

# **W4111 – Introduction to Databases**

## **Final Lecture**

### **Various Topics**



# Chapter 7: Normalization

Database System Concepts, 7<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Outline

- Features of Good Relational Design
- Functional Dependencies
- Decomposition Using Functional Dependencies
- Normal Forms
- Functional Dependency Theory
- Algorithms for Decomposition using Functional Dependencies
- Decomposition Using Multivalued Dependencies
- More Normal Form
- Atomic Domains and First Normal Form
- Database-Design Process
- Modeling Temporal Data

This is all the material in the new textbook;  
Will not cover all of it.



# Overview of Normalization



# Features of Good Relational Designs

- Hypothetically, assume our schema originally had one relation that provided information about faculty and departments.
- Suppose we combine *instructor* and *department* into *in\_dep*, which represents the natural join on the relations *instructor* and *department*

<i>ID</i>	<i>name</i>	<i>salary</i>	<i>dept_name</i>	<i>building</i>	<i>budget</i>
22222	Einstein	95000	Physics	Watson	70000
12121	Wu	90000	Finance	Painter	120000
32343	El Said	60000	History	Painter	50000
45565	Katz	75000	Comp. Sci.	Taylor	100000
98345	Kim	80000	Elec. Eng.	Taylor	85000
76766	Crick	72000	Biology	Watson	90000
10101	Srinivasan	65000	Comp. Sci.	Taylor	100000
58583	Califieri	62000	History	Painter	50000
83821	Brandt	92000	Comp. Sci.	Taylor	100000
15151	Mozart	40000	Music	Packard	80000
33456	Gold	87000	Physics	Watson	70000
76543	Singh	80000	Finance	Painter	120000

- There is repetition of information
- Need to use null values (if we add a new department with no instructors)



# Decomposition

- The only way to avoid the repetition-of-information problem in the `in_dep` schema is to decompose it into two schemas – *instructor* and *department* schemas.
- Not all decompositions are good. Suppose we decompose

*employee*(*ID*, *name*, *street*, *city*, *salary*)

into

*employee1* (*ID*, *name*)

*employee2* (*name*, *street*, *city*, *salary*)

- The problem arises when we have two employees with the same name
- The next slide shows how we lose information -- we cannot reconstruct the original *employee* relation -- and so, this is a **lossy decomposition**.



# *The Evils of Redundancy*

- ❖ *Redundancy* is at the root of several problems associated with relational schemas:
  - redundant storage, insert/delete/update anomalies
- ❖ Integrity constraints, in particular *functional dependencies*, can be used to identify schemas with such problems and to suggest refinements.
- ❖ Main refinement technique: *decomposition* (replacing ABCD with, say, AB and BCD, or ACD and ABD).
- ❖ Decomposition should be used judiciously:
  - Is there reason to decompose a relation?
  - What problems (if any) does the decomposition cause?



## Example (Contd.)

- ❖ Problems due to  $R \rightarrow W$ :
  - Update anomaly: Can we change  $W$  in just the 1st tuple of SNLRWH?
  - Insertion anomaly: What if we want to insert an employee and don't know the hourly wage for his rating?
  - Deletion anomaly: If we delete all employees with rating 5, we lose the information about the wage for rating 5!

Will 2 smaller tables be better?

Wages

R	W
8	10
5	7

Hourly\_Emps2

S	N	L	R	H
123-22-3666	Attishoo	48	8	40
231-31-5368	Smiley	22	8	30
131-24-3650	Smethurst	35	5	30
434-26-3751	Guldu	35	5	32
612-67-4134	Madayan	35	8	40

S	N	L	R	W	H
123-22-3666	Attishoo	48	8	10	40
231-31-5368	Smiley	22	8	10	30
131-24-3650	Smethurst	35	5	7	30
434-26-3751	Guldu	35	5	7	32
612-67-4134	Madayan	35	8	10	40



# Functional Dependencies (FDs)

- ❖ A functional dependency  $X \rightarrow Y$  holds over relation R if, for every allowable instance  $r$  of R:
  - $t_1 \in r, t_2 \in r, \pi_X(t_1) = \pi_X(t_2)$  implies  $\pi_Y(t_1) = \pi_Y(t_2)$
  - i.e., given two tuples in  $r$ , if the X values agree, then the Y values must also agree. (X and Y are *sets* of attributes.)
- ❖ An FD is a statement about *all* allowable relations.
  - Must be identified based on semantics of application.
  - Given some allowable instance  $r_1$  of R, we can check if it violates some FD  $f$ , but we cannot tell if  $f$  holds over R!
- ❖ K is a candidate key for R means that  $K \rightarrow R$ 
  - However,  $K \rightarrow R$  does not require K to be *minimal*!



## *Example: Constraints on Entity Set*

- ❖ Consider relation obtained from Hourly\_Emps:
  - Hourly\_Emps (*ssn, name, lot, rating, hrly\_wages, hrs\_worked*)
- ❖ Notation: We will denote this relation schema by listing the attributes: **SNLRWH**
  - This is really the *set* of attributes {S,N,L,R,W,H}.
  - Sometimes, we will refer to all attributes of a relation by using the relation name. (e.g., Hourly\_Emps for SNLRWH)
- ❖ Some FDs on Hourly\_Emps:
  - *ssn* is the key:  $S \rightarrow \text{SNLRWH}$
  - *rating* determines *hrly\_wages*:  $R \rightarrow W$



## Example: Constraints on Entity Set

- ❖ Consider relation obtained from Hourly\_Emps:
  - Hourly\_Emps (ssn, name, lot, rating, hrly\_wages, hrs\_worked)

❖ Notation: We will denote this relation schema by

A more intuitive example – Consider a table of addresses.

- Addresses(id, street\_no, street\_name, city, country, zipcode)
- (city, country) are *functionally dependent* on zipcode.

❖ Some FDs on Hourly\_Emps:

- *ssn is the key:*  $S \rightarrow SNLRWH$
- *rating determines hrly\_wages:*  $R \rightarrow W$



# Reasoning About FDs

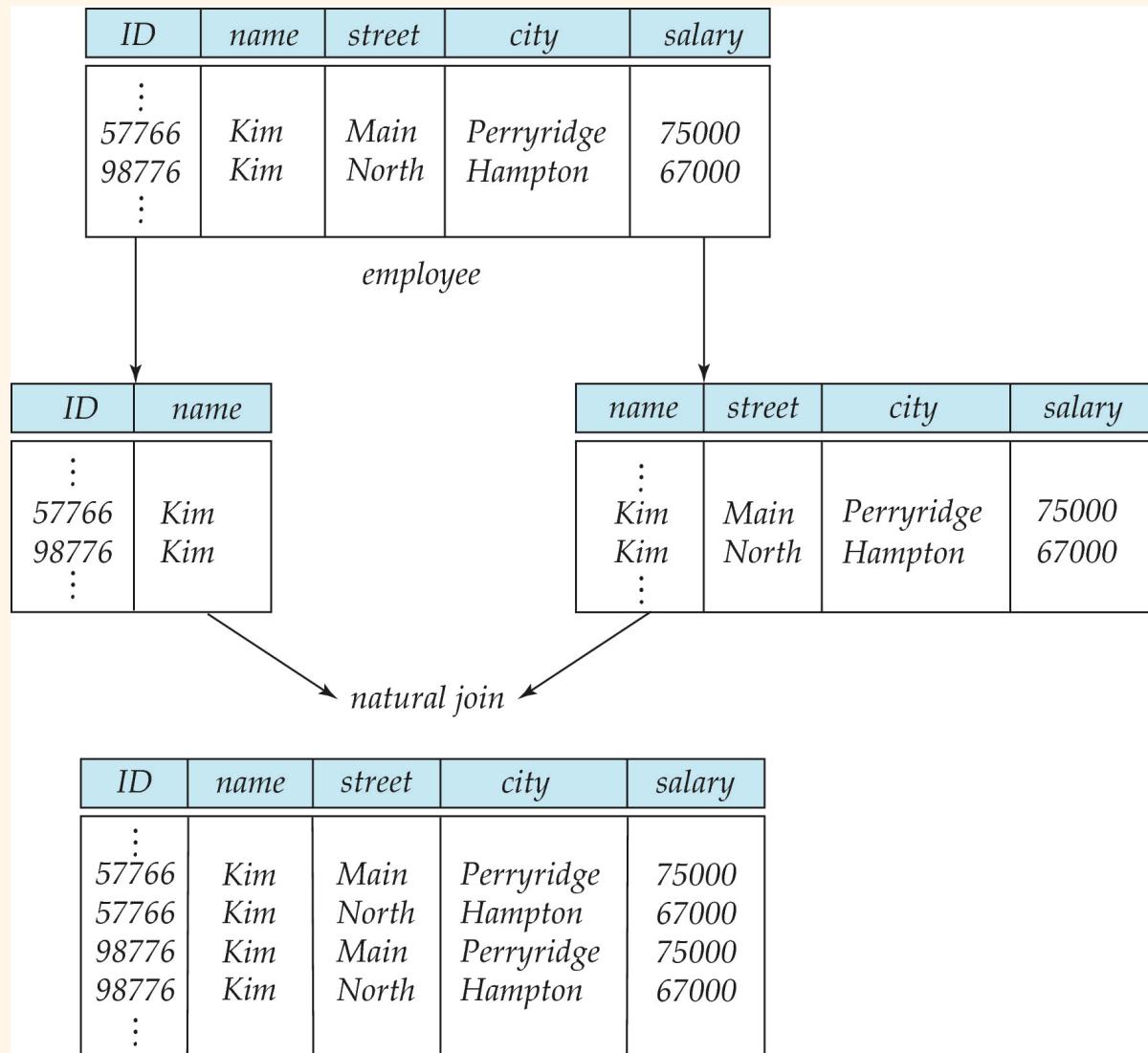
- ❖ Given some FDs, we can usually infer additional FDs:
  - $ssn \rightarrow did$ ,  $did \rightarrow lot$  implies  $ssn \rightarrow lot$
- ❖ An FD  $f$  is *implied by* a set of FDs  $F$  if  $f$  holds whenever all FDs in  $F$  hold.
  - $F^+ = \text{closure of } F$  is the set of all FDs that are implied by  $F$ .
- ❖ Armstrong's Axioms ( $X, Y, Z$  are sets of attributes):
  - Reflexivity: If  $X \subseteq Y$ , then  $Y \rightarrow X$
  - Augmentation: If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$  for any  $Z$
  - Transitivity: If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$
- ❖ These are *sound* and *complete* inference rules for FDs!

