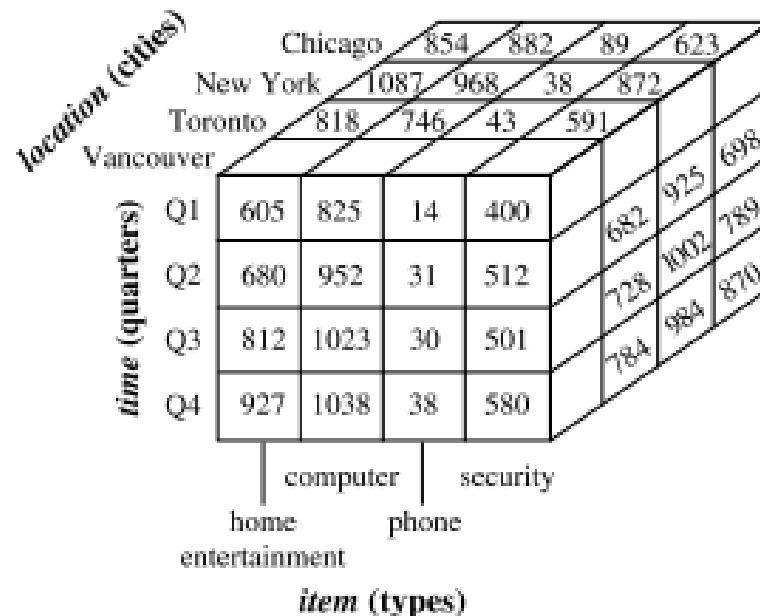
location = "Chicago"

location = "New York"

location = "Toronto"

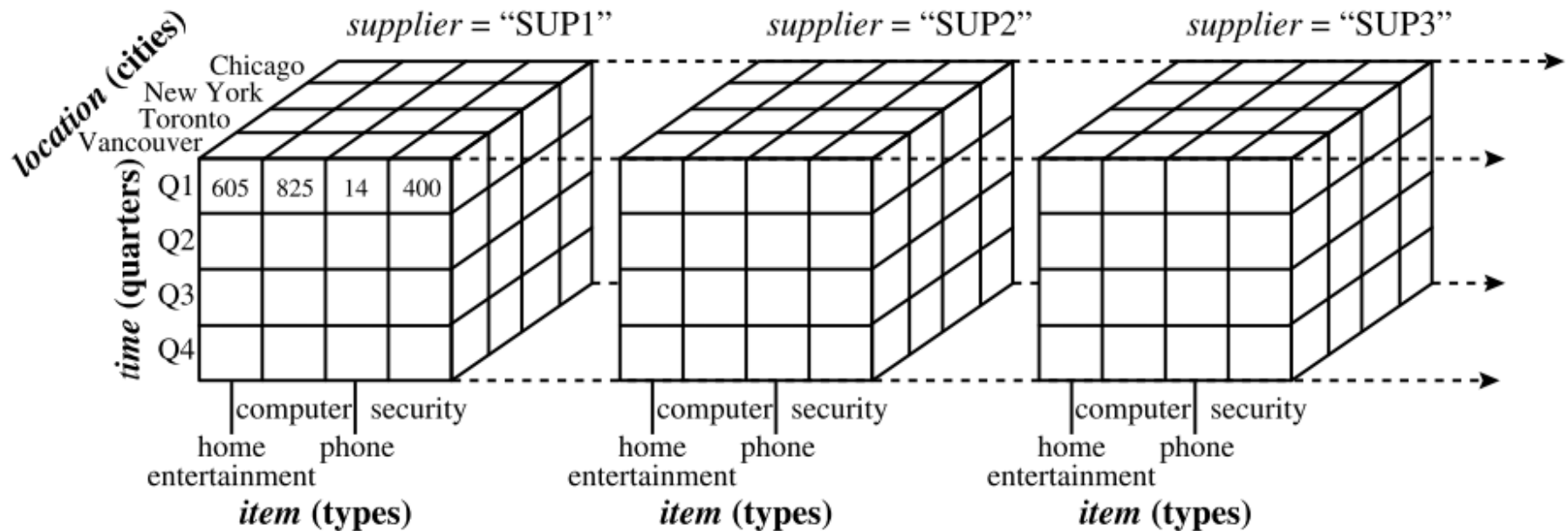
location = "Vancouver"

item					item					item					item				
time	home				home				home				home						
	ent.	comp.	phone	sec.	ent.	comp.	phone	sec.	ent.	comp.	phone	sec.	ent.	comp.	phone	sec.			
Q1	854	882	89	623	1087	968	38	872	818	746	43	591	605	825	14	400			
Q2	943	890	64	698	1130	1024	41	925	894	769	52	682	680	952	31	512			
Q3	1032	924	59	789	1034	1048	45	1002	940	795	58	728	812	1023	30	501			
Q4	1129	992	63	870	1142	1091	54	984	978	864	59	784	927	1038	38	580			





