



University of Colorado  
Boulder

# CSCI 4502/5502

## Data Mining

---

Fall 2020  
Lecture 06 (Sep 10)

# Reminders

---

- ◆ Homework I
- ◆ due at 9:30am, Th, Sep 10
- ◆ SUBMIT your attempt in Canvas before deadline
- ◆ Computing and Software Career & Internship Fair
- ◆ 11am-4pm, Tu, Sep 15, virtual on Handshake



# Announcement 2

---

- ◆ Homework 2
- ◆ posted in Canvas, due at 9:30am, Th, Sep 17
- ◆ HW2 for both CSCI 4502 and CSCI 5502
- ◆ Jupyter Notebook for HW2
- ◆ SUBMIT your attempt in Canvas before deadline
- ◆ check syllabus for office hours



# Review: Chap 3: Data Preprocessing

---

- ◆ Data preprocessing overview
- ◆ data quality
- ◆ major tasks in data preprocessing
- ◆ Data cleaning
- ◆ Data integration
- ◆ Data reduction
- ◆ Data transformation and discretization





University of Colorado  
Boulder

# Chapter 4: Data Warehouse and OLAP

---

# Chapter 5: Data Cube Technology

# What is A Data Warehouse?

---

- ◆ A decision support database that is maintained **separately** from an 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



# Subject-Oriented

---

- ◆ Organized around **major subjects**
  - ◆ e.g., customers, product, sales; students, courses, departments
- ◆ Focus on the modeling & analysis of data for **decision making**, not on daily operations or transaction processing
- ◆ Provide a simple & concise view around particular subject issues by **excluding data that are not useful in the decision support process**



# Integrated

---

- ◆ Integrate multiple, heterogeneous data sources
  - ◆ relational database, flat files, on-line transaction records
- ◆ Data cleaning and data integration techniques applied
  - ◆ ensure consistency in naming, encoding, attribute measures, etc.
  - ◆ e.g., hotel price: currency, tax, breakfast, ...



# Time-Variant

---

- ◆ Significantly longer time span
- ◆ operational database: current data
- ◆ data warehouse: historical perspective
  - ◆ e.g., past 5-10 years
- ◆ Every key structure in a data warehouse
  - ◆ contains time info, explicitly or implicitly
- ◆ key of operational data may not contain time info



# Nonvolatile

---

- ◆ A physically separate store of data transformed from operational environments
- ◆ no transaction processing, recovery, concurrency control
- ◆ Only two operations in data accessing
- ◆ No operational update of data
- ◆ initial loading of data
- ◆ access of data



# Data Warehouse vs. Operational DBMS

---

- ◆ OLTP (on-line transaction processing)
- ◆ major task of traditional relational DBMS
- ◆ day-to-day operations
- ◆ OLAP (on-line analytical processing)
- ◆ major task of data warehouse system
- ◆ data analysis and decision making



# OLTP vs. OLAP

---

	OLTP	OLAP
<b>users</b>	clerk, IT professional	knowledge worker
<b>function</b>	day to day operations	decision support
<b>DB design</b>	application-oriented	subject-oriented
<b>data</b>	current, up-to-date detailed, flat relational isolated	historical, summarized, multidimensional integrated, consolidated
<b>usage</b>	repetitive	ad-hoc
<b>access</b>	read/write index/hash on prim. key	lots of scans
<b>unit of work</b>	short, simple transaction	complex query
<b># records accessed</b>	tens	millions
<b>#users</b>	thousands	hundreds
<b>DB size</b>	100MB-GB	100GB-TB
<b>metric</b>	transaction throughput	query throughput, response



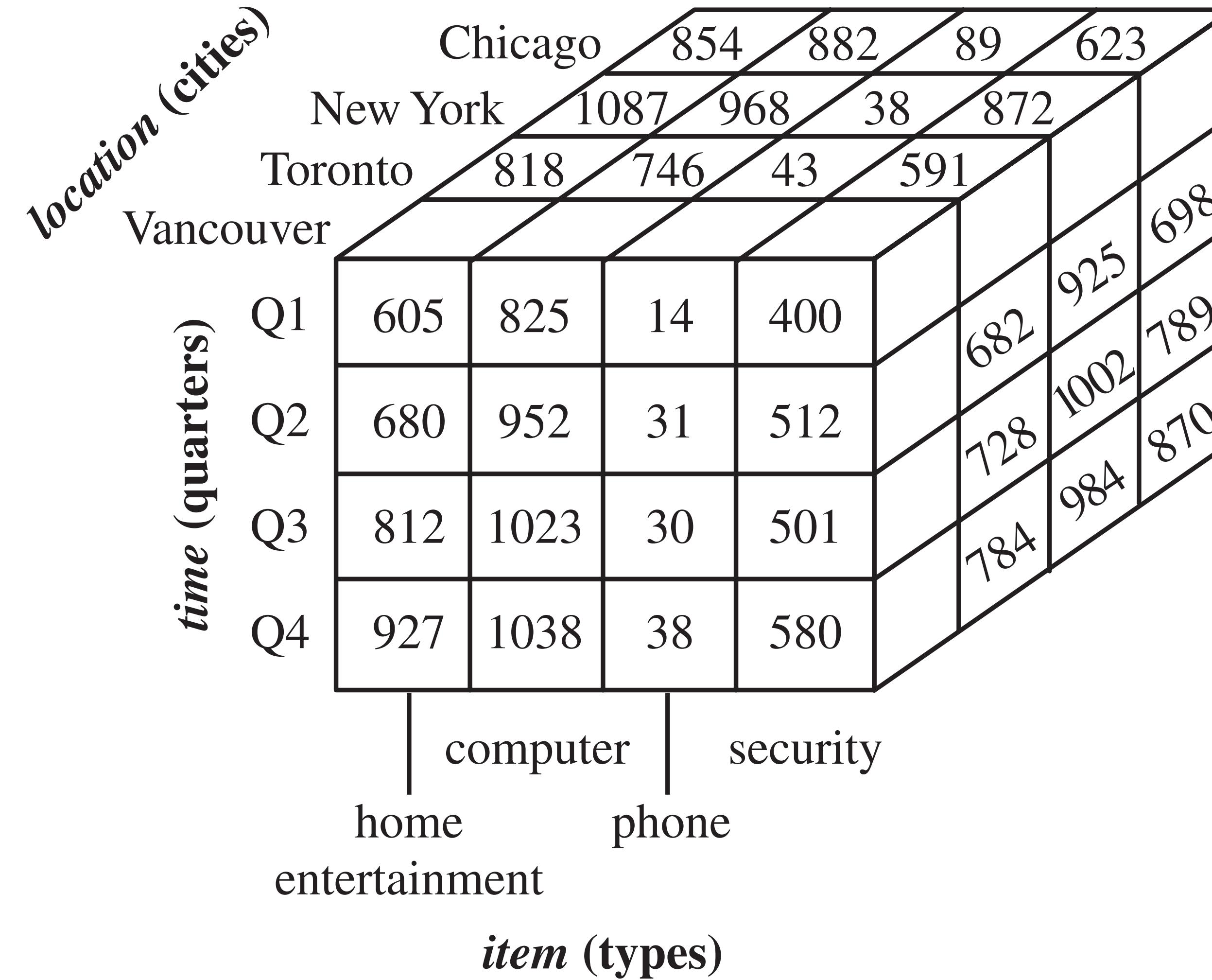
# What is A Data Cube?