# Decomposition

## ■ *Northwind Database Examples*



# A Lossy Decomposition





# *Lossless Decomposition*

- ❖ Let  $R$  be a relation schema and let  $R_1$  and  $R_2$  form a decomposition of  $R$ . That is  $R = R_1 \cup R_2$
- ❖ We say that the decomposition is a **lossless decomposition** if there is no loss of information by replacing  $R$  with the two relation schemas  $R_1 \cup R_2$
- ❖ Formally,

$$\Pi_{R_1}(r) \quad \Pi_{R_2}(r) = r$$

- ❖ And, conversely a decomposition is lossy if

$$r \subset \Pi_{R_1}(r) \quad \Pi_{R_2}(r) = r$$

$\bowtie$



# Normalization Theory

- Decide whether a particular relation  $R$  is in “good” form.
- In the case that a relation  $R$  is not in “good” form, decompose it into set of relations  $\{R_1, R_2, \dots, R_n\}$  such that
  - Each relation is in good form
  - The decomposition is a lossless decomposition
- Our theory is based on:
  - Functional dependencies
  - Multivalued dependencies



# Functional Dependencies

- There are usually a variety of constraints (rules) on the data in the real world.
- For example, some of the constraints that are expected to hold in a university database are:
  - Students and instructors are uniquely identified by their ID.
  - Each student and instructor has only one name.
  - Each instructor and student is (primarily) associated with only one department.
  - Each department has only one value for its budget, and only one associated building.



# Functional Dependencies (Cont.)

- An instance of a relation that satisfies all such real-world constraints is called a **legal instance** of the relation;
- A legal instance of a database is one where all the relation instances are legal instances
- Constraints on the set of legal relations.
- Require that the value for a certain set of attributes determines uniquely the value for another set of attributes.
- A functional dependency is a generalization of the notion of a *key*.



# Functional Dependencies Definition

- Let  $R$  be a relation schema

$$\alpha \subseteq R \text{ and } \beta \subseteq R$$

- The **functional dependency**

$$\alpha \rightarrow \beta$$

**holds on  $R$**  if and only if for any legal relations  $r(R)$ , whenever any two tuples  $t_1$  and  $t_2$  of  $r$  agree on the attributes  $\alpha$ , they also agree on the attributes  $\beta$ . That is,

$$t_1[\alpha] = t_2[\alpha] \Rightarrow t_1[\beta] = t_2[\beta]$$

- Example: Consider  $r(A,B)$  with the following instance of  $r$ .

1	4
1	5
3	7

- On this instance,  $B \rightarrow A$  hold;  $A \rightarrow B$  does **NOT** hold,



# Closure of a Set of Functional Dependencies

- Given a set  $F$  set of functional dependencies, there are certain other functional dependencies that are logically implied by  $F$ .
  - If  $A \rightarrow B$  and  $B \rightarrow C$ , then we can infer that  $A \rightarrow C$
  - etc.
- The set of **all** functional dependencies logically implied by  $F$  is the **closure** of  $F$ .
- We denote the *closure* of  $F$  by  $F^+$ .



# Keys and Functional Dependencies

- $K$  is a superkey for relation schema  $R$  if and only if  $K \rightarrow R$
- $K$  is a candidate key for  $R$  if and only if
  - $K \rightarrow R$ , and
  - for no  $\alpha \subset K$ ,  $\alpha \rightarrow R$
- Functional dependencies allow us to express constraints that cannot be expressed using superkeys. Consider the schema:

*in\_dep (ID, name, salary, dept\_name, building, budget ).*

We expect these functional dependencies to hold:

$dept\_name \rightarrow building$

$ID \rightarrow building$

but would not expect the following to hold:

$dept\_name \rightarrow salary$

DFF Note:

- In the current data:  $ID \rightarrow dept\_name \rightarrow building$
- But
- In the schema,  $dept\_name$  may be NULL..



# Use of Functional Dependencies

- We use functional dependencies to:
  - To test relations to see if they are legal under a given set of functional dependencies.
    - ▶ If a relation  $r$  is legal under a set  $F$  of functional dependencies, we say that  $r$  **satisfies**  $F$ .
  - To specify constraints on the set of legal relations
    - ▶ We say that  $F$  **holds on**  $R$  if all legal relations on  $R$  satisfy the set of functional dependencies  $F$ .
- Note: A specific instance of a relation schema may satisfy a functional dependency even if the functional dependency does not hold on all legal instances.
  - For example, a specific instance of *instructor* may, by chance, satisfy  $name \rightarrow ID$ .



# Trivial Functional Dependencies

- A functional dependency is **trivial** if it is satisfied by all instances of a relation
- Example:
  - $ID, name \rightarrow ID$
  - $name \rightarrow name$
- In general,  $\alpha \rightarrow \beta$  is trivial if  $\beta \subseteq \alpha$



# Lossless Decomposition

- We can use functional dependencies to show when certain decompositions are lossless.
- For the case of  $R = (R_1, R_2)$ , we require that for all possible relations  $r$  on schema  $R$

$$r = \Pi_{R1}(r) \bowtie \Pi_{R2}(r)$$

- A decomposition of  $R$  into  $R_1$  and  $R_2$  is lossless decomposition if at least one of the following dependencies is in  $F^+$ :
  - $R_1 \cap R_2 \rightarrow R_1$
  - $R_1 \cap R_2 \rightarrow R_2$
- The above functional dependencies are a sufficient condition for lossless join decomposition; the dependencies are a necessary condition only if all constraints are functional dependencies



# Example

- $R = (A, B, C)$   
 $F = \{A \rightarrow B, B \rightarrow C\}$
- $R_1 = (A, B), R_2 = (B, C)$ 
  - Lossless decomposition:  
$$R_1 \cap R_2 = \{B\} \text{ and } B \rightarrow BC$$
- $R_1 = (A, B), R_2 = (A, C)$ 
  - Lossless decomposition:  
$$R_1 \cap R_2 = \{A\} \text{ and } A \rightarrow AB$$
- Note:
  - $B \rightarrow BC$   
is a shorthand notation for
  - $B \rightarrow \{B, C\}$



# Dependency Preservation

- Testing functional dependency constraints each time the database is updated can be costly
- It is useful to design the database in a way that constraints can be tested efficiently.
- If testing a functional dependency can be done by considering just one relation, then the cost of testing this constraint is low
- When decomposing a relation it is possible that it is no longer possible to do the testing without having to perform a Cartesian Product.
- A decomposition that makes it computationally hard to enforce functional dependency is said to be NOT **dependency preserving**.



# Dependency Preservation Example

- Consider a schema:

$\text{dept\_advisor}(s\_ID, i\_ID, \text{department\_name})$

- With function dependencies:

$i\_ID \rightarrow \text{dept\_name}$

$s\_ID, \text{dept\_name} \rightarrow i\_ID$

- In the above design we are forced to repeat the department name once for each time an instructor participates in a  $\text{dept\_advisor}$  relationship.
- To fix this, we need to decompose  $\text{dept\_advisor}$
- Any decomposition will not include all the attributes in
  - $s\_ID, \text{dept\_name} \rightarrow i\_ID$
- Thus, the composition NOT be dependency preserving

DFF Note:

This is not intuitive until you realize the intent is that all students in a department have the same advisor.



# Normal Forms



# Boyce-Codd Normal Form

- A relation schema  $R$  is in BCNF with respect to a set  $F$  of functional dependencies if for all functional dependencies in  $F^+$  of the form

$$\alpha \rightarrow \beta$$

where  $\alpha \subseteq R$  and  $\beta \subseteq R$ , at least one of the following holds:

- $\alpha \rightarrow \beta$  is trivial (i.e.,  $\beta \subseteq \alpha$ )
- $\alpha$  is a superkey for  $R$

## DFF Note:

- The theoretical treatment is no conveying the practical intent.
- If  $\alpha$  is a superkey, I can set a primary key/unique constraint on  $\alpha$ .
- Consider tables with address info in the rows.



# Decomposing a Schema into BCNF

- Let  $R$  be a schema  $R$  that is not in BCNF. Let  $\alpha \rightarrow \beta$  be the FD that causes a violation of BCNF.
- We decompose  $R$  into:
  - $(\alpha \cup \beta)$
  - $(R - (\beta - \alpha))$
- In our example of  $in\_dep$ ,
  - $\alpha = dept\_name$
  - $\beta = building, budget$and  $in\_dep$  is replaced by
  - $(\alpha \cup \beta) = (dept\_name, building, budget)$
  - $(R - (\beta - \alpha)) = (ID, name, dept\_name, salary)$

DFF Note – again this is baffling

- $R = (id, name\_last, name\_first, street, city, state, zipcode)$
- $\alpha \rightarrow \beta$  means  $zipcode \rightarrow (city, state)$
- $(\alpha \cup \beta) = (zipcode, city, state)$ , which we call Address
- $R - (\beta - \alpha) = R - \beta + \alpha = (id, name\_last, name\_first, zipcode)$ , which we call Person
- Setting primary key ID in Address and zipcode on Address preserves the dependency and is lossless.