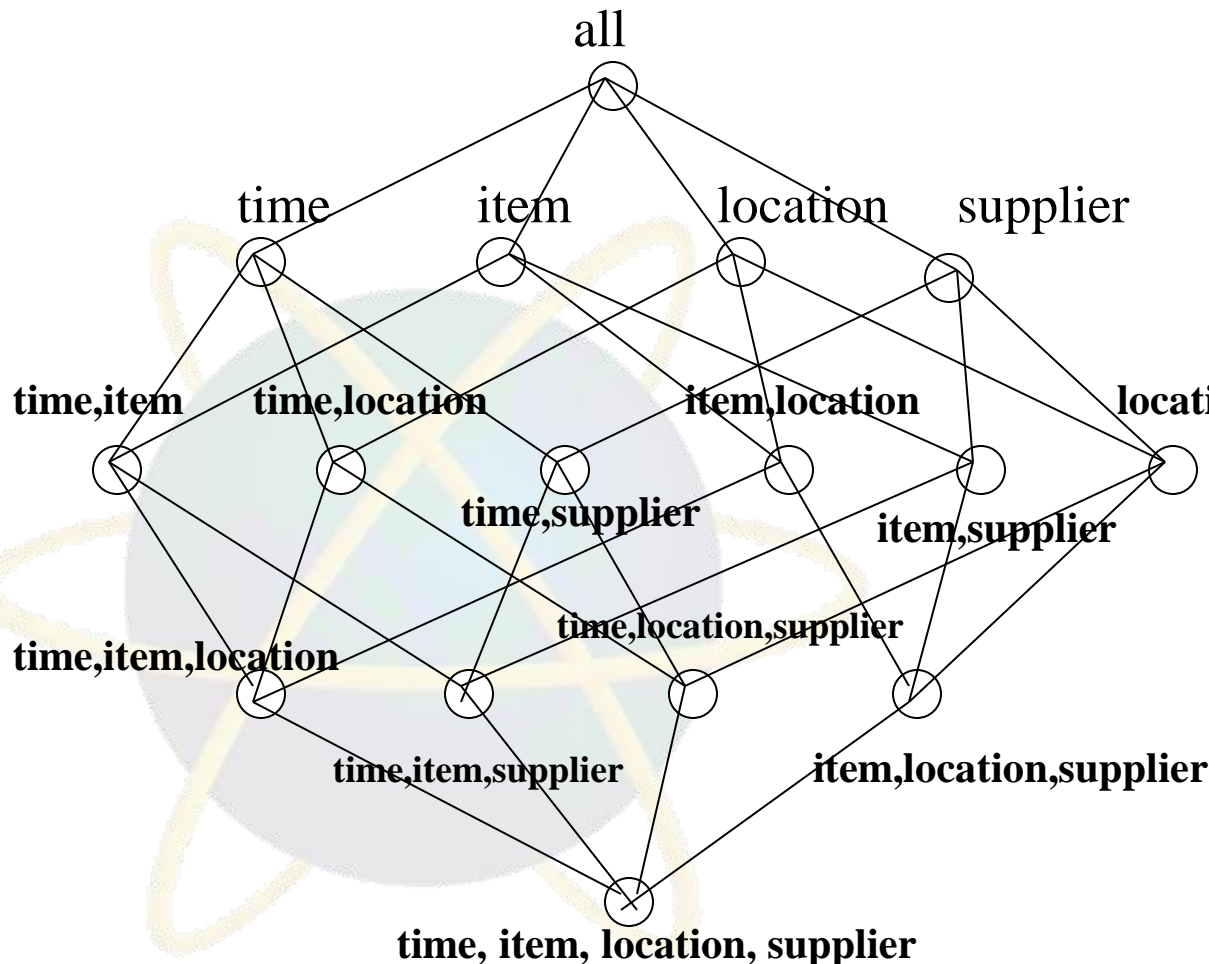
4 – D View of sales data



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**.