---

- ◆ Data warehouses and OLAP are based on
  - ◆ a multi-dimensional data model
- ◆ Data cube
  - ◆ allow data to be modeled and viewed in multiple dimensions (e.g., sales)
  - ◆ dimensions: e.g., time, item, branch, location
  - ◆ facts: numerical measures, e.g., items\_sold, dollars\_sold



# Data Cube Example



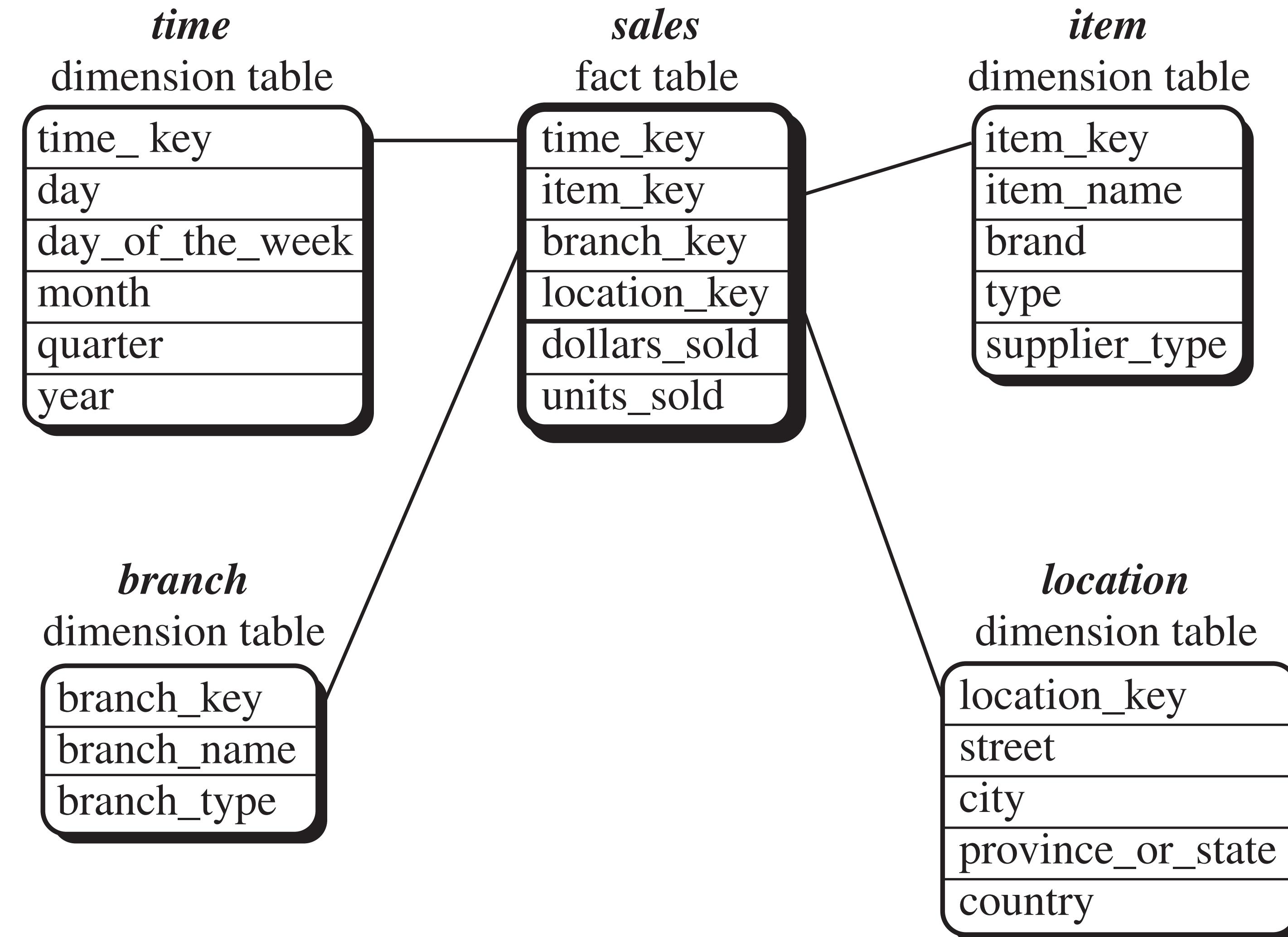
# Conceptual Modeling

---

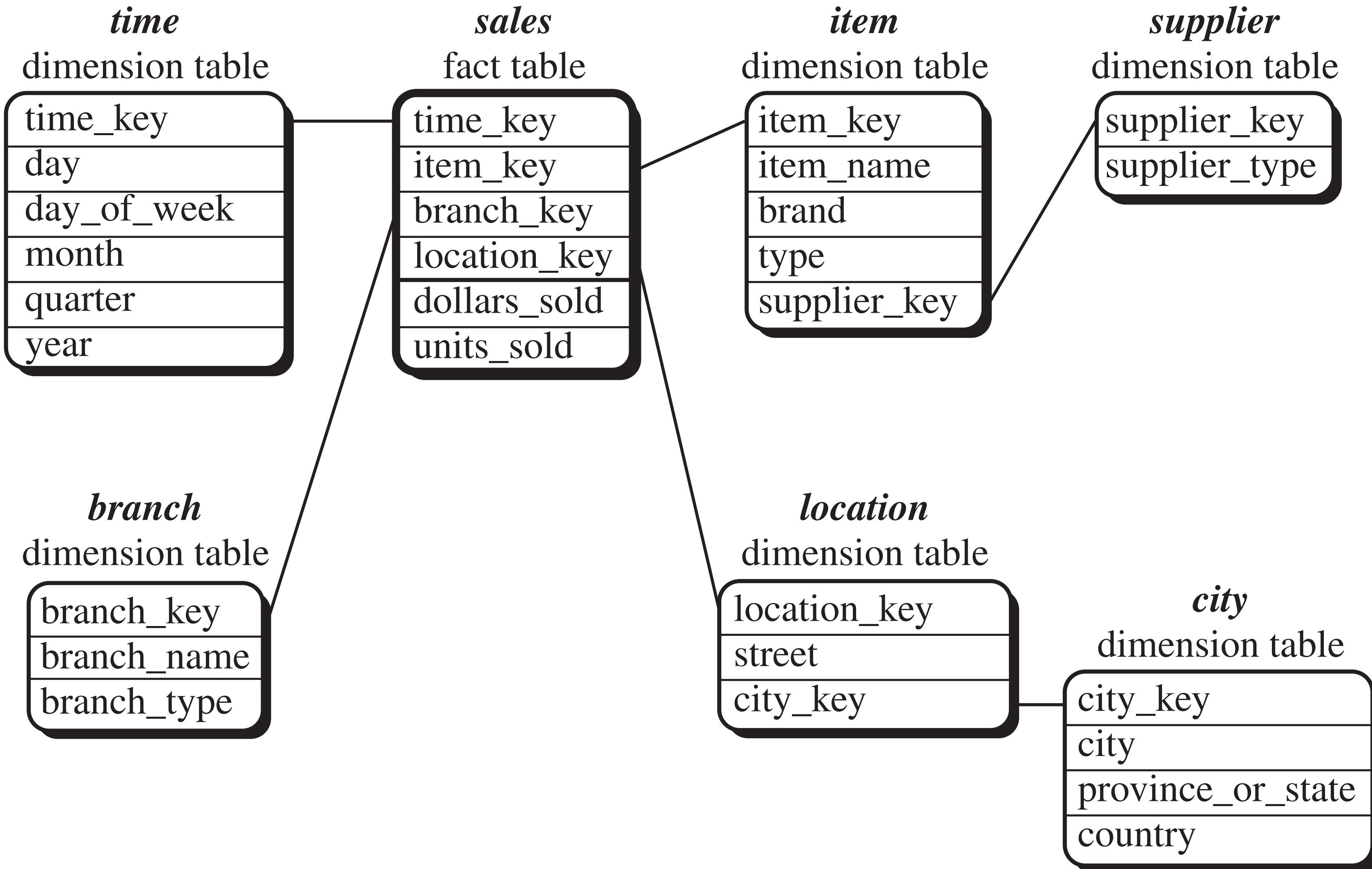
- ◆ Modeling data warehouses: Dimensions & facts
- ◆ Star schema
  - ◆ a fact table, a set of dimension tables
- ◆ Snowflake schema
  - ◆ a fact table, a hierarchy of dimension tables
- ◆ Fact constellations
  - ◆ multiple fact tables share dimension tables