# BCNF and Dependency Preservation

- It is not always possible to achieve both BCNF and dependency preservation
- Consider a schema:  
$$\text{dept\_advisor}(s\_ID, i\_ID, \text{department\_name})$$
- With function dependencies:  
$$i\_ID \rightarrow \text{dept\_name}$$
  
$$s\_ID, \text{dept\_name} \rightarrow i\_ID$$
- $\text{dept\_advisor}$  is not in BCNF
  - $i\_ID$  is not a superkey.
- Any decomposition of  $\text{dept\_advisor}$  will not include all the attributes in  
$$s\_ID, \text{dept\_name} \rightarrow i\_ID$$
- Thus, the composition is NOT be dependency preserving



# Third Normal Form

- A relation schema  $R$  is in **third normal form (3NF)** if for all:

$$\alpha \rightarrow \beta \text{ in } F^+$$

at least one of the following holds:

- $\alpha \rightarrow \beta$  is trivial (i.e.,  $\beta \in \alpha$ )
- $\alpha$  is a superkey for  $R$
- Each attribute  $A$  in  $\beta - \alpha$  is contained in a candidate key for  $R$ .

(NOTE: each attribute may be in a different candidate key)

- If a relation is in BCNF it is in 3NF (since in BCNF one of the first two conditions above must hold).
- Third condition is a minimal relaxation of BCNF to ensure dependency preservation (will see why later).



# 3NF Example

- Consider a schema:

$\text{dept\_advisor}(s\_ID, i\_ID, \text{dept\_name})$

- With function dependencies:

$i\_ID \rightarrow \text{dept\_name}$

$s\_ID, \text{dept\_name} \rightarrow i\_ID$

- Two candidate keys =  $\{s\_ID, \text{dept\_name}\}, \{s\_ID, i\_ID\}$

- We have seen before that  $\text{dept\_advisor}$  is not in BCNF

- $R$ , however, is in 3NF

- $s\_ID, \text{dept\_name}$  is a superkey

- $i\_ID \rightarrow \text{dept\_name}$  and  $i\_ID$  is NOT a superkey, but:

- ▶  $\{\text{dept\_name}\} - \{i\_ID\} = \{\text{dept\_name}\}$  and

- ▶  $\text{dept\_name}$  is contained in a candidate key



# Redundancy in 3NF

- Consider the schema R below, which is in 3NF

- $R = (J, K, L)$
- $F = \{JK \rightarrow L, L \rightarrow K\}$
- And an instance table:

$J$	$L$	$K$
$j_1$	$l_1$	$k_1$
$j_2$	$l_1$	$k_1$
$j_3$	$l_1$	$k_1$
<i>null</i>	$l_2$	$k_2$

- What is wrong with the table?

- Repetition of information
- Need to use null values (e.g., to represent the relationship  $l_2, k_2$  where there is no corresponding value for  $J$ )



# Comparison of BCNF and 3NF

- Advantages to 3NF over BCNF. It is always possible to obtain a 3NF design without sacrificing losslessness or dependency preservation.
- Disadvantages to 3NF.
  - We may have to use null values to represent some of the possible meaningful relationships among data items.
  - There is the problem of repetition of information.



# Goals of Normalization

- Let  $R$  be a relation scheme with a set  $F$  of functional dependencies.
- Decide whether a relation scheme  $R$  is in “good” form.
- In the case that a relation scheme  $R$  is not in “good” form, need to decompose it into a set of relation scheme  $\{R_1, R_2, \dots, R_n\}$  such that:
  - Each relation scheme is in good form
  - The decomposition is a lossless decomposition
  - Preferably, the decomposition should be dependency preserving.



# How good is BCNF?

- There are database schemas in BCNF that do not seem to be sufficiently normalized
- Consider a relation

*inst\_info (ID, child\_name, phone)*

- where an instructor may have more than one phone and can have multiple children
- Instance of *inst\_info*

<i>ID</i>	<i>child_name</i>	<i>phone</i>
99999	David	512-555-1234
99999	David	512-555-4321
99999	William	512-555-1234
99999	William	512-555-4321

DFF Note – again this is baffling

- The FD is  $(\text{phone}) \rightarrow (\text{ID}, \text{child\_name})$ , which is a weird was to model the intent.
- Modeling it this weird was causes weird behavior.



# How good is BCNF? (Cont.)

- There are no non-trivial functional dependencies and therefore the relation is in BCNF
- Insertion anomalies – i.e., if we add a phone 981-992-3443 to 99999, we need to add two tuples

(99999, David, 981-992-3443)  
(99999, William, 981-992-3443)

DFF Note – again this is baffling

- The FD is  $(\text{phone}) \rightarrow (\text{ID}, \text{child\_name})$ , which is a weird was to model the intent.
- Modeling it this weird was causes weird behavior.



# Higher Normal Forms

- It is better to decompose *inst\_info* into:

- *inst\_child*:

<i>ID</i>	<i>child_name</i>
99999	David
99999	William

- *inst\_phone*:

<i>ID</i>	<i>phone</i>
99999	512-555-1234
99999	512-555-4321

- This suggests the need for higher normal forms, such as Fourth Normal Form (4NF), which we shall see later

DFF Note:

- We are not going to see this later.
- This is the next step on the road to  $\infty$ MF, which is also known as “Normal Form H\*II”



# Functional-Dependency Theory

Oh H\*ell No!  
They can make me cover topics, but  
I have limits!



# Algorithm for Decomposition Using Functional Dependencies

Also known as,  
“I need a topic for my Ph.D. Thesis.”



# Multivalued Dependencies



# Multivalued Dependencies (MVDs)

- Suppose we record names of children, and phone numbers for instructors:
  - $inst\_child(ID, child\_name)$
  - $inst\_phone(ID, phone\_number)$
- If we were to combine these schemas to get
  - $inst\_info(ID, child\_name, phone\_number)$
  - Example data:
    - (99999, David, 512-555-1234)
    - (99999, David, 512-555-4321)
    - (99999, William, 512-555-1234)
    - (99999, William, 512-555-4321)
- This relation is in BCNF
  - Why?



# Multivalued Dependencies

- Let  $R$  be a relation schema and let  $\alpha \subseteq R$  and  $\beta \subseteq R$ . The **multivalued dependency**

$$\alpha \rightarrow\!\!\!\rightarrow \beta$$

holds on  $R$  if in any legal relation  $r(R)$ , for all pairs of tuples  $t_1$  and  $t_2$  in  $r$  such that  $t_1[\alpha] = t_2[\alpha]$ , there exist tuples  $t_3$  and  $t_4$  in  $r$  such that:

$$\begin{aligned}t_1[\alpha] &= t_2[\alpha] = t_3[\alpha] = t_4[\alpha] \\t_3[\beta] &= t_1[\beta] \\t_3[R - \beta] &= t_2[R - \beta] \\t_4[\beta] &= t_2[\beta] \\t_4[R - \beta] &= t_1[R - \beta]\end{aligned}$$

In the real world, if you are thinking about things like this when designing a schema, you need to get help.



# MVD -- Tabular representation

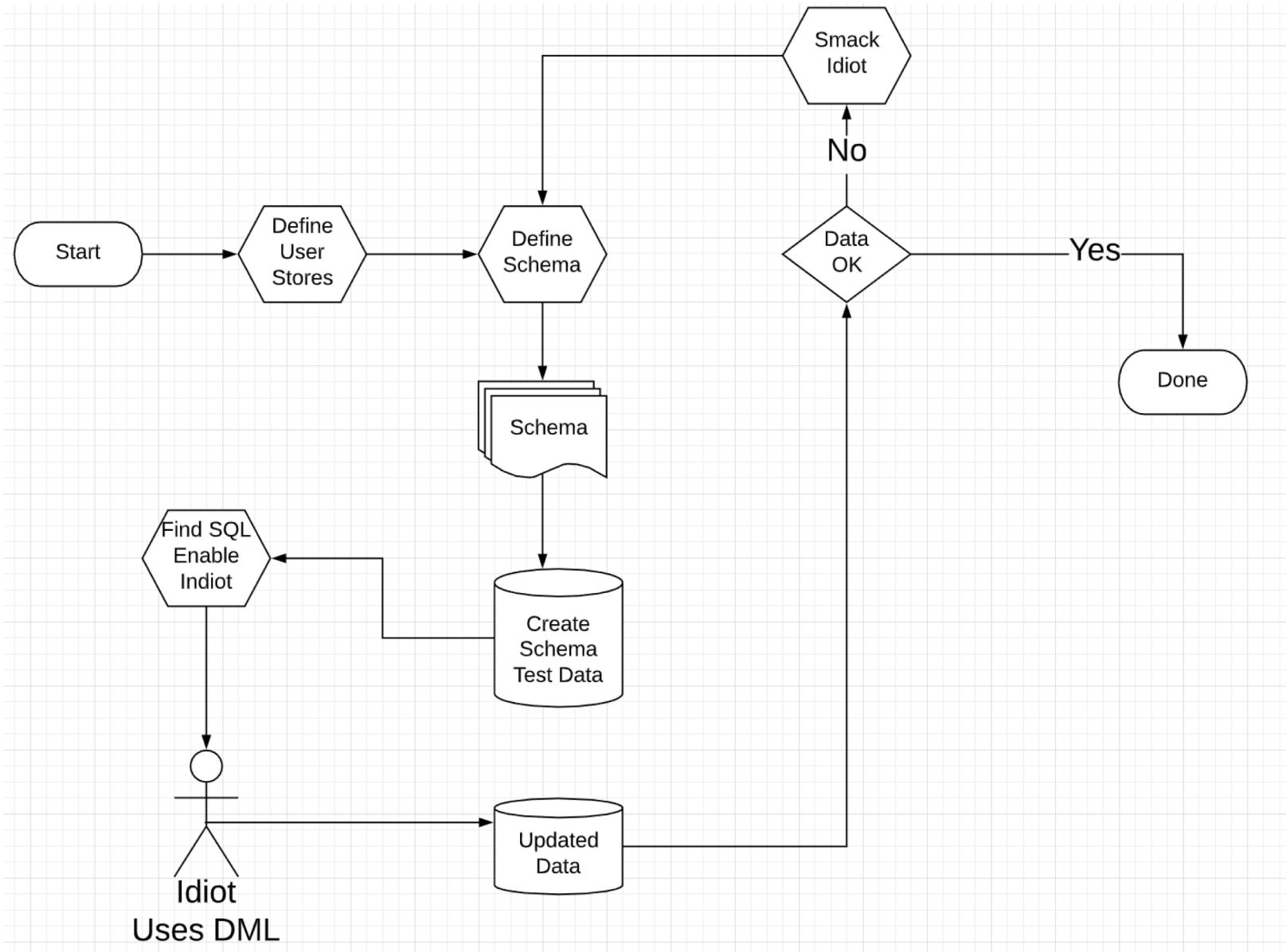
- Tabular representation of  $\alpha \rightarrow\!\!\!\rightarrow \beta$