Cube: A Lattice of Cuboids



0-D(apex) cuboid

1-D cuboids

2-D cuboids

3-D cuboids

4-D(base) cuboid

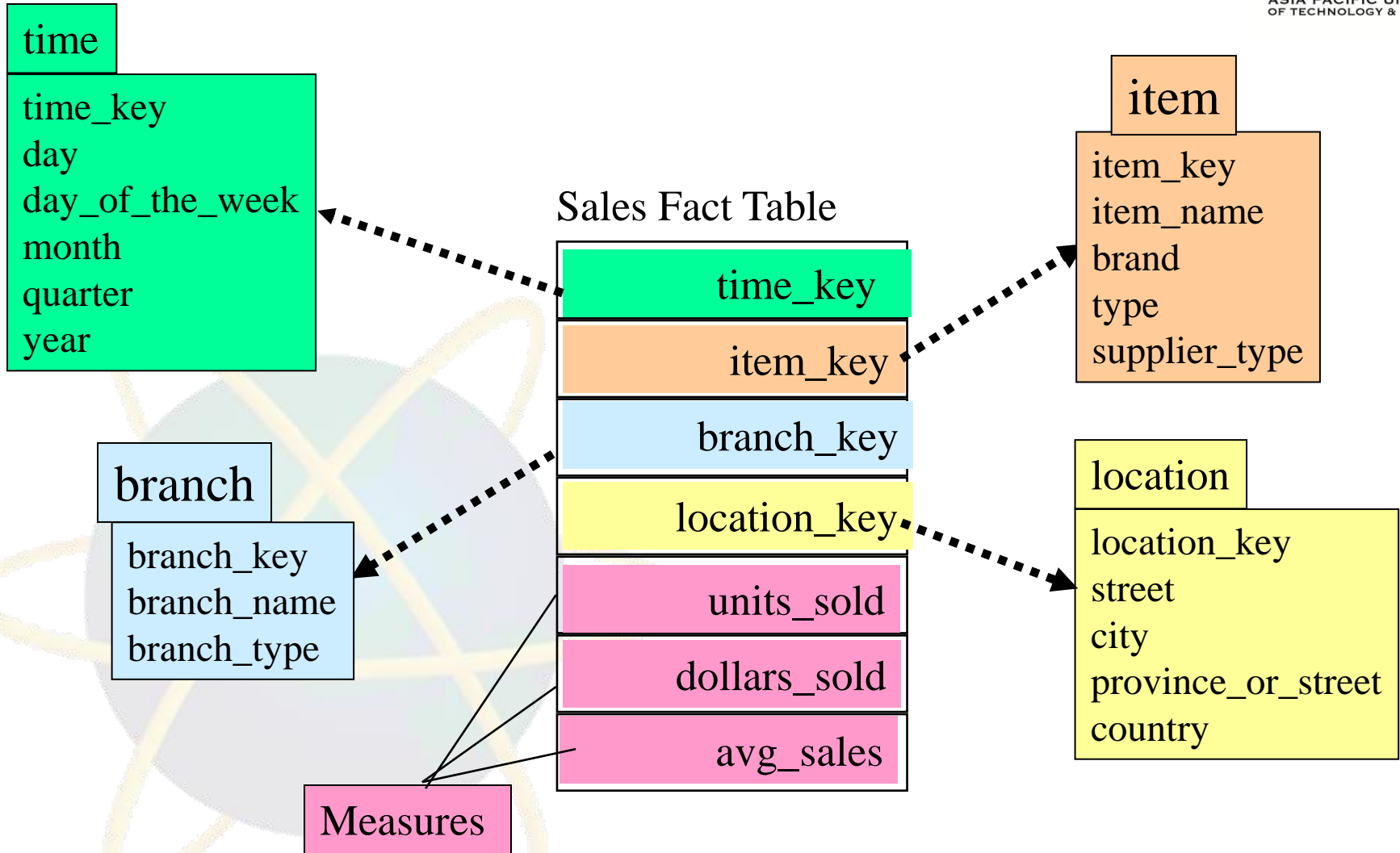
Conceptual Modeling of Data Warehouses

- 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

Star schema

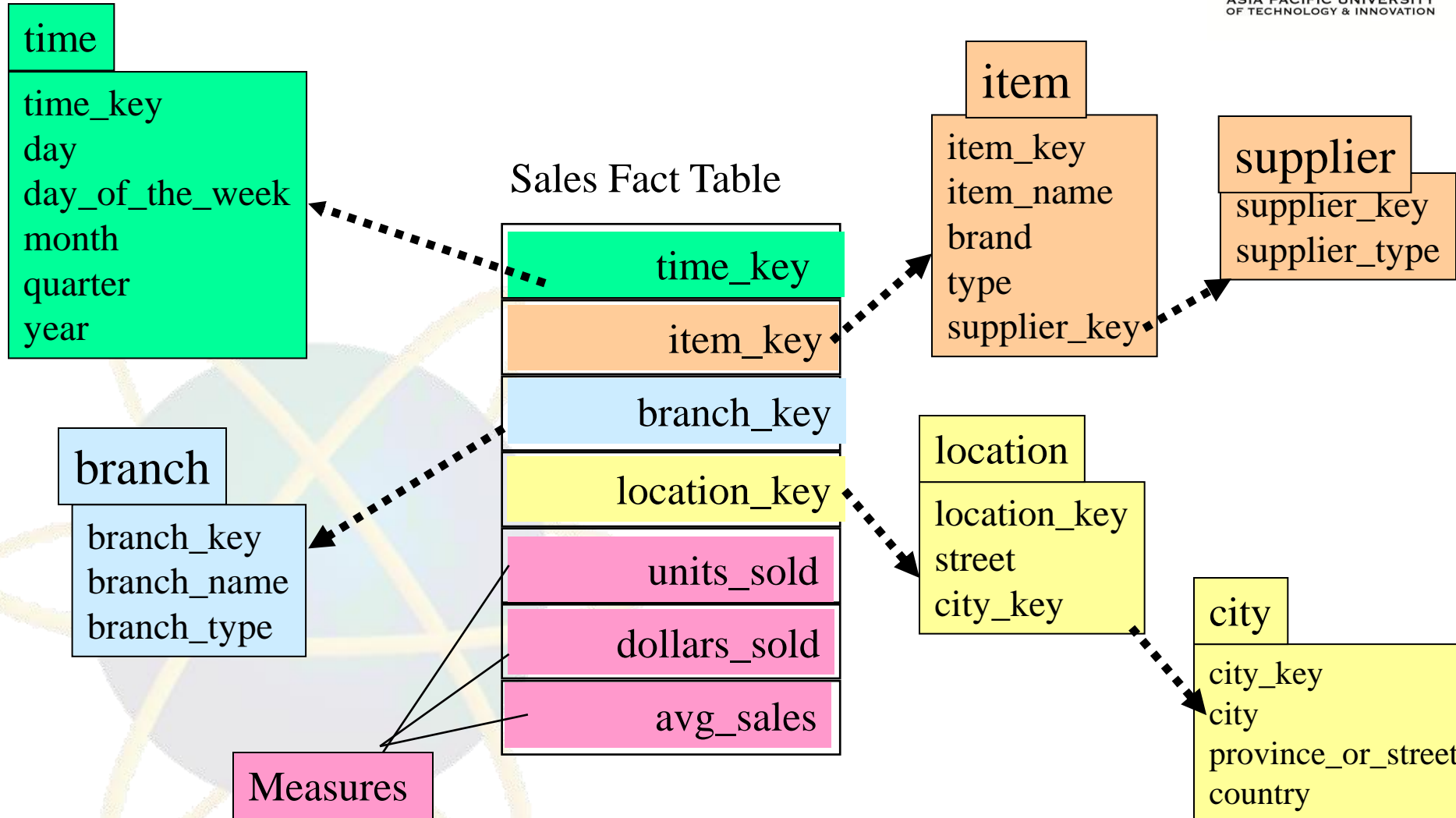
- The most common modeling paradigm is the star schema, in which the data warehouse contains :
 - (1) a large central table (fact table) containing the bulk of the data, with no redundancy, and
 - (2) a set of smaller attendant tables (dimension tables), one for each dimension. The schema graph resembles a starburst, with the dimension tables displayed in a radial pattern around the central fact table.