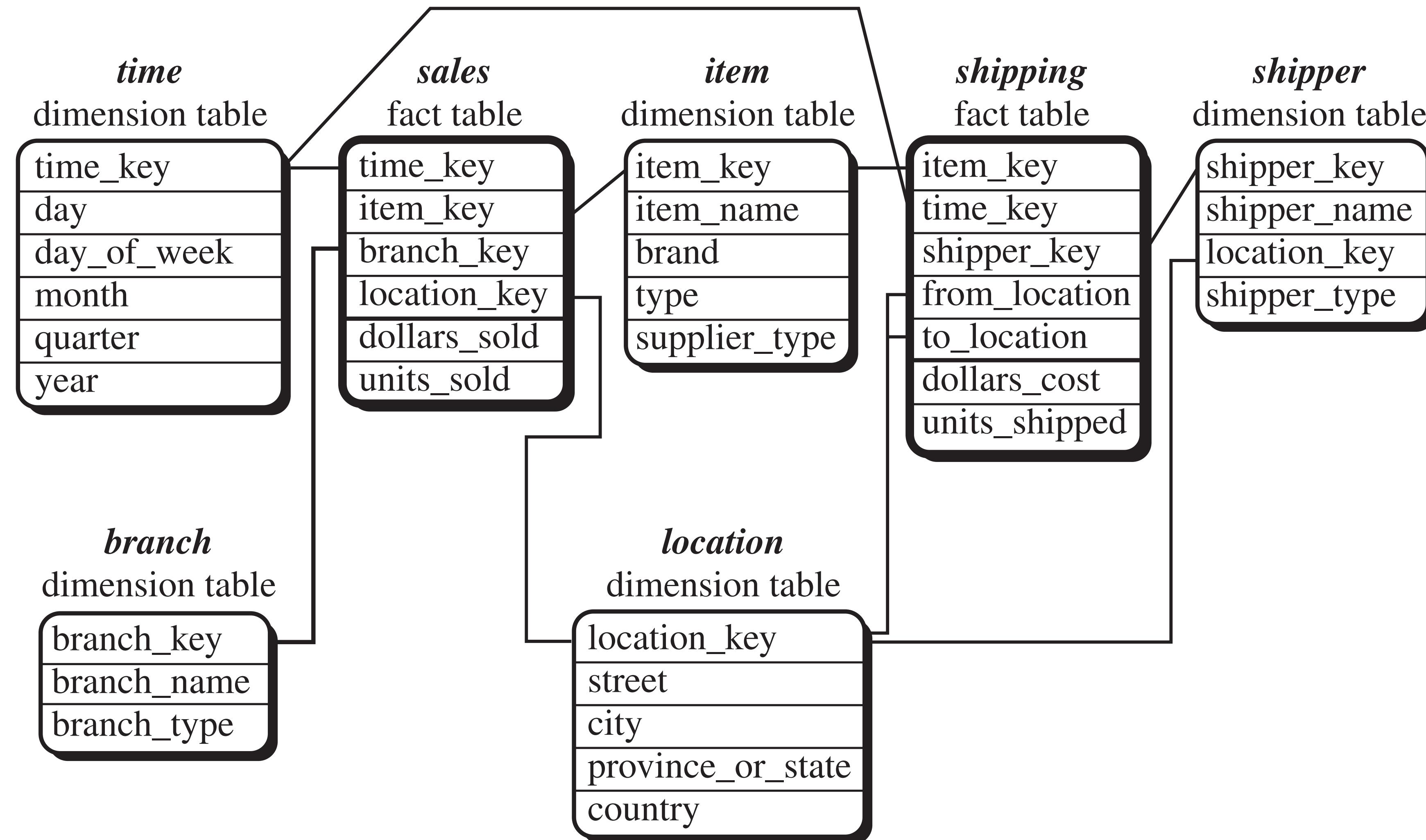
# Example of Star Schema



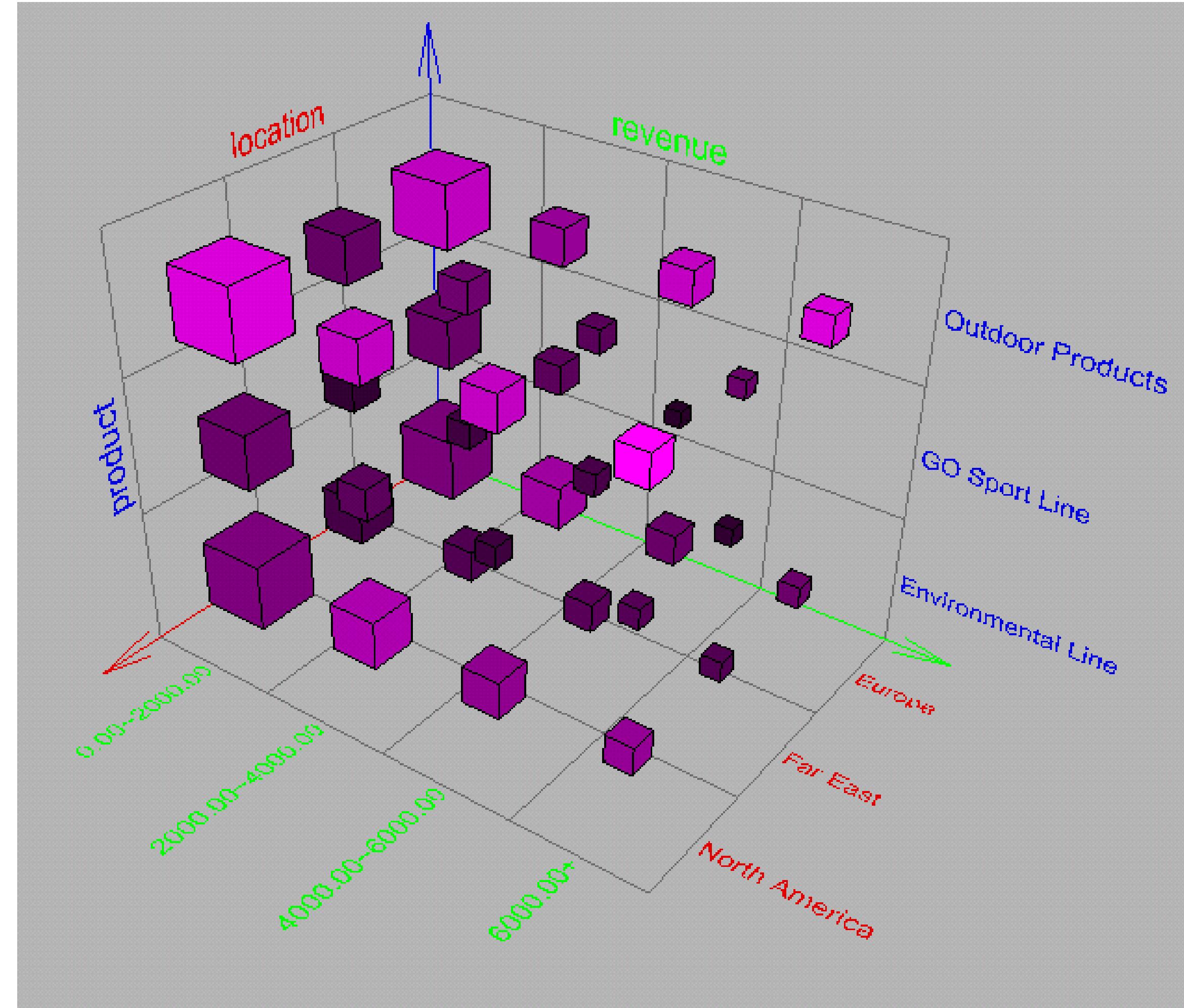
# Example of Snowflake Schema



# Example of Fact Constellation



# Browsing a Data Cube

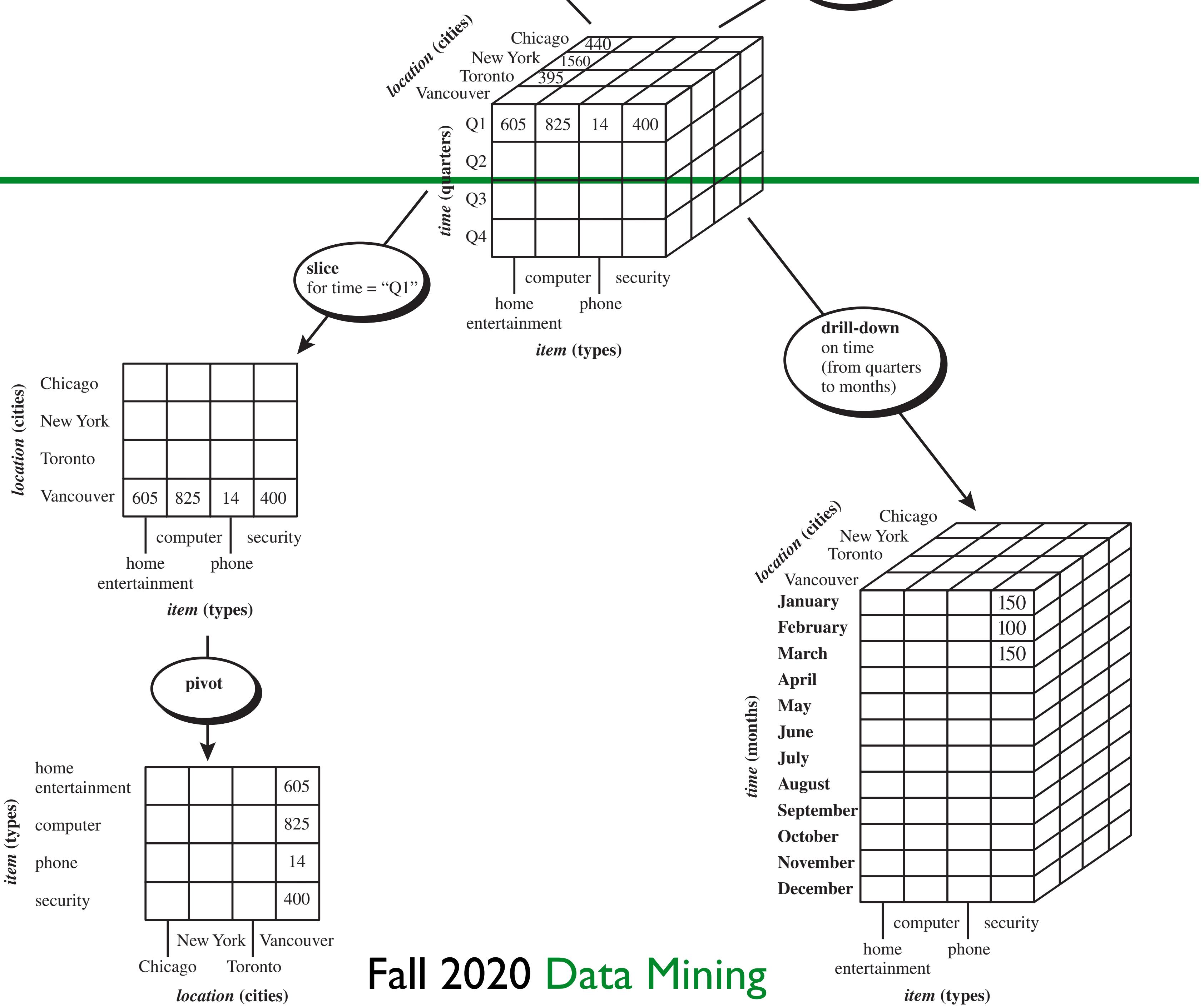


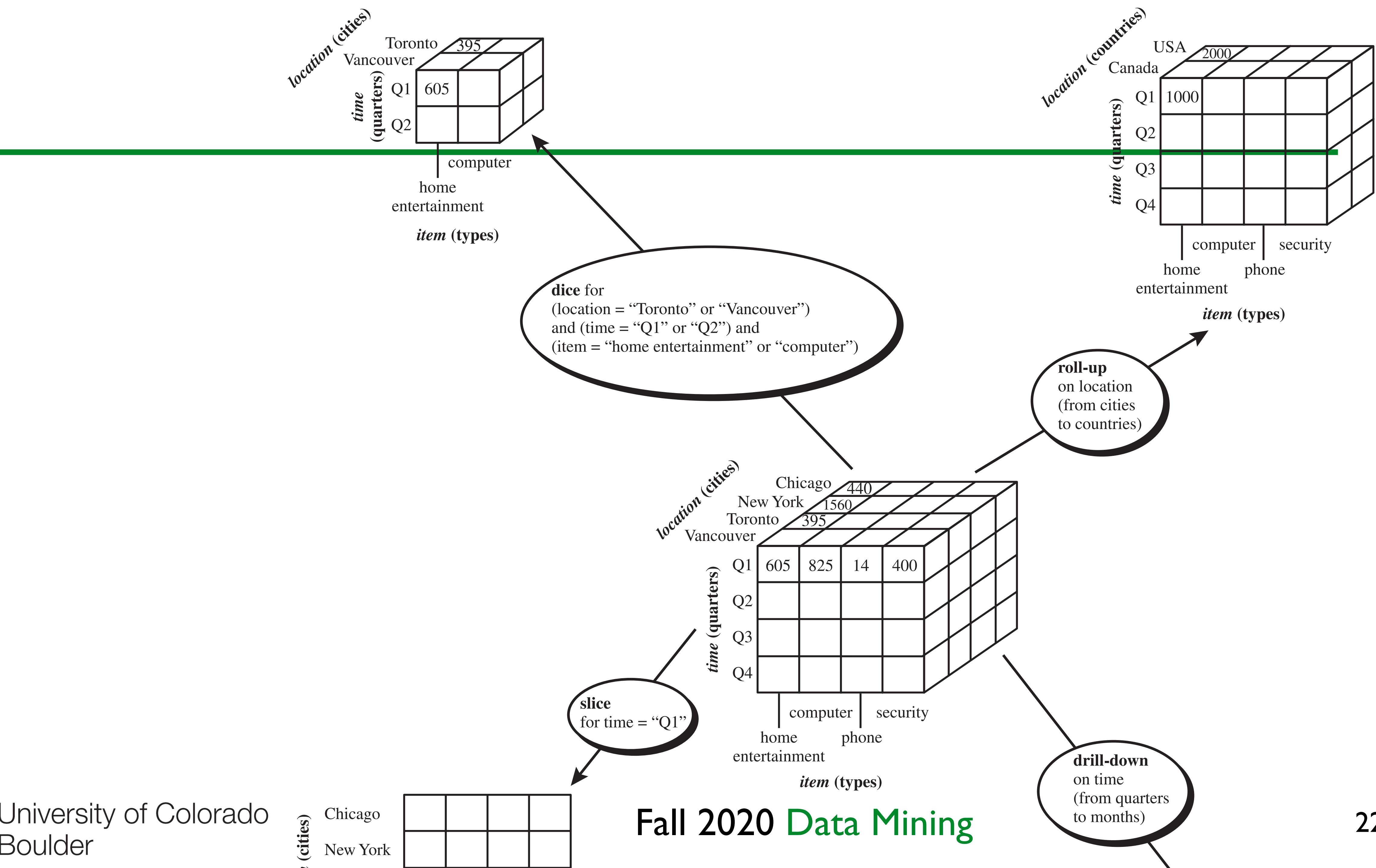
# Typical OLAP Operations

---

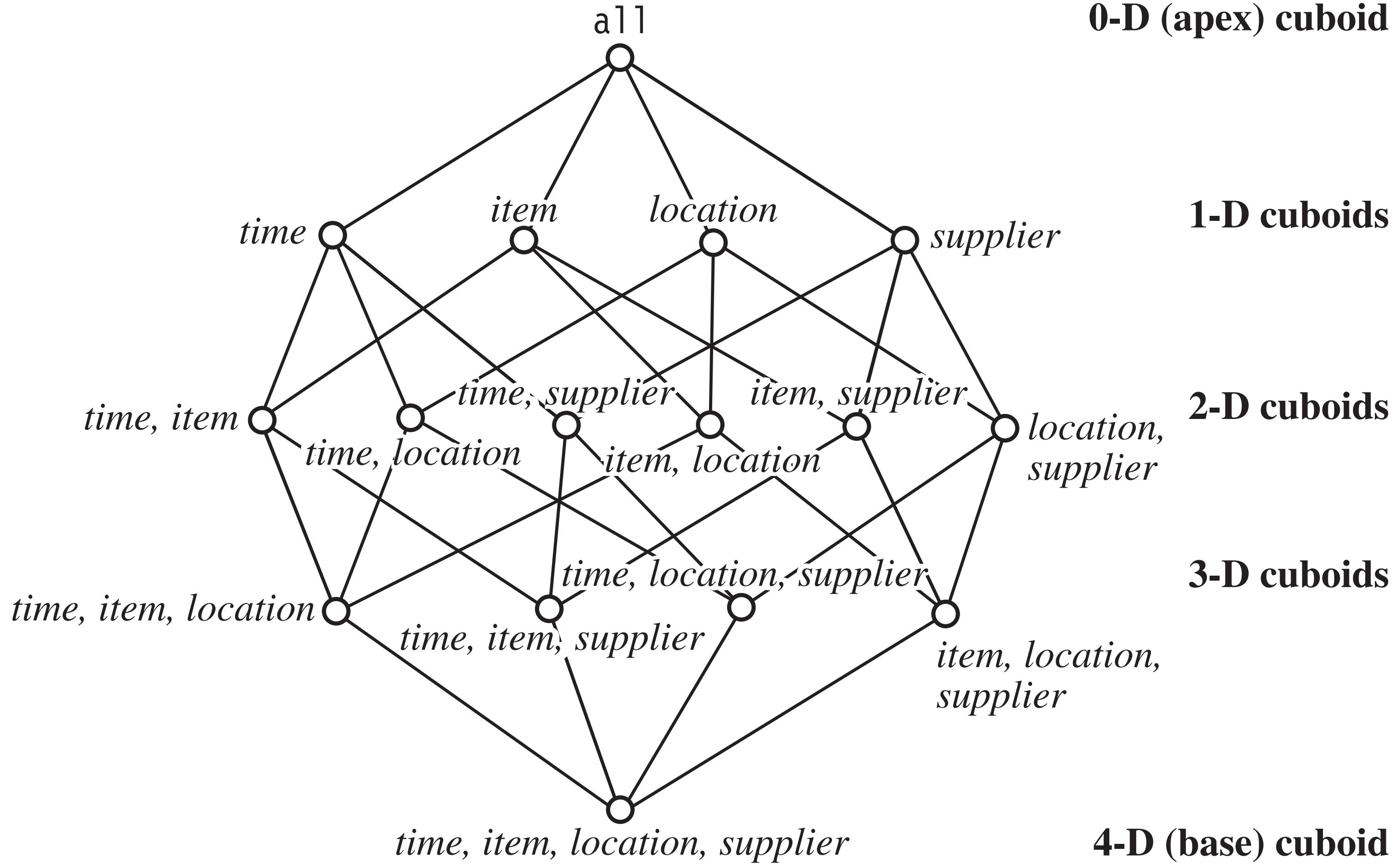