	$\alpha$	$\beta$	$R - \alpha - \beta$
$t_1$	$a_1 \dots a_i$	$a_{i+1} \dots a_j$	$a_{j+1} \dots a_n$
$t_2$	$a_1 \dots a_i$	$b_{i+1} \dots b_j$	$b_{j+1} \dots b_n$
$t_3$	$a_1 \dots a_i$	$a_{i+1} \dots a_j$	$b_{j+1} \dots b_n$
$t_4$	$a_1 \dots a_i$	$b_{i+1} \dots b_j$	$a_{j+1} \dots a_n$

"Oh, H\*ll no!"  
We're done here.

# **Universal – Ferguson's Law of Useful Normalization and Konstraints (U–FLUNK)**

# U-FLUNK is a Workflow Process





# Further Normal Forms

- **Join dependencies** generalize multivalued dependencies
  - lead to **project-join normal form (PJNF)** (also called **fifth normal form**)
- A class of even more general constraints, leads to a normal form called **domain-key normal form**.
- Problem with these generalized constraints: are hard to reason with, and no set of sound and complete set of inference rules exists.
- Hence rarely used

“Sometimes, I really hate my life.”



# Denormalization for Performance

- May want to use non-normalized schema for performance
- For example, displaying *prereqs* along with *course\_id*, and *title* requires join of *course* with *prereq*
- Alternative 1: Use denormalized relation containing attributes of *course* as well as *prereq* with all above attributes
  - faster lookup
  - extra space and extra execution time for updates
  - extra coding work for programmer and possibility of error in extra code
- Alternative 2: use a materialized view defined a  $course \bowtie prereq$ 
  - Benefits and drawbacks same as above, except no extra coding work for programmer and avoids possible errors



# Other Design Issues

- Some aspects of database design are not caught by normalization
- Examples of bad database design, to be avoided:

Instead of *earnings* (*company\_id*, *year*, *amount*), use

- *earnings\_2004*, *earnings\_2005*, *earnings\_2006*, etc., all on the schema (*company\_id*, *earnings*).
  - ▶ Above are in BCNF, but make querying across years difficult and needs new table each year
- *company\_year* (*company\_id*, *earnings\_2004*, *earnings\_2005*, *earnings\_2006*)
  - ▶ Also in BCNF, but also makes querying across years difficult and requires new attribute each year.
  - ▶ Is an example of a **crosstab**, where values for one attribute become column names
  - ▶ Used in spreadsheets, and in data analysis tools



# Modeling Temporal Data

- **Temporal data** have an association time interval during which the data are *valid*.
- A **snapshot** is the value of the data at a particular point in time
- Several proposals to extend ER model by adding valid time to
  - attributes, e.g., address of an instructor at different points in time
  - entities, e.g., time duration when a student entity exists
  - relationships, e.g., time during which an instructor was associated with a student as an advisor.
- But no accepted standard
- Adding a temporal component results in functional dependencies like
$$ID \rightarrow street, city$$
not holding, because the address varies over time
- A **temporal functional dependency**  $X \rightarrow Y$  holds on schema  $R$  if the functional dependency  $X \rightarrow Y$  holds on all snapshots for all legal instances  $r$  ( $R$ ).



# Modeling Temporal Data (Cont.)

- In practice, database designers may add start and end time attributes to relations
  - E.g.,  $course(course\_id, course\_title)$  is replaced by  
 $course(course\_id, course\_title, start, end)$
  - Constraint: no two tuples can have overlapping valid times
    - ▶ Hard to enforce efficiently
- Foreign key references may be to current version of data, or to data at a point in time
  - E.g., student transcript should refer to course information at the time the course was taken



# **Proof of Correctness of 3NF Decomposition Algorithm**

**Just Kidding.**



# First Normal Form

- Domain is **atomic** if its elements are considered to be indivisible units
  - Examples of non-atomic domains:
    - ▶ Set of names, composite attributes
    - ▶ Identification numbers like CS101 that can be broken up into parts
- A relational schema R is in **first normal form** if the domains of all attributes of R are atomic
- Non-atomic values complicate storage and encourage redundant (repeated) storage of data
  - Example: Set of accounts stored with each customer, and set of owners stored with each account
  - We assume all relations are in first normal form (and revisit this in Chapter 22: Object Based Databases)



# First Normal Form (Cont.)

- Atomicity is actually a property of how the elements of the domain are used.
  - Example: Strings would normally be considered indivisible
  - Suppose that students are given roll numbers which are strings of the form *CS0012* or *EE1127*
  - If the first two characters are extracted to find the department, the domain of roll numbers is not atomic.
  - Doing so is a bad idea: leads to encoding of information in application program rather than in the database.



# Chapter 20: Data Analysis

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Chapter 20: Data Analysis

- Decision Support Systems
- Data Warehousing
- Data Mining
- Classification
- Association Rules
- Clustering



# Decision Support Systems

- **Decision-support systems** are used to make business decisions, often based on data collected by on-line transaction-processing systems.
- Examples of business decisions:
  - What items to stock?
  - What insurance premium to change?
  - To whom to send advertisements?
- Examples of data used for making decisions
  - Retail sales transaction details
  - Customer profiles (income, age, gender, etc.)



# Decision-Support Systems: Overview

- **Data analysis** tasks are simplified by specialized tools and SQL extensions
  - Example tasks
    - ▶ For each product category and each region, what were the total sales in the last quarter and how do they compare with the same quarter last year
    - ▶ As above, for each product category and each customer category
- **Statistical analysis** packages (e.g., : S++) can be interfaced with databases
  - Statistical analysis is a large field, but not covered here
- **Data mining** seeks to discover knowledge automatically in the form of statistical rules and patterns from large databases.
- A **data warehouse** archives information gathered from multiple sources, and stores it under a unified schema, at a single site.
  - Important for large businesses that generate data from multiple divisions, possibly at multiple sites
  - Data may also be purchased externally

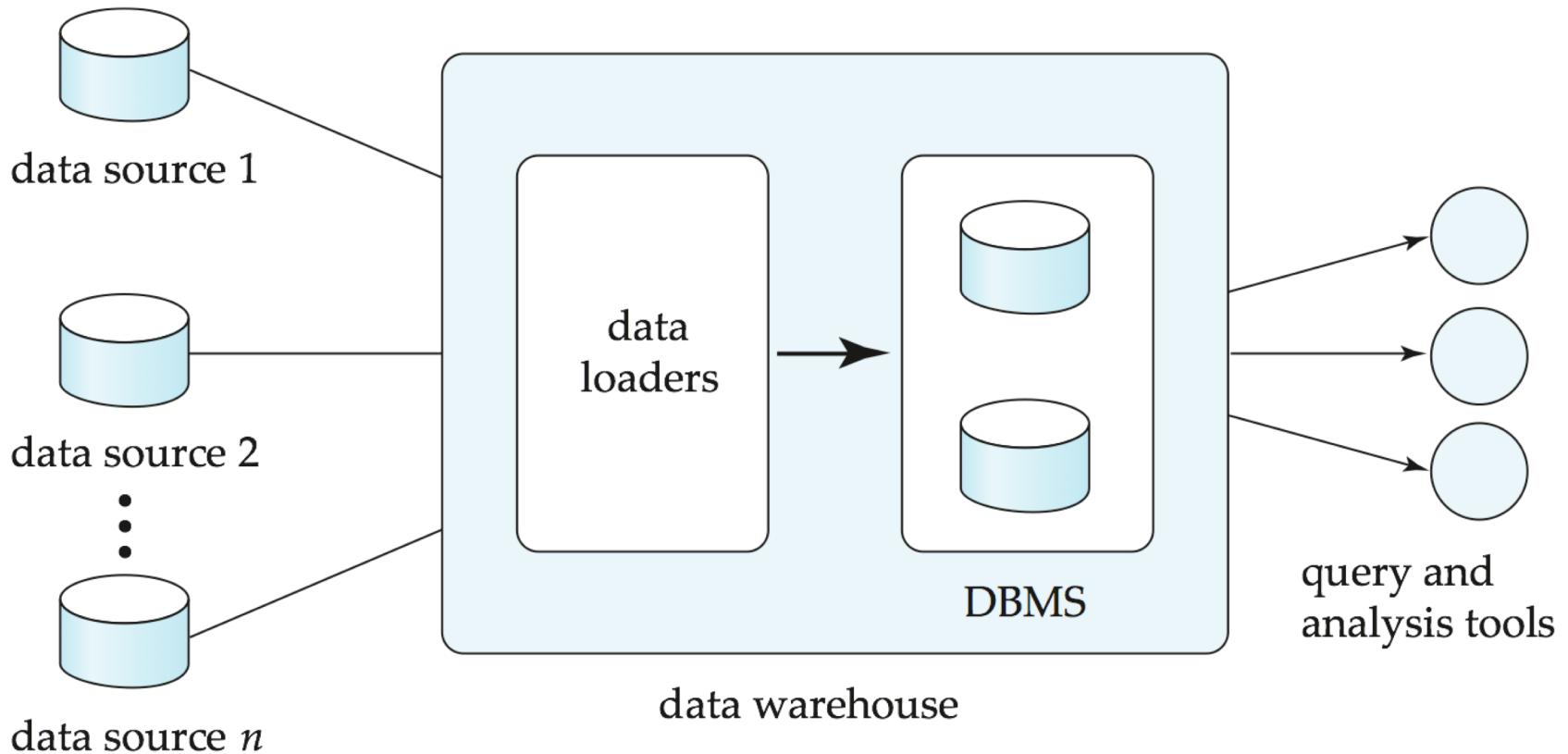


# Data Warehousing

- Data sources often store only current data, not historical data
- Corporate decision making requires a unified view of all organizational data, including historical data
- A **data warehouse** is a repository (archive) of information gathered from multiple sources, stored under a unified schema, at a single site
  - Greatly simplifies querying, permits study of historical trends
  - Shifts decision support query load away from transaction processing systems