Example of Star Schema





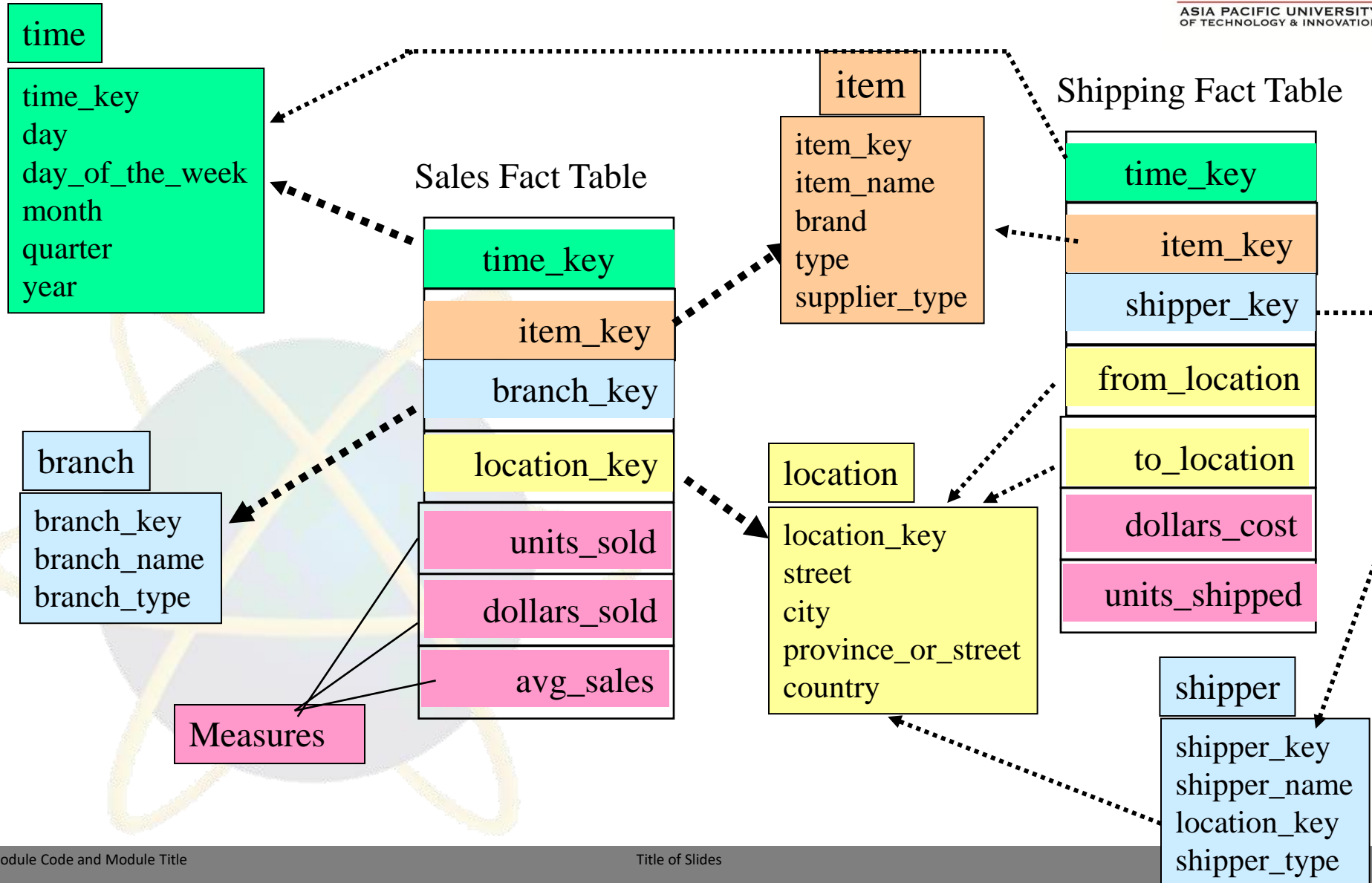
Example of Snowflake Schema



Example of Fact Constellation



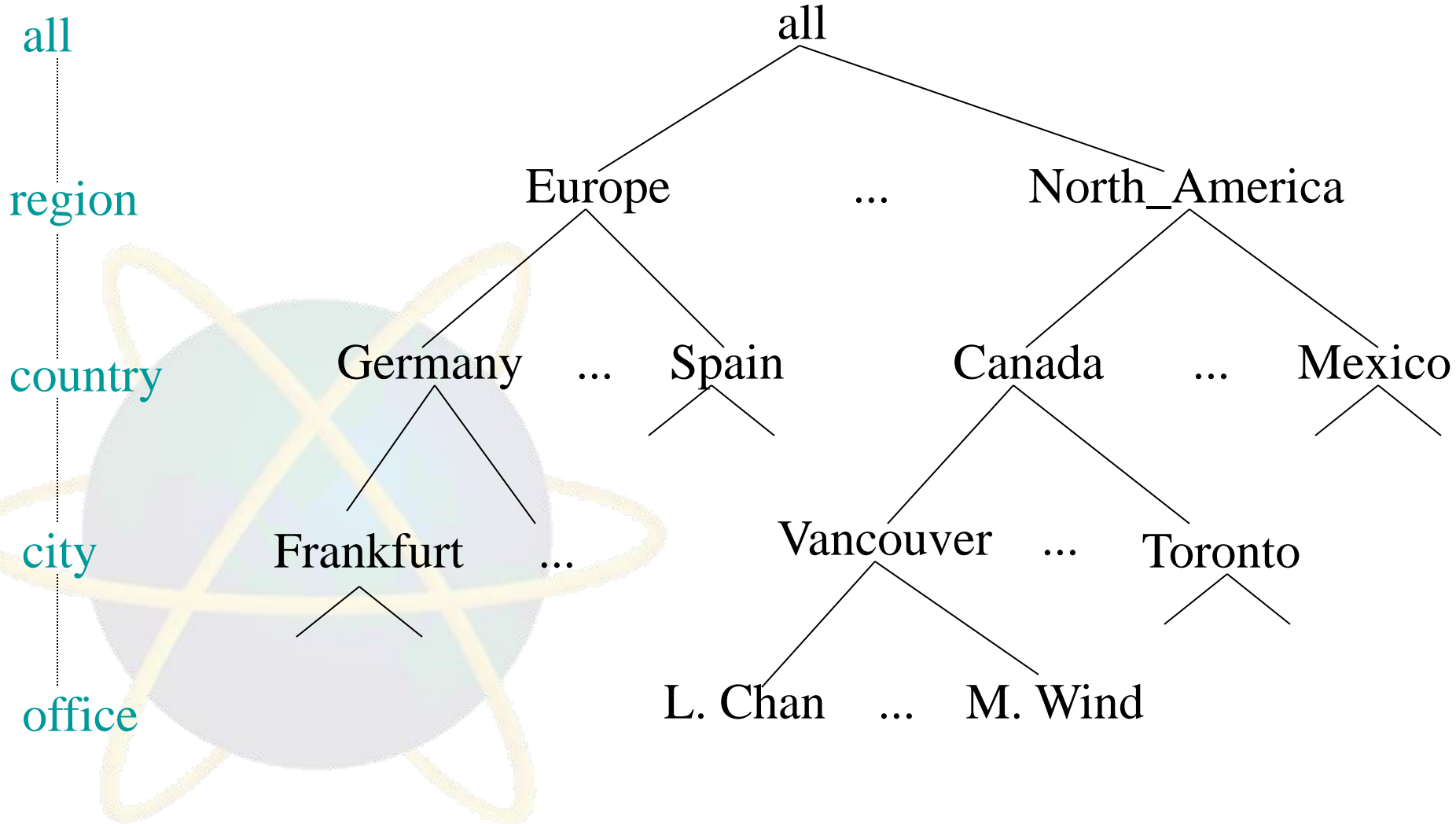
A · P · U
ASIA PACIFIC UNIVERSITY
OF TECHNOLOGY & INNOVATION



A Concept Hierarchy: Dimension (location)