- ◆ **Roll-up (drill-up)**: summarization
- ◆ **Drill-down**: reverse of roll-up
- ◆ **Slice and dice**: project and select (sub-cube)
- ◆ **Pivot (rotate)**: visualization, 3D to 2Ds
- ◆ **drill-across**: more than one fact tables
- ◆ **drill-through**: to the back-end relational tables







# Cube: Lattice of Cuboids



# Cuboid Cells

---

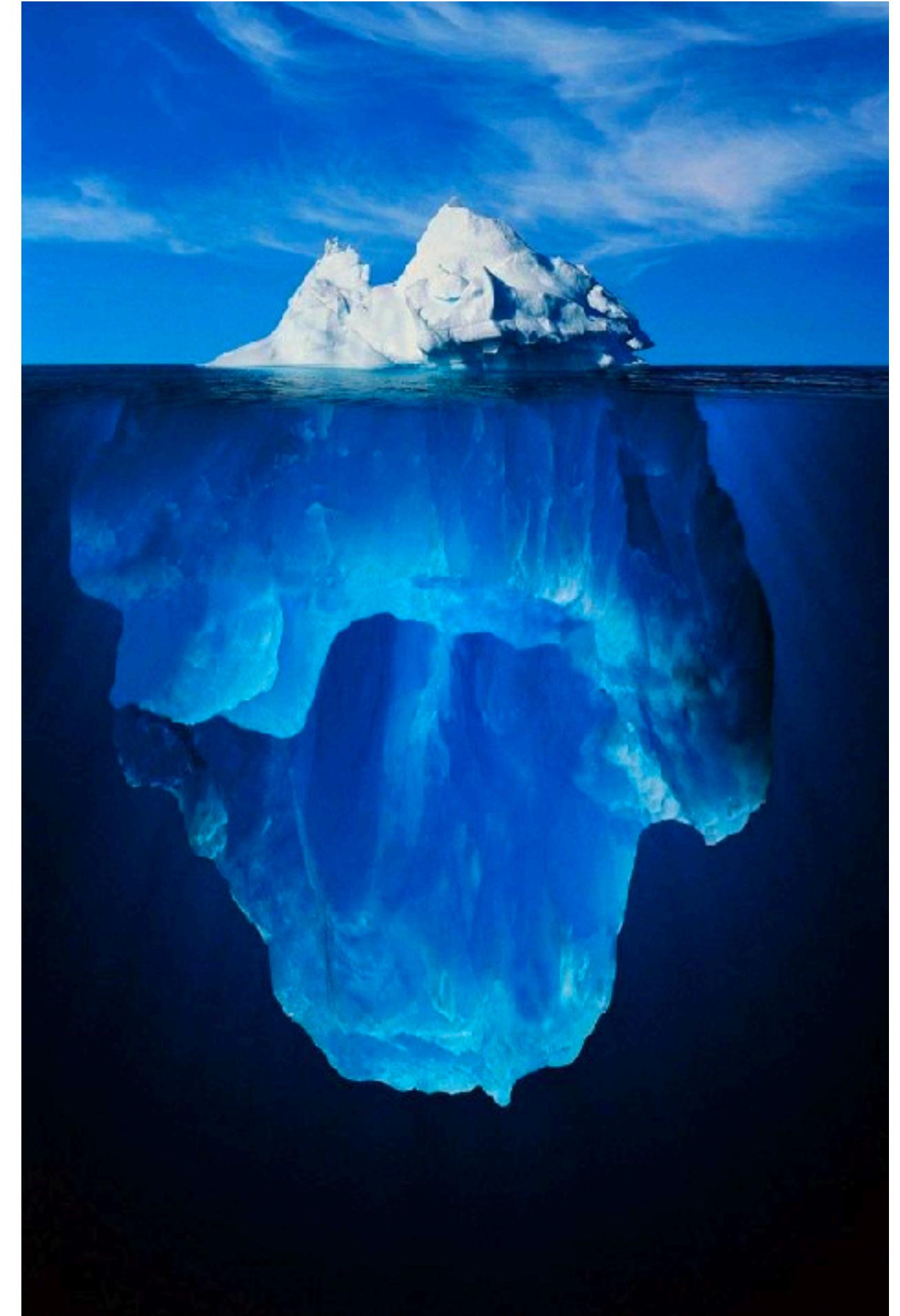
- ◆ **Cuboid cells**
  - ◆ base vs. aggregate, ancestor vs. descendant
  - ◆ E.g., (**month, location, customer\_group, price**)
    - ◆ a = (Jan, \*, \*, 2800)
    - ◆ b = (Jan, \*, Business, 150)
    - ◆ c = (Jan, Toronto, Business, 45)
- ◆ **Materialization** of data cube
  - ◆ full, partial, or no materialization