# Data Warehousing





# Design Issues

## ■ *When and how to gather data*

- **Source driven architecture**: data sources transmit new information to warehouse, either continuously or periodically (e.g., at night)
- **Destination driven architecture**: warehouse periodically requests new information from data sources
- Keeping warehouse exactly synchronized with data sources (e.g., using two-phase commit) is too expensive
  - ▶ Usually OK to have slightly out-of-date data at warehouse
  - ▶ Data/updates are periodically downloaded from online transaction processing (OLTP) systems.

## ■ *What schema to use*

- Schema integration



# More Warehouse Design Issues

- *Data cleansing*
  - E.g., correct mistakes in addresses (misspellings, zip code errors)
  - **Merge** address lists from different sources and **purge** duplicates
- *How to propagate updates*
  - Warehouse schema may be a (materialized) view of schema from data sources
- *What data to summarize*
  - Raw data may be too large to store on-line
  - Aggregate values (totals/subtotals) often suffice
  - Queries on raw data can often be transformed by query optimizer to use aggregate values

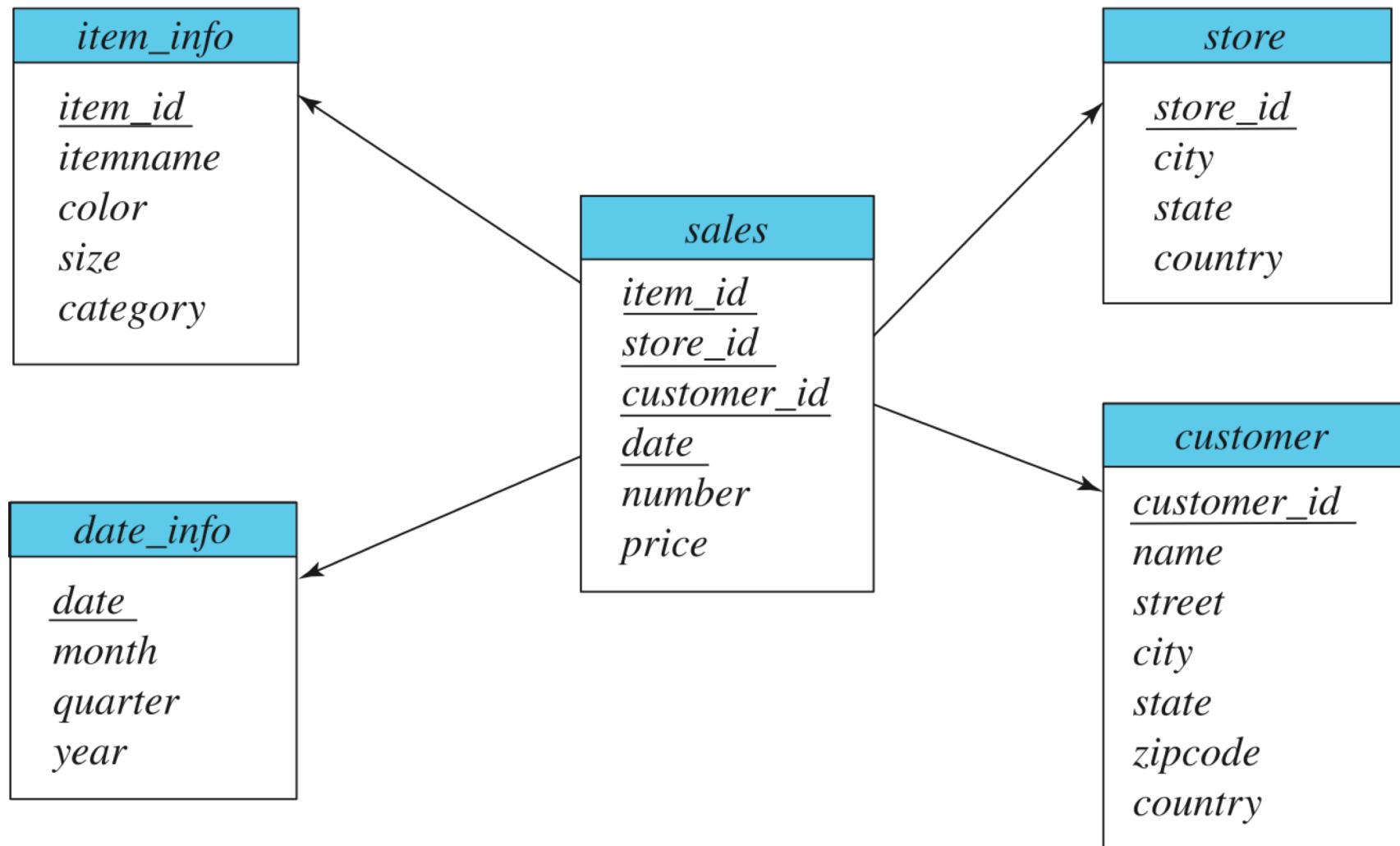


# Warehouse Schemas

- Dimension values are usually encoded using small integers and mapped to full values via dimension tables
- Resultant schema is called a **star schema**
  - More complicated schema structures
    - ▶ **Snowflake schema**: multiple levels of dimension tables
    - ▶ **Constellation**: multiple fact tables



# Data Warehouse Schema





# Data Mining

- Data mining is the process of semi-automatically analyzing large databases to find useful patterns
- **Prediction** based on past history
  - Predict if a credit card applicant poses a good credit risk, based on some attributes (income, job type, age, ...) and past history
  - Predict if a pattern of phone calling card usage is likely to be fraudulent
- Some examples of prediction mechanisms:
  - **Classification**
    - ▶ Given a new item whose class is unknown, predict to which class it belongs
  - **Regression** formulae
    - ▶ Given a set of mappings for an unknown function, predict the function result for a new parameter value



# Data Mining (Cont.)

## ■ Descriptive Patterns

- **Associations**

- ▶ Find books that are often bought by “similar” customers. If a new such customer buys one such book, suggest the others too.
- Associations may be used as a first step in detecting **causation**
  - ▶ E.g., association between exposure to chemical X and cancer,

- **Clusters**

- ▶ E.g., typhoid cases were clustered in an area surrounding a contaminated well
- ▶ Detection of clusters remains important in detecting epidemics