A · P · U
ASIA PACIFIC UNIVERSITY
OF TECHNOLOGY & INNOVATION



View of Warehouses and Hierarchies



A · P · U
ASIA PACIFIC UNIVERSITY
OF TECHNOLOGY & INNOVATION

dbminer

File Edit Query View Window Help

WareHouse Dimensions

DemoWH

- SCHEMAS
 - MasterDemoDB.dbo.SalesD
 - COLUMNS
 - DIMENSIONS
 - Product
 - Region
 - revenue
 - cost
 - profit
 - order_qty
 - MEASUREMENTS
 - CUBES
 - SalesData_Cube
 - Small_Cube
 - DMQLs
 - stockdata.dbo.stock
 - COLUMNS
 - DIMENSIONS
 - date
 - price
 - price1
 - MEASUREMENTS
 - CUBES
 - DMQLs

Level Name

- region
- country
- branch_r
- rep_name

WareHouse Dimensions

ANY

- Europe
 - Belgium
 - France
 - Germany
 - Essen
 - Frankfurt
 - Spain
 - Sweden
 - United Kingdom
- Far East
- North America
 - Canada
 - Montreal
 - Toronto
 - Vancouver
 - Charles Loo Nam
 - Hari Krain
 - Kaley Gregson
 - Lee Chan
 - Malcom Young
 - Marthe Whiteduck
 - Torey Wandiko
 - Mexico
 - United States

Description

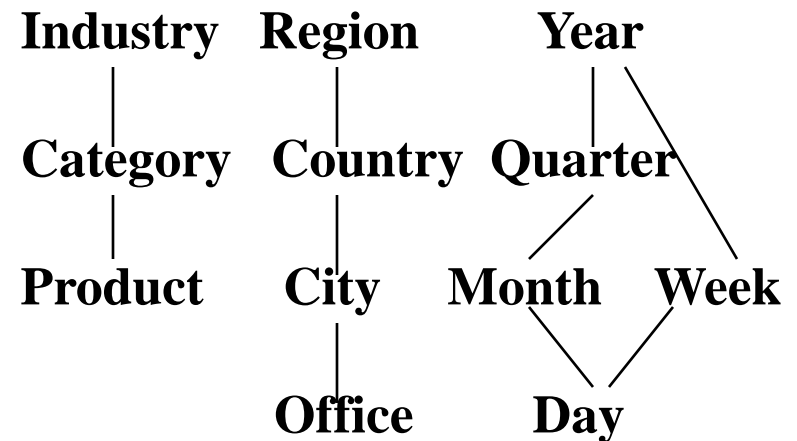
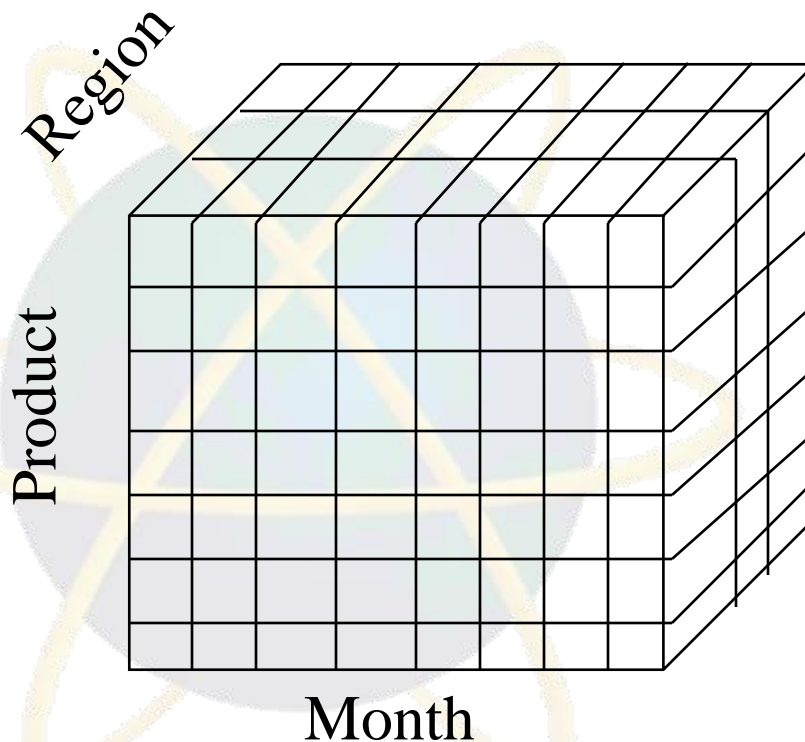
For Help, press F1