# Iceberg Cube

---

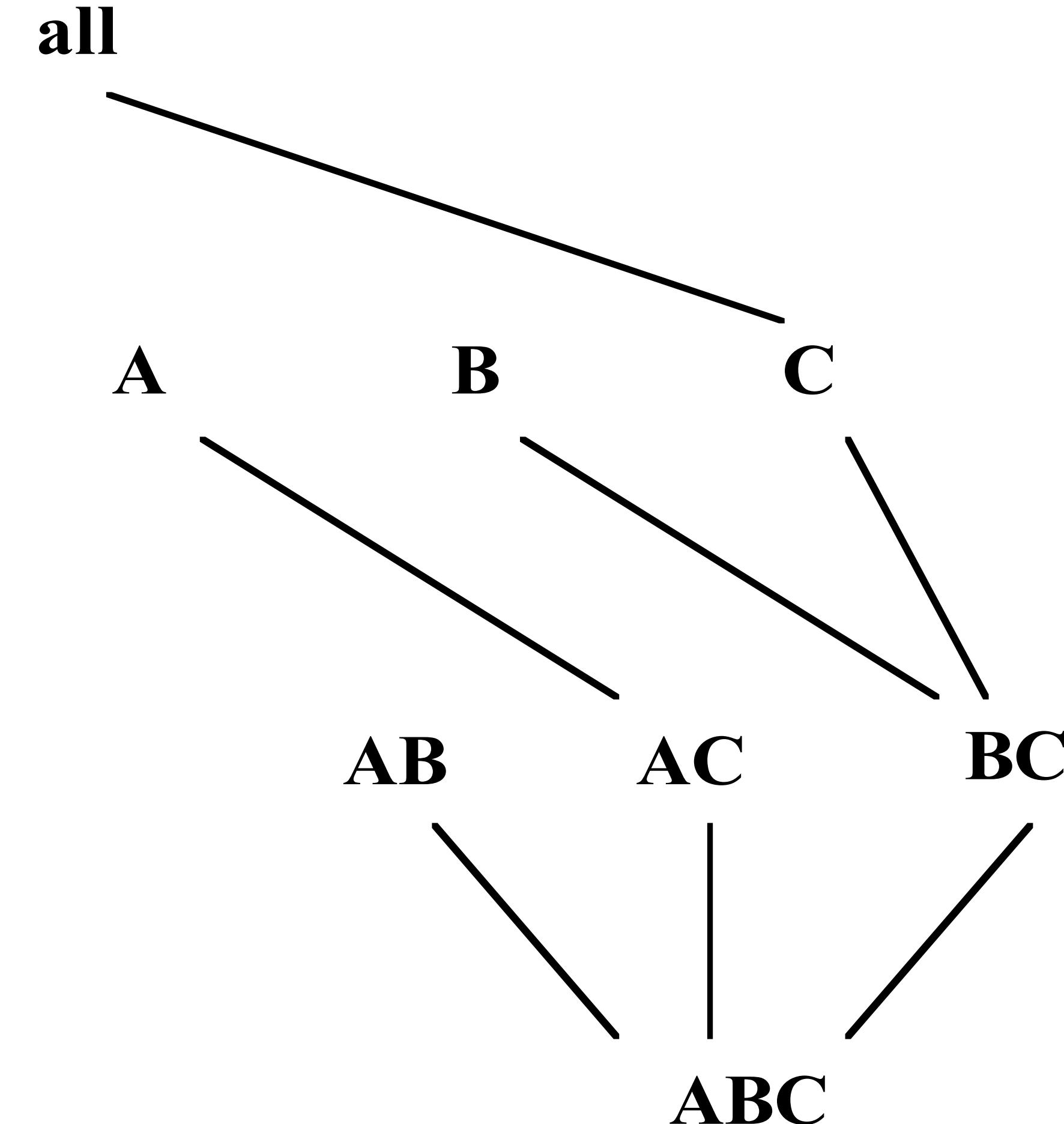
- ◆ Only a small portion may be “**above the water**” in a sparse cube
- ◆ Compute only the cuboid cells whose aggregate (e.g., count) is above a threshold
  - ◆ **minimum support**
- ◆ Avoid explosive growth of the cube