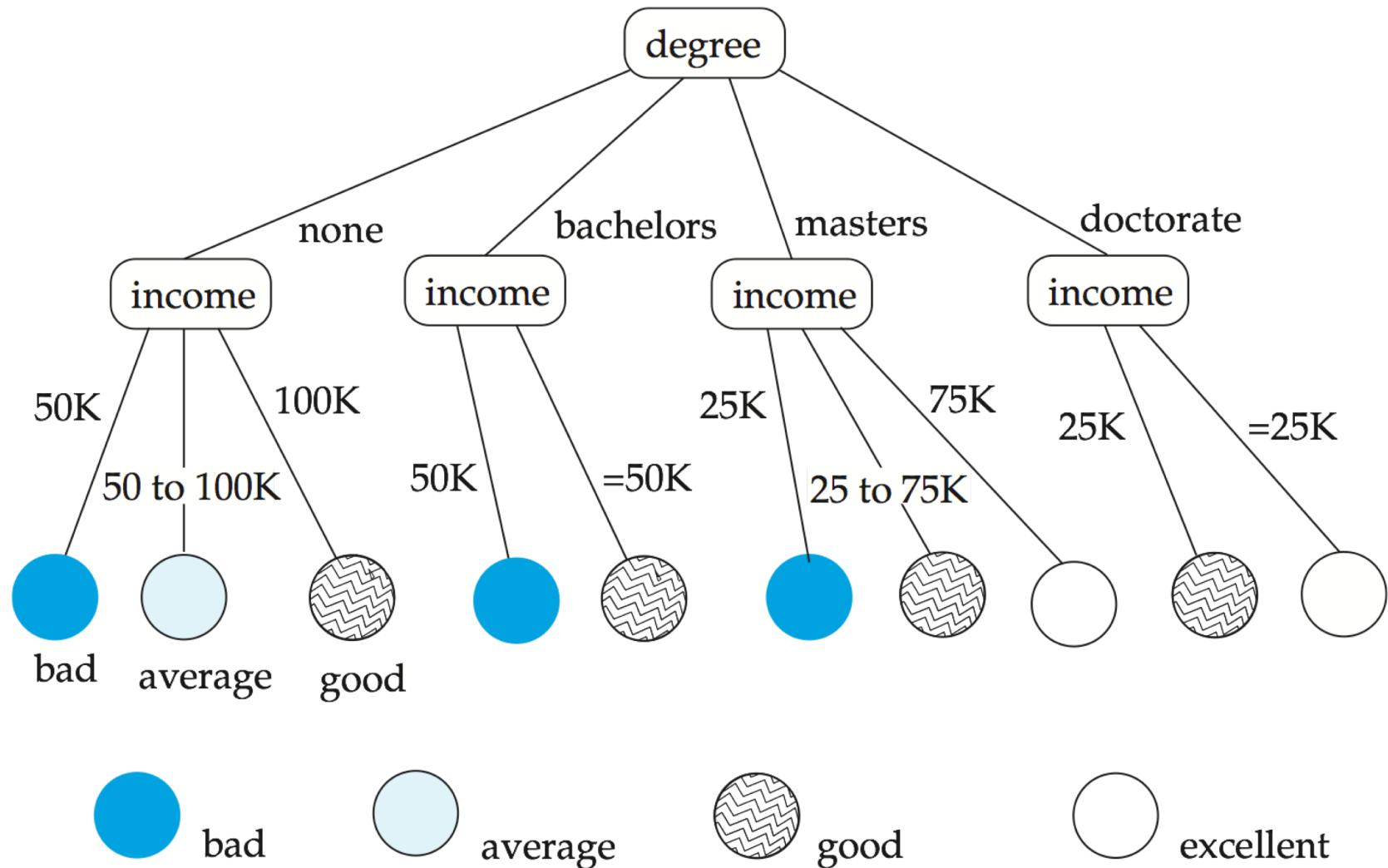


# Classification Rules

- Classification rules help assign new objects to classes.
  - E.g., given a new automobile insurance applicant, should he or she be classified as low risk, medium risk or high risk?
- Classification rules for above example could use a variety of data, such as educational level, salary, age, etc.
  - $\forall \text{ person } P, P.\text{degree} = \text{masters} \text{ and } P.\text{income} > 75,000$   
 $\Rightarrow P.\text{credit} = \text{excellent}$
  - $\forall \text{ person } P, P.\text{degree} = \text{bachelors} \text{ and } (P.\text{income} \geq 25,000 \text{ and } P.\text{income} \leq 75,000)$   
 $\Rightarrow P.\text{credit} = \text{good}$
- Rules are not necessarily exact: there may be some misclassifications
- Classification rules can be shown compactly as a decision tree.



# Decision Tree





# Construction of Decision Trees

- **Training set:** a data sample in which the classification is already known.
- **Greedy** top down generation of decision trees.
  - Each internal node of the tree partitions the data into groups based on a **partitioning attribute**, and a **partitioning condition** for the node
  - **Leaf** node:
    - ▶ all (or most) of the items at the node belong to the same class, or
    - ▶ all attributes have been considered, and no further partitioning is possible.



# Other Types of Classifiers

- Neural net classifiers are studied in artificial intelligence and are not covered here
- Bayesian classifiers use **Bayes theorem**, which says

$$p(c_j | d) = \frac{p(d | c_j) p(c_j)}{p(d)}$$

where

$p(c_j | d)$  = probability of instance  $d$  being in class  $c_j$ ,

$p(d | c_j)$  = probability of generating instance  $d$  given class  $c_j$ ,

$p(c_j)$  = probability of occurrence of class  $c_j$ , and

$p(d)$  = probability of instance  $d$  occurring



# Regression

- Regression deals with the prediction of a value, rather than a class.
  - Given values for a set of variables,  $X_1, X_2, \dots, X_n$ , we wish to predict the value of a variable Y.
- One way is to infer coefficients  $a_0, a_1, a_2, \dots, a_n$  such that
$$Y = a_0 + a_1 * X_1 + a_2 * X_2 + \dots + a_n * X_n$$
- Finding such a linear polynomial is called **linear regression**.
  - In general, the process of finding a curve that fits the data is also called **curve fitting**.
- The fit may only be approximate
  - because of noise in the data, or
  - because the relationship is not exactly a polynomial
- Regression aims to find coefficients that give the best possible fit.



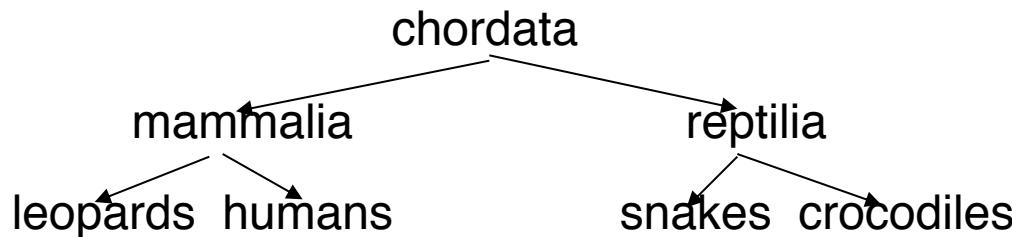
# Clustering

- Clustering: Intuitively, finding clusters of points in the given data such that similar points lie in the same cluster
- Can be formalized using distance metrics in several ways
  - Group points into  $k$  sets (for a given  $k$ ) such that the average distance of points from the centroid of their assigned group is minimized
    - ▶ Centroid: point defined by taking average of coordinates in each dimension.
  - Another metric: minimize average distance between every pair of points in a cluster
- Has been studied extensively in statistics, but on small data sets
  - Data mining systems aim at clustering techniques that can handle very large data sets
  - E.g., the Birch clustering algorithm (more shortly)



# Hierarchical Clustering

- Example from biological classification
  - (the word classification here does not mean a prediction mechanism)



- Other examples: Internet directory systems (e.g., Yahoo, more on this later)
- **Agglomerative clustering algorithms**

- Build small clusters, then cluster small clusters into bigger clusters, and so on

## **Divisive clustering algorithms**

- Start with all items in a single cluster, repeatedly refine (break) clusters into smaller ones



# Other Types of Mining

- **Text mining:** application of data mining to textual documents
  - cluster Web pages to find related pages
  - cluster pages a user has visited to organize their visit history
  - classify Web pages automatically into a Web directory
- **Data visualization** systems help users examine large volumes of data and detect patterns visually
  - Can visually encode large amounts of information on a single screen
  - Humans are very good at detecting visual patterns



# Chapter 19: Distributed Databases

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Chapter 19: Distributed Databases

- Heterogeneous and Homogeneous Databases
- Distributed Data Storage
- Distributed Transactions
- Commit Protocols
- Concurrency Control in Distributed Databases
- Availability
- Distributed Query Processing
- Heterogeneous Distributed Databases
- Directory Systems



# Distributed Database System

- A distributed database system consists of loosely coupled sites that share no physical component
- Database systems that run on each site are independent of each other
- Transactions may access data at one or more sites



# Homogeneous Distributed Databases

- In a homogeneous distributed database
  - All sites have identical software
  - Are aware of each other and agree to cooperate in processing user requests.
  - Each site surrenders part of its autonomy in terms of right to change schemas or software
  - Appears to user as a single system
- In a heterogeneous distributed database
  - Different sites may use different schemas and software
    - ▶ Difference in schema is a major problem for query processing
    - ▶ Difference in software is a major problem for transaction processing
  - Sites may not be aware of each other and may provide only limited facilities for cooperation in transaction processing



# Distributed Data Storage

- Assume relational data model
- Replication
  - System maintains multiple copies of data, stored in different sites, for faster retrieval and fault tolerance.
- Fragmentation
  - Relation is partitioned into several fragments stored in distinct sites
- Replication and fragmentation can be combined
  - Relation is partitioned into several fragments: system maintains several identical replicas of each such fragment.



# Data Replication

- A relation or fragment of a relation is **replicated** if it is stored redundantly in two or more sites.
- **Full replication** of a relation is the case where the relation is stored at all sites.
- Fully redundant databases are those in which every site contains a copy of the entire database.



# Data Replication (Cont.)

## ■ Advantages of Replication

- **Availability:** failure of site containing relation  $r$  does not result in unavailability of  $r$  if replicas exist.
- **Parallelism:** queries on  $r$  may be processed by several nodes in parallel.
- **Reduced data transfer:** relation  $r$  is available locally at each site containing a replica of  $r$ .

## ■ Disadvantages of Replication

- Increased cost of updates: each replica of relation  $r$  must be updated.
- Increased complexity of concurrency control: concurrent updates to distinct replicas may lead to inconsistent data unless special concurrency control mechanisms are implemented.
  - ▶ One solution: choose one copy as **primary copy** and apply concurrency control operations on primary copy



# Data Fragmentation

- Division of relation  $r$  into fragments  $r_1, r_2, \dots, r_n$  which contain sufficient information to reconstruct relation  $r$ .
- **Horizontal fragmentation:** each tuple of  $r$  is assigned to one or more fragments
- **Vertical fragmentation:** the schema for relation  $r$  is split into several smaller schemas
  - All schemas must contain a common candidate key (or superkey) to ensure lossless join property.
  - A special attribute, the tuple-id attribute may be added to each schema to serve as a candidate key.



# Horizontal Fragmentation of *account* Relation

<i>branch_name</i>	<i>account_number</i>	<i>balance</i>
Hillside	A-305	500
Hillside	A-226	336
Hillside	A-155	62

$$account_1 = \sigma_{branch\_name="Hillside"}(account)$$

<i>branch_name</i>	<i>account_number</i>	<i>balance</i>
Valleyview	A-177	205
Valleyview	A-402	10000
Valleyview	A-408	1123
Valleyview	A-639	750

$$account_2 = \sigma_{branch\_name="Valleyview"}(account)$$



## Vertical Fragmentation of *employee\_info* Relation

<i>branch_name</i>	<i>customer_name</i>	<i>tuple_id</i>
Hillside	Lowman	1
Hillside	Camp	2
Valleyview	Camp	3
Valleyview	Kahn	4
Hillside	Kahn	5
Valleyview	Kahn	6
Valleyview	Green	7

$\text{deposit}_1 = \Pi_{\text{branch\_name}, \text{customer\_name}, \text{tuple\_id}}(\text{employee\_info})$

<i>account_number</i>	<i>balance</i>	<i>tuple_id</i>
A-305	500	1
A-226	336	2
A-177	205	3
A-402	10000	4
A-155	62	5
A-408	1123	6
A-639	750	7

$\text{deposit}_2 = \Pi_{\text{account\_number}, \text{balance}, \text{tuple\_id}}(\text{employee\_info})$



# Advantages of Fragmentation

- Horizontal:
  - allows parallel processing on fragments of a relation
  - allows a relation to be split so that tuples are located where they are most frequently accessed
- Vertical:
  - allows tuples to be split so that each part of the tuple is stored where it is most frequently accessed
  - tuple-id attribute allows efficient joining of vertical fragments
  - allows parallel processing on a relation
- Vertical and horizontal fragmentation can be mixed.
  - Fragments may be successively fragmented to an arbitrary depth.