NUM

Multidimensional Data

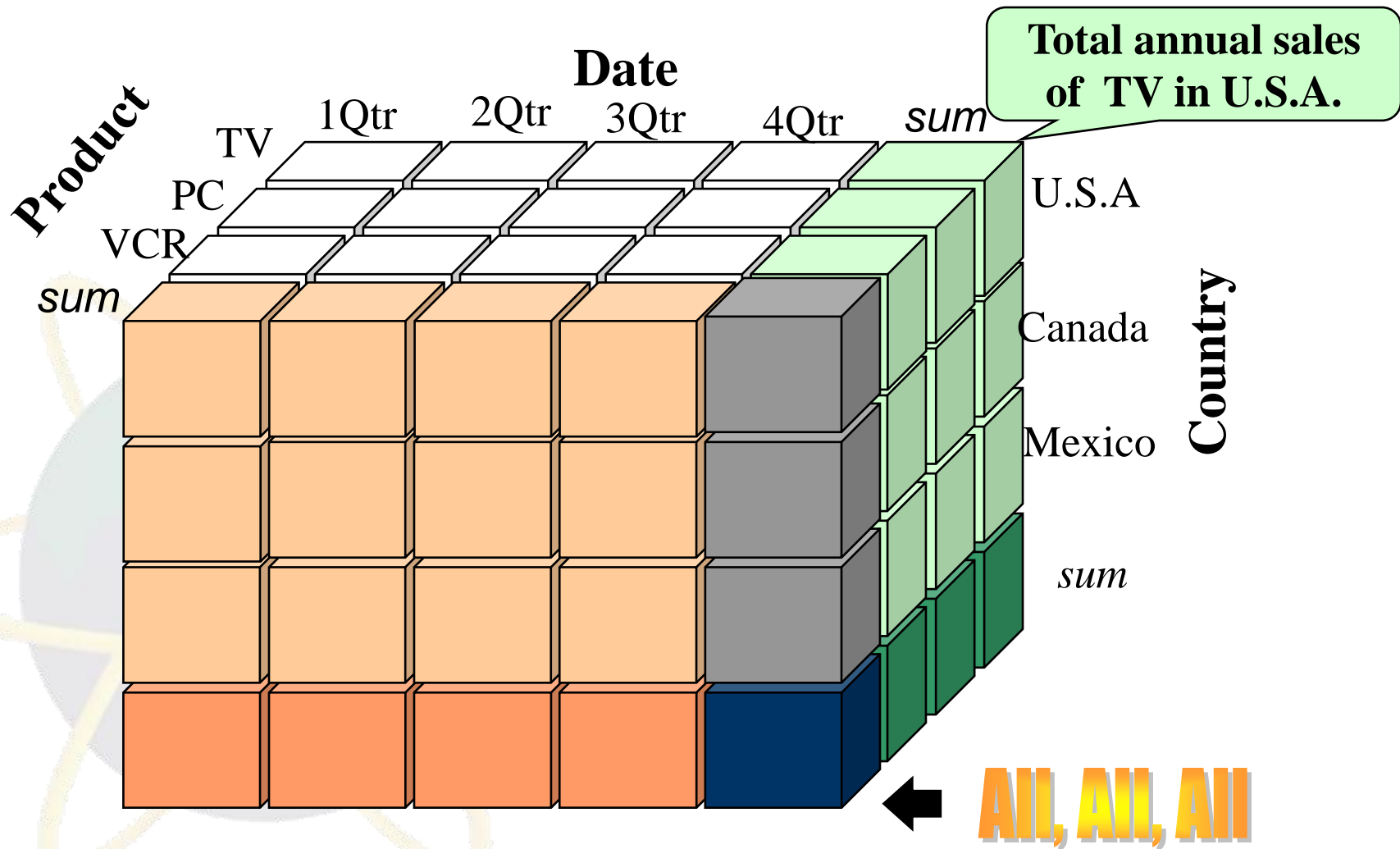
- Sales volume as a function of product, month, and region