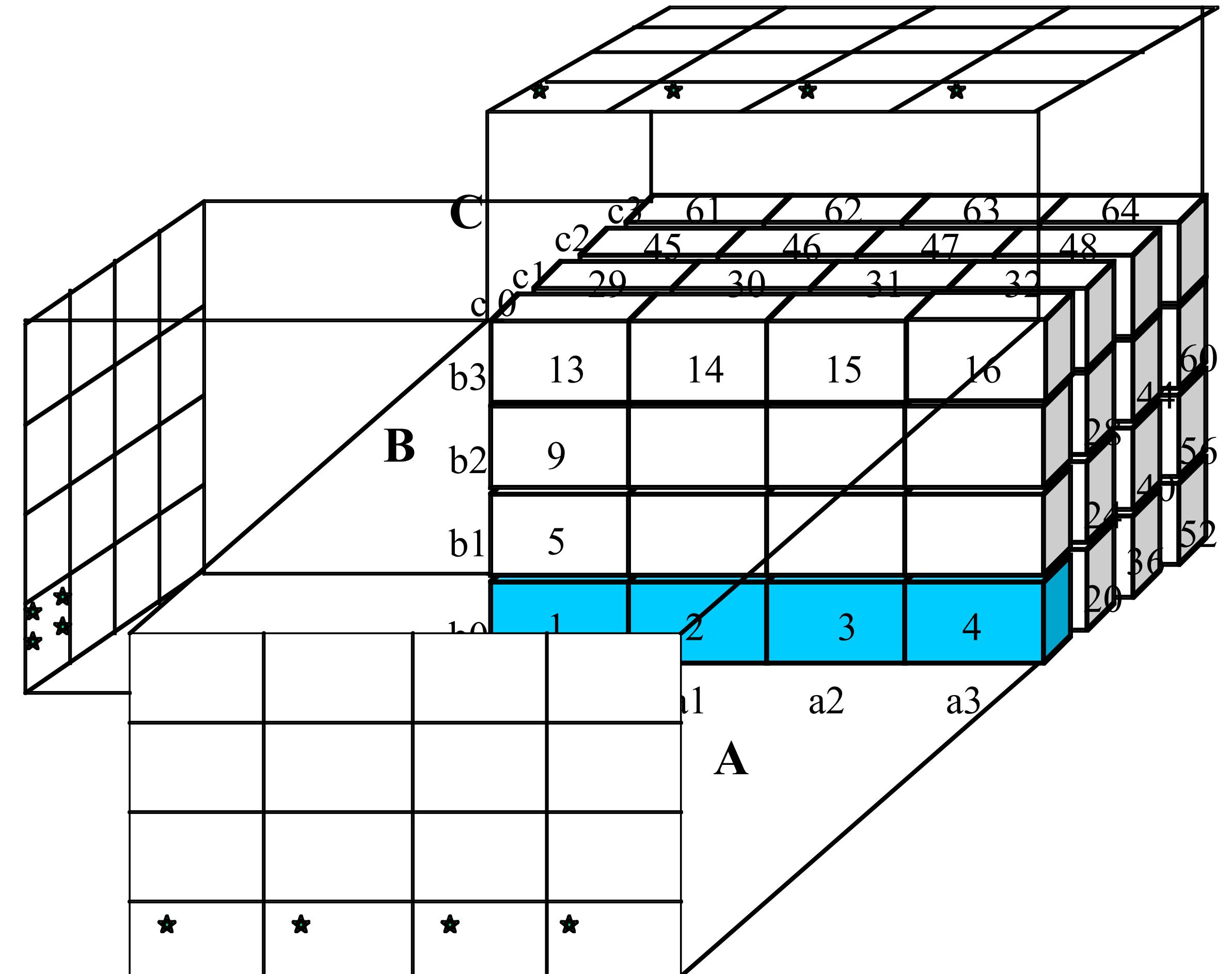
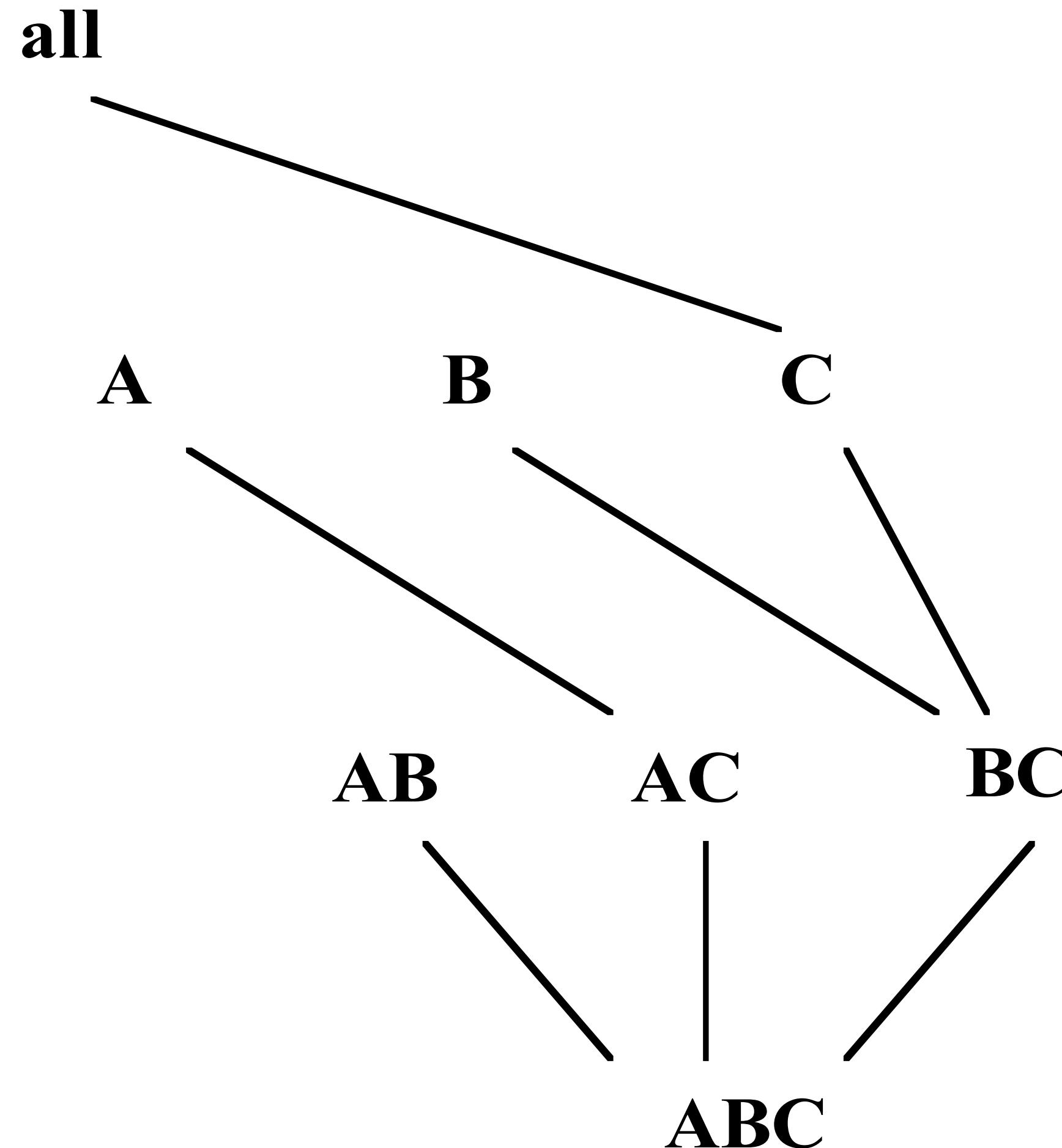
# Multi-Way Array Aggregation

---

- ◆ Full cube computation
- ◆ Array-based “bottom-up” algorithm
- ◆ Multi-dimensional chunks
- ◆ Simultaneous aggregation on multiple dimensions
- ◆ Cannot do Apriori pruning
- ◆ Not for high dimensions