# Data Transparency

- **Data transparency:** Degree to which system user may remain unaware of the details of how and where the data items are stored in a distributed system
- Consider transparency issues in relation to:
  - Fragmentation transparency
  - Replication transparency
  - Location transparency

# (Microservice) Scaling Dimensions

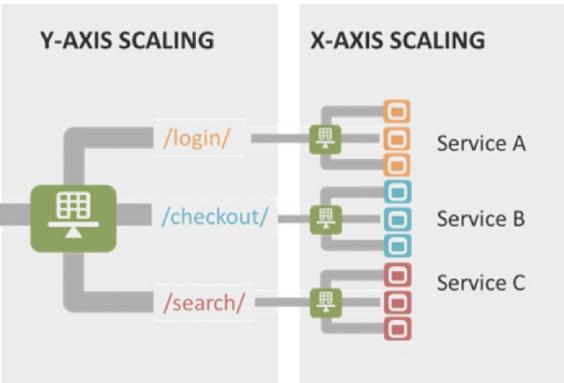
## X-AXIS SCALING

Network name: Horizontal scaling, scale out



## Y-AXIS SCALING

Network name: Layer 7 Load Balancing, Content switching, HTTP Message Steering



## Z-AXIS SCALING

Network name: Layer 7 Load Balancing, Content switching, HTTP Message Steering



- Concept of how to structure applications and large databases.
- The concept also applies when thinking “just about.”
- When you think about it, an RDB sort of has “code” that implements integrity. The code is just declarative.



# Distributed Transactions and 2 Phase Commit

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use

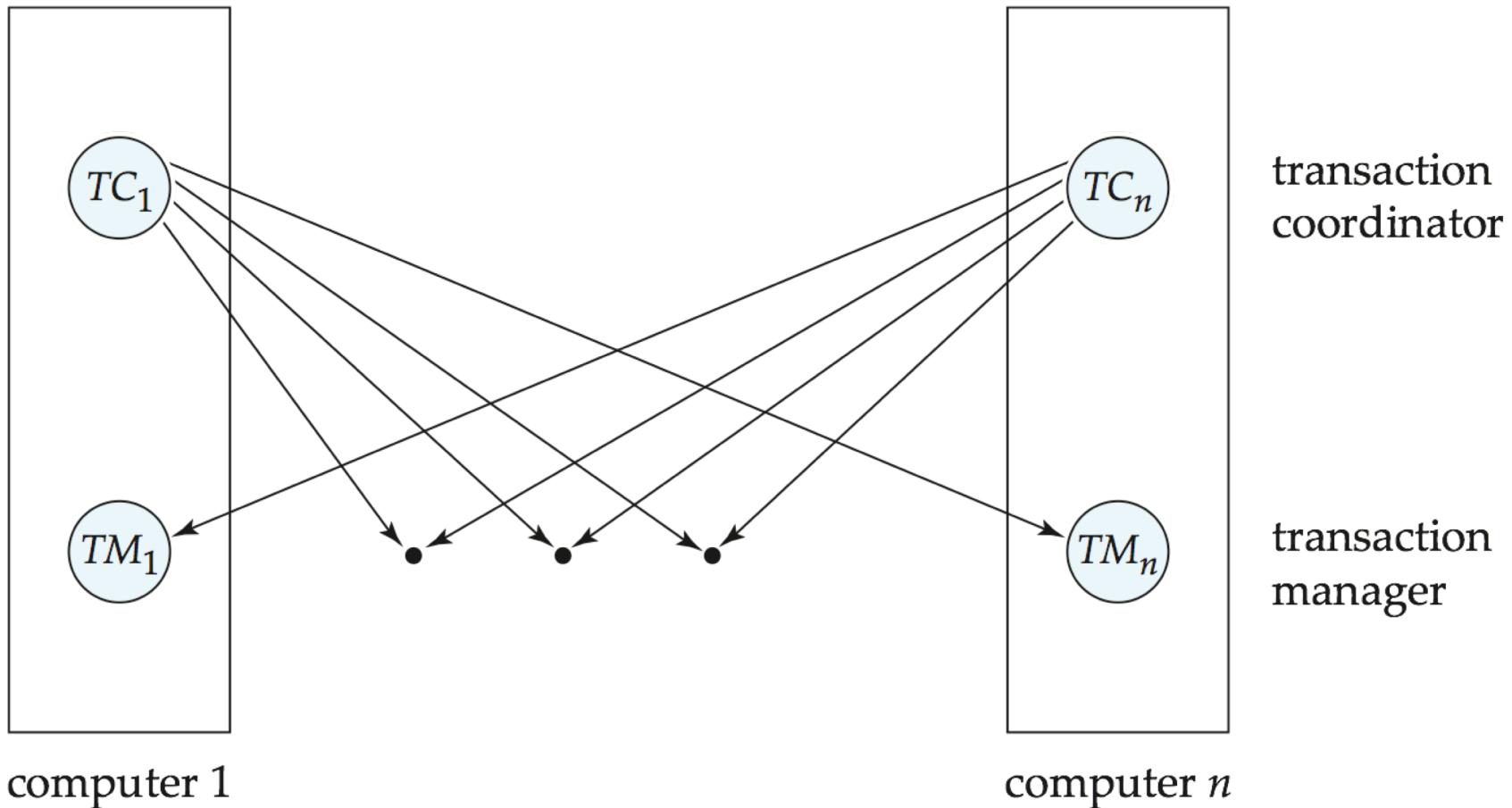


# Distributed Transactions

- Transaction may access data at several sites.
- Each site has a local **transaction manager** responsible for:
  - Maintaining a log for recovery purposes
  - Participating in coordinating the concurrent execution of the transactions executing at that site.
- Each site has a **transaction coordinator**, which is responsible for:
  - Starting the execution of transactions that originate at the site.
  - Distributing subtransactions at appropriate sites for execution.
  - Coordinating the termination of each transaction that originates at the site, which may result in the transaction being committed at all sites or aborted at all sites.



# Transaction System Architecture





# System Failure Modes

- Failures unique to distributed systems:
  - Failure of a site.
  - Loss of messages
    - ▶ Handled by network transmission control protocols such as TCP-IP
  - Failure of a communication link
    - ▶ Handled by network protocols, by routing messages via alternative links
  - **Network partition**
    - ▶ A network is said to be **partitioned** when it has been split into two or more subsystems that lack any connection between them
      - Note: a subsystem may consist of a single node
- Network partitioning and site failures are generally indistinguishable.



# Commit Protocols

- Commit protocols are used to ensure atomicity across sites
  - a transaction which executes at multiple sites must either be committed at all the sites, or aborted at all the sites.
  - not acceptable to have a transaction committed at one site and aborted at another
- The *two-phase commit* (2PC) protocol is widely used
- The *three-phase commit* (3PC) protocol is more complicated and more expensive, but avoids some drawbacks of two-phase commit protocol. This protocol is not used in practice.