Dimensions: Product, Location, Time
Hierarchical summarization paths

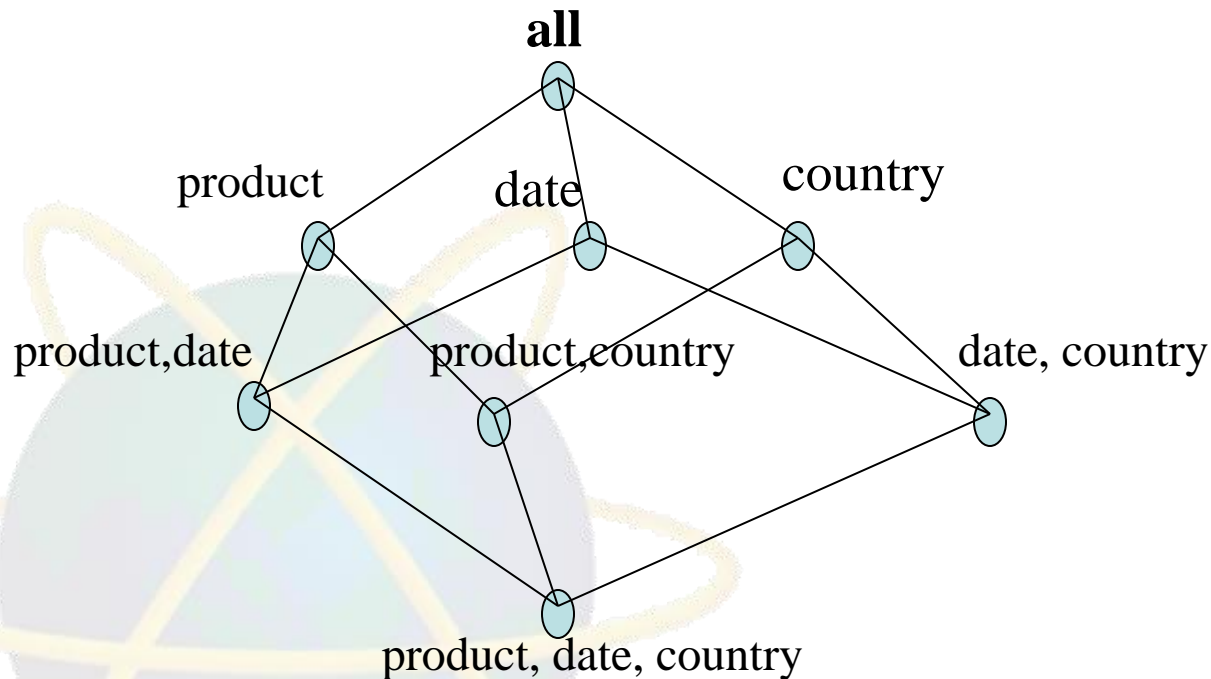




A Sample Data Cube



Cuboids Corresponding to the Cube



0-D(apex) cuboid

1-D cuboids

2-D cuboids

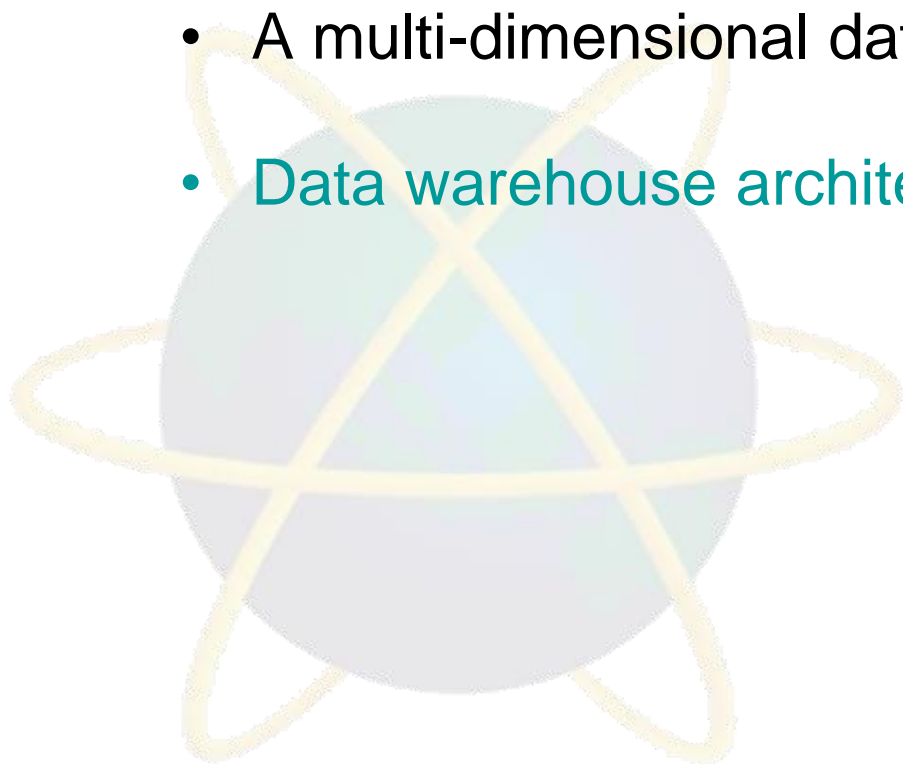
3-D(base) cuboid

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*

Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture

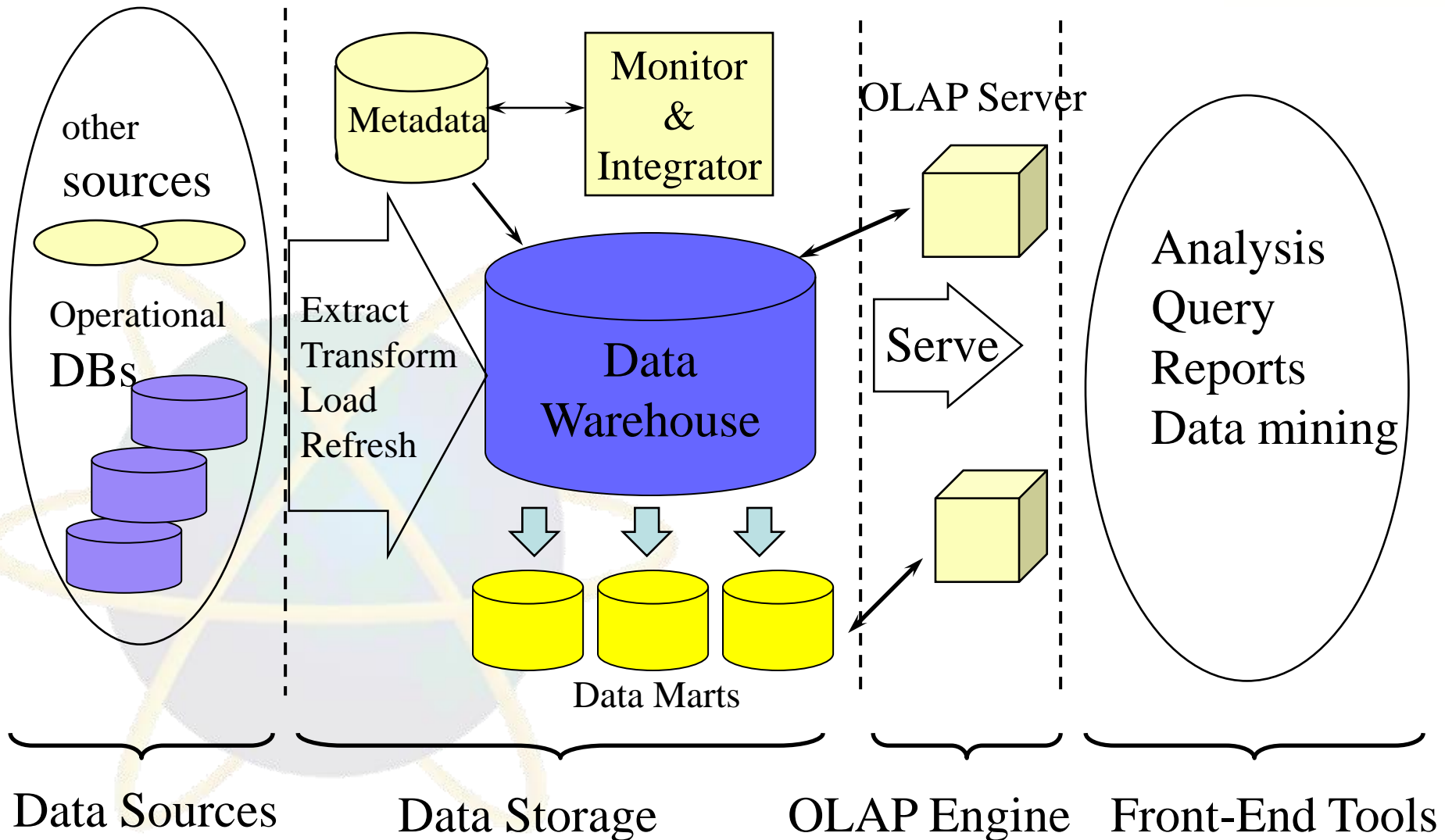