# Multi-Way Array Aggregation

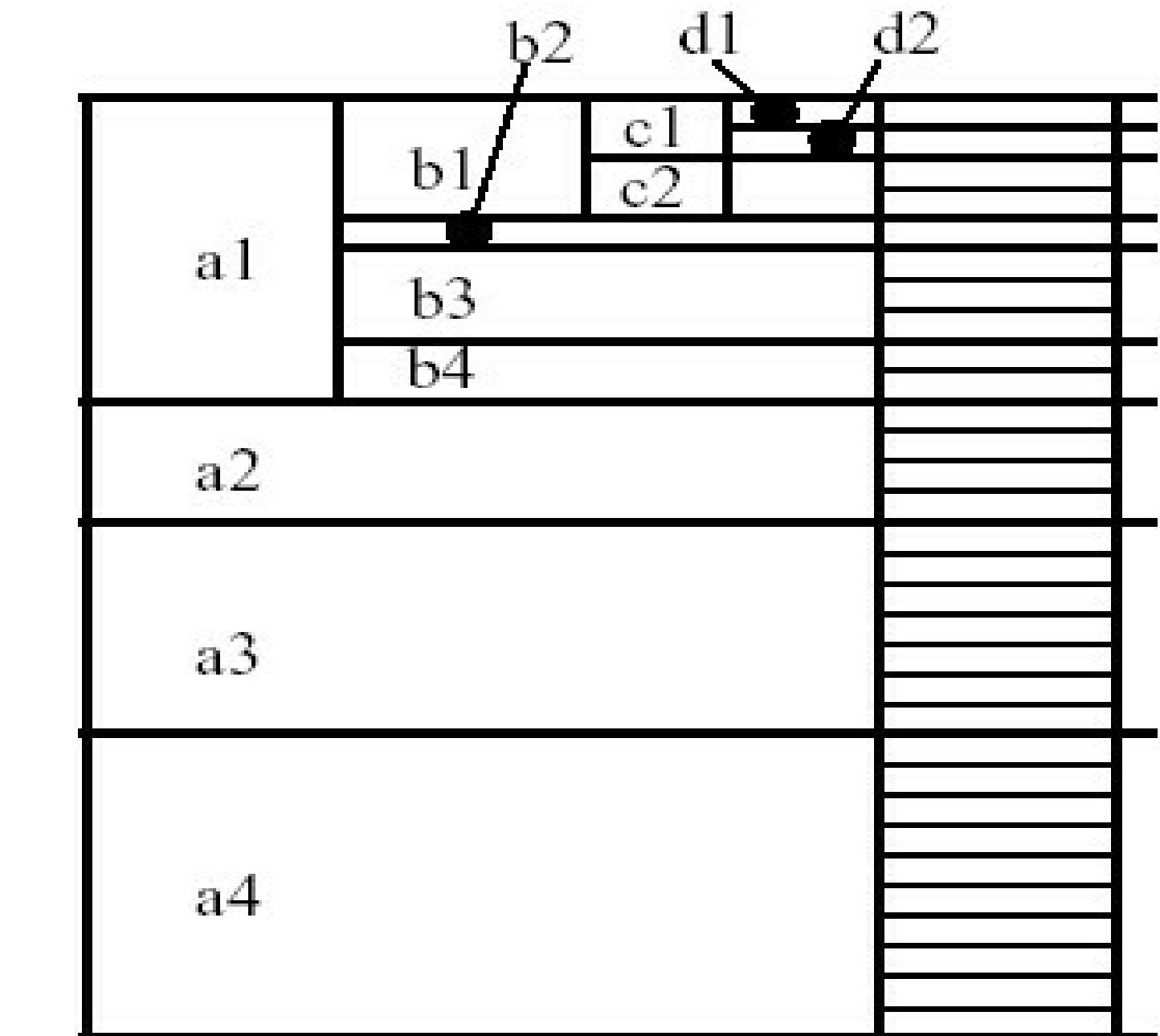
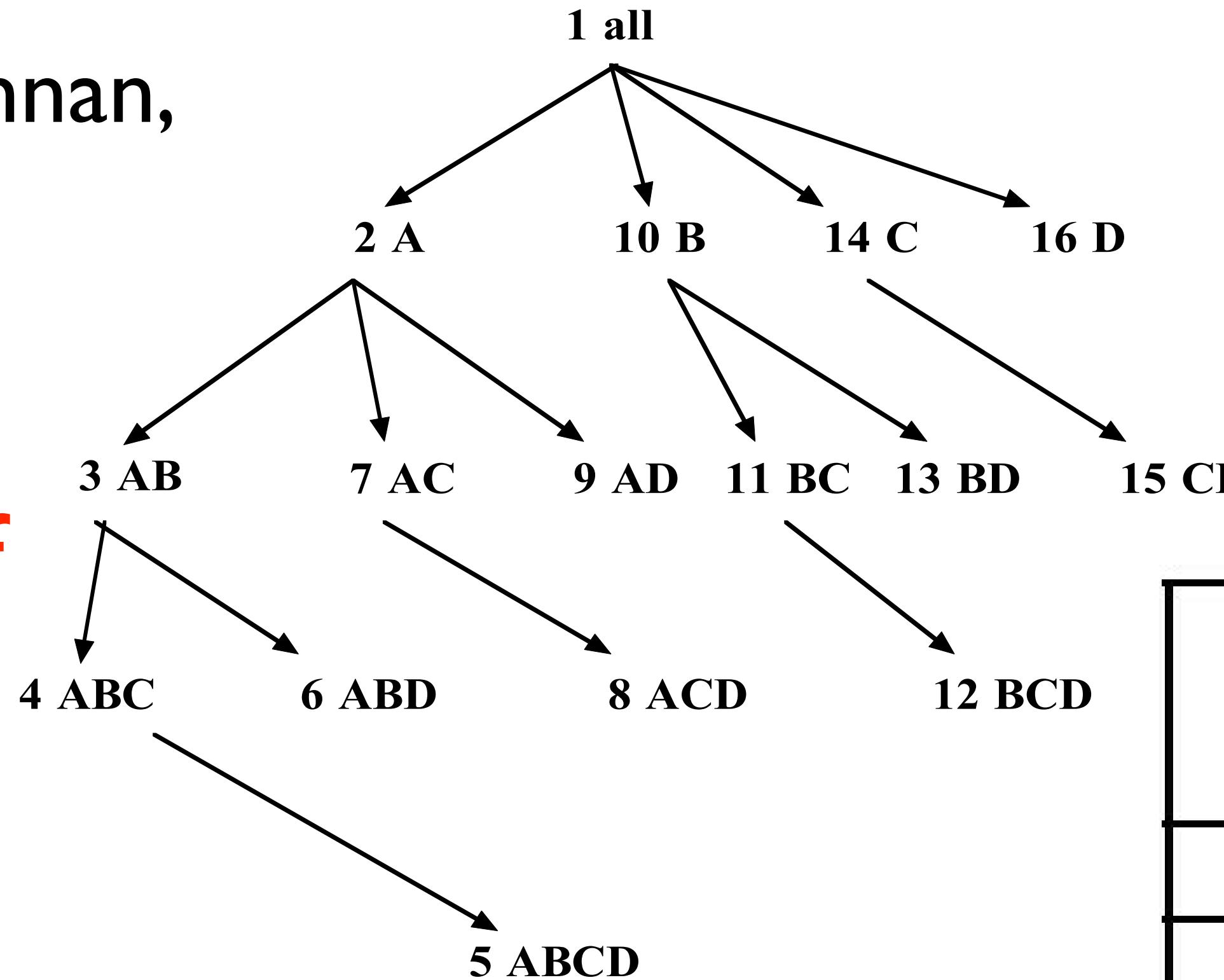


# Bottom-Up Computation (BUC)

- ◆ [Beyer & Ramakrishnan, SIGMOD'99]

- ◆ Top-down computation of iceberg cubes

- ◆ Divides dimensions into partitions and facilitates iceberg pruning



# Summary

---

- ◆ **Chap 4 & 5: Data Warehouse, Data Cube**
  - ◆ what is data warehouse?
  - ◆ OLTP vs. OLAP
  - ◆ what is data cube?
  - ◆ data cube operations
  - ◆ data cube computation



# Review: Part I

---

- ◆ Chapter 1: Introduction
  - ◆ Chapter 2: Getting to Know your Data
  - ◆ Chapter 3: Data Preprocessing
  - ◆ Chapter 4: Data Warehousing & Online Analytical Processing
  - ◆ Chapter 5: Data Cube Technology
- 
- ◆ **Part 2: Core DM Techniques**
  - ◆ **Part 3: Mining Complex Data, DM Trends**



# Course Project

---

- ◆ 40% of overall grade
- ◆ A self-contained project related to DM
- ◆ Work in teams (3-4 students)
- ◆ Pick your own project idea
- ◆ Project discussion at piazza, online resources
- ◆ Discuss w/ instructor & other students
- ◆ **Start early!!**



# Choose Your Project

---

- ◆ What are you interested in?
- ◆ Who are on your team?
- ◆ What data set(s) are available?
- ◆ What problem do you want to answer using the data?
- ◆ Why is the problem interesting/important?
- ◆ What are the challenges?
- ◆ What have been done before? Limitations?



# Define Your Project

---

- ◆ Project title
- ◆ Project team (4502/5502, expertise, tasks)
- ◆ Problem statement
- ◆ Motivation, literature survey
- ◆ Significance & difference given prior work
- ◆ Proposed work (data set, approaches)
- ◆ Evaluation: metrics, existing solutions
- ◆ Milestones



# Course Project Proposal

---

- ◆ Week 6
- ◆ Team formation & project identification
  - ◆ discuss w/ instructor & other students
- ◆ Project forum announcement
  - ◆ team, title, brief description
- ◆ Project proposal report
- ◆ Checkpoint (Week 12)
- ◆ Final report (Week 16)