Data Warehouse Design Process

- 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

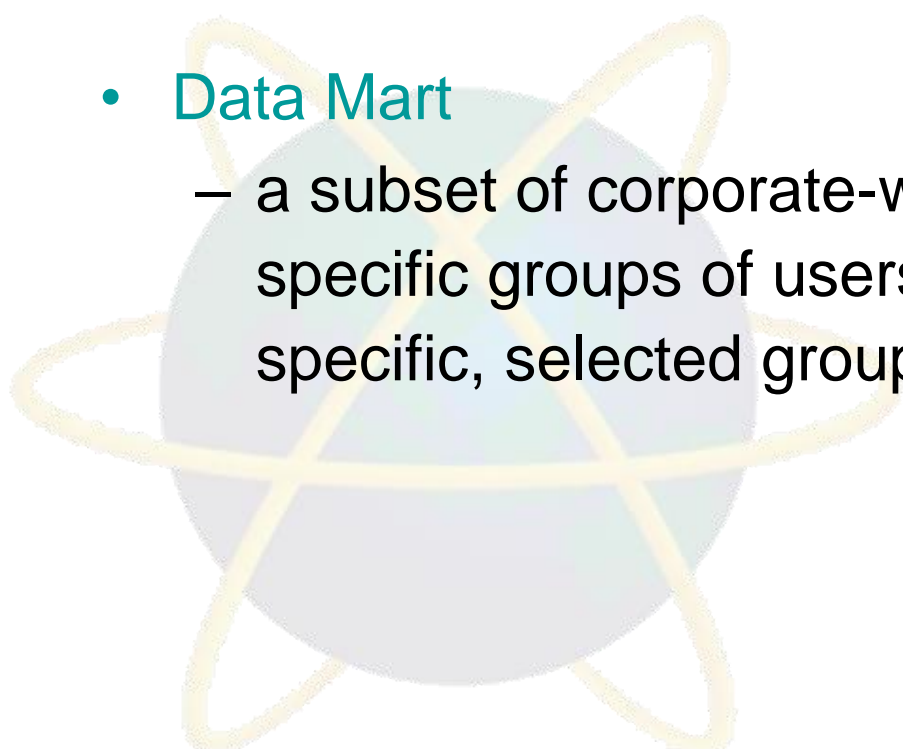


Multi-Tiered Architecture

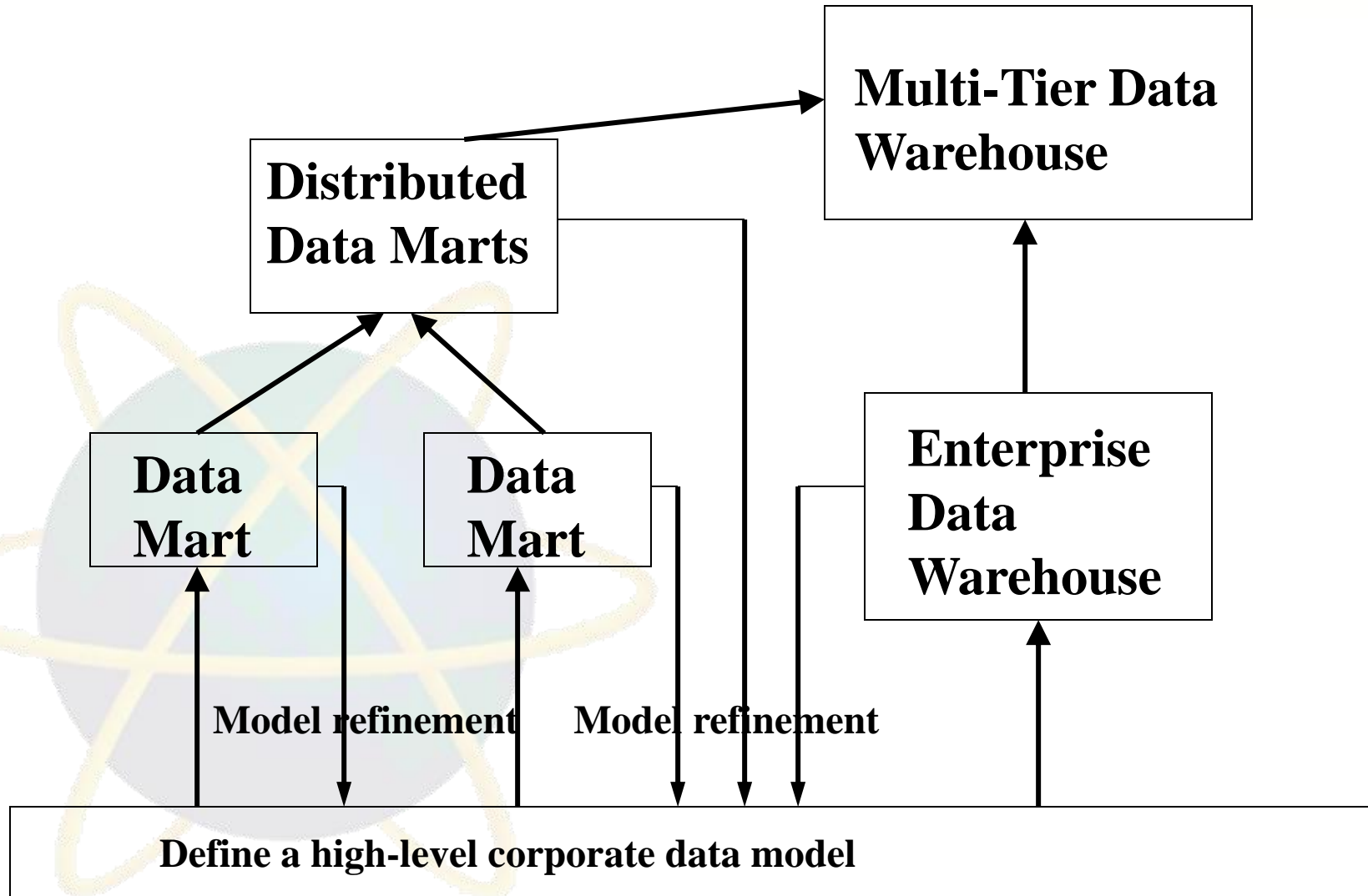


Two Data Warehouse Models

- **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



Data Warehouse Development: A Recommended Approach



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.

Q & A