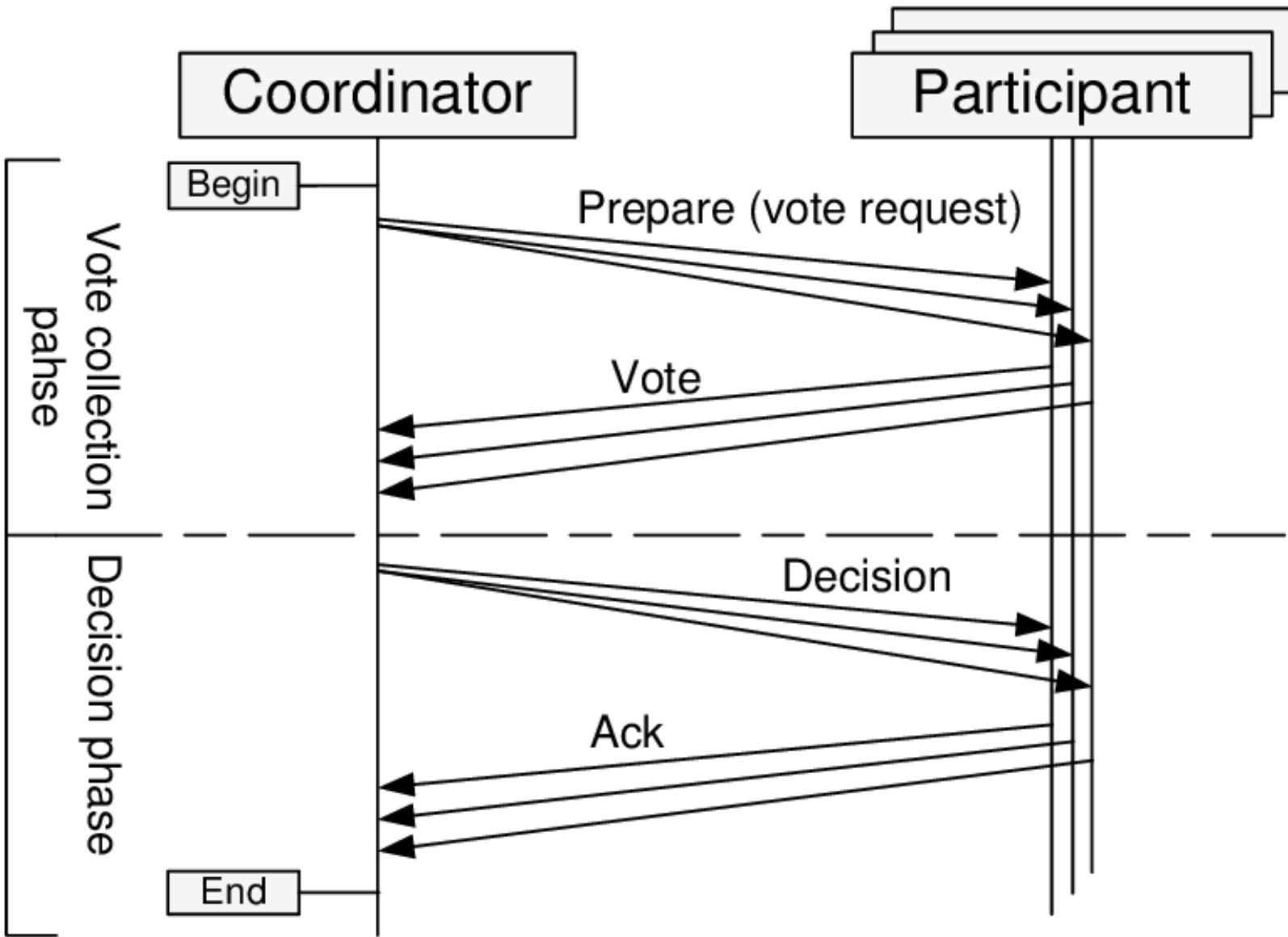


# Commit Protocols

- Commit protocols are used to ensure atomicity across sites
  - a transaction which executes at multiple sites must either be committed at all the sites, or aborted at all the sites.
  - not acceptable to have a transaction committed at one site and aborted at another
- The *two-phase commit* (2PC) protocol is **widely used**
- The *three-phase commit* (3PC) protocol is more complicated and more expensive, but avoids some drawbacks of two-phase commit protocol. This protocol is not used in practice.

**DFF Note:**  
**I am not too sure about this ...**

# Two-Phase Commit





## DFF Note:

- There are a lot of advanced, distributed concurrency control (isolation) architecture.
- Will just discuss a couple.

# Concurrency Control

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Concurrency Control

- Modify concurrency control schemes for use in distributed environment.
- We assume that each site participates in the execution of a commit protocol to ensure global transaction atomicity.
- We assume all replicas of any item are updated
  - Will see how to relax this in case of site failures later



# Single-Lock-Manager Approach

- System maintains a *single* lock manager that resides in a *single* chosen site, say  $S_i$
- When a transaction needs to lock a data item, it sends a lock request to  $S_i$  and lock manager determines whether the lock can be granted immediately
  - If yes, lock manager sends a message to the site which initiated the request
  - If no, request is delayed until it can be granted, at which time a message is sent to the initiating site



# Single-Lock-Manager Approach (Cont.)

- The transaction can read the data item from *any* one of the sites at which a replica of the data item resides.
- Writes must be performed on all replicas of a data item
- Advantages of scheme:
  - Simple implementation
  - Simple deadlock handling
- Disadvantages of scheme are:
  - Bottleneck: lock manager site becomes a bottleneck
  - Vulnerability: system is vulnerable to lock manager site failure.



# Distributed Lock Manager

- In this approach, functionality of locking is implemented by lock managers at each site
  - Lock managers control access to local data items
    - ▶ But special protocols may be used for replicas
- Advantage: work is distributed and can be made robust to failures
- Disadvantage: deadlock detection is more complicated
  - Lock managers cooperate for deadlock detection
    - ▶ More on this later
- Several variants of this approach
  - Primary copy
  - Majority protocol
  - Biased protocol
  - Quorum consensus



# Primary Copy

- Choose one replica of data item to be the **primary copy**.
  - Site containing the replica is called the **primary site** for that data item
  - Different data items can have different primary sites
- When a transaction needs to lock a data item  $Q$ , it requests a lock at the primary site of  $Q$ .
  - Implicitly gets lock on all replicas of the data item
- Benefit
  - Concurrency control for replicated data handled similarly to unreplicated data - simple implementation.
- Drawback
  - If the primary site of  $Q$  fails,  $Q$  is inaccessible even though other sites containing a replica may be accessible.



# Majority Protocol

- Local lock manager at each site administers lock and unlock requests for data items stored at that site.
- When a transaction wishes to lock an unreplicated data item  $Q$  residing at site  $S_i$ , a message is sent to  $S_i$  's lock manager.
  - If  $Q$  is locked in an incompatible mode, then the request is delayed until it can be granted.
  - When the lock request can be granted, the lock manager sends a message back to the initiator indicating that the lock request has been granted.



# Majority Protocol (Cont.)

## ■ In case of replicated data

- If  $Q$  is replicated at  $n$  sites, then a lock request message must be sent to more than half of the  $n$  sites in which  $Q$  is stored.
- The transaction does not operate on  $Q$  until it has obtained a lock on a majority of the replicas of  $Q$ .
- When writing the data item, transaction performs writes on *all* replicas.

## ■ Benefit

- Can be used even when some sites are unavailable
  - ▶ details on how handle writes in the presence of site failure later

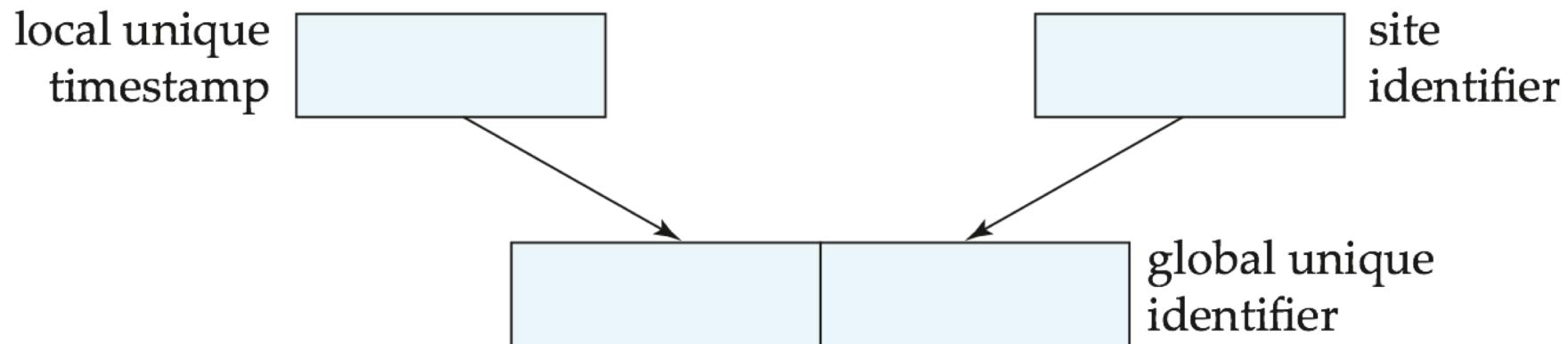
## ■ Drawback

- Requires  $2(n/2 + 1)$  messages for handling lock requests, and  $(n/2 + 1)$  messages for handling unlock requests.
- Potential for deadlock even with single item - e.g., each of 3 transactions may have locks on 1/3rd of the replicas of a data.



# Timestamping

- Timestamp based concurrency-control protocols can be used in distributed systems
- Each transaction must be given a unique timestamp
- Main problem: how to generate a timestamp in a distributed fashion
  - Each site generates a unique local timestamp using either a logical counter or the local clock.
  - Global unique timestamp is obtained by concatenating the unique local timestamp with the unique identifier.





# Timestamping (Cont.)

- A site with a slow clock will assign smaller timestamps
  - Still logically correct: serializability not affected
  - But: “disadvantages” transactions
- To fix this problem
  - Define within each site  $S_i$ , a **logical clock** ( $LC_i$ ), which generates the unique local timestamp
  - Require that  $S_i$  advance its logical clock whenever a request is received from a transaction  $T_i$  with timestamp  $\langle x, y \rangle$  and  $x$  is greater than the current value of  $LC_i$ .
  - In this case, site  $S_i$  advances its logical clock to the value  $x + 1$ .



# Replication with Weak Consistency

- Many commercial databases support replication of data with weak degrees of consistency (i.e., without a guarantee of serializability)
- E.g.: **master-slave replication**: updates are performed at a single “master” site, and propagated to “slave” sites.
  - Propagation is not part of the update transaction: it is decoupled
    - ▶ May be immediately after transaction commits
    - ▶ May be periodic
  - Data may only be read at slave sites, not updated
    - ▶ No need to obtain locks at any remote site
  - Particularly useful for distributing information
    - ▶ E.g. from central office to branch-office
  - Also useful for running read-only queries offline from the main database



# Replication with Weak Consistency (Cont.)

- Replicas should see a **transaction-consistent snapshot** of the database
  - That is, a state of the database reflecting all effects of all transactions up to some point in the serialization order, and no effects of any later transactions.
- E.g. Oracle provides a **create snapshot** statement to create a snapshot of a relation or a set of relations at a remote site
  - snapshot refresh either by recomputation or by incremental update
  - Automatic refresh (continuous or periodic) or manual refresh



# Multimaster and Lazy Replication

- With multimaster replication (also called update-anywhere replication) updates are permitted at any replica, and are automatically propagated to all replicas
  - Basic model in distributed databases, where transactions are unaware of the details of replication, and database system propagates updates as part of the same transaction
    - ▶ Coupled with 2 phase commit
- Many systems support **lazy propagation** where updates are transmitted after transaction commits
  - Allows updates to occur even if some sites are disconnected from the network, but at the cost of consistency



# Trading Consistency for Availability

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# What is Consistency?

- Consistency in Databases (ACID):
  - Database has a set of integrity constraints
  - A consistent database state is one where all integrity constraints are satisfied
  - Each transaction run individually on a consistent database state must leave the database in a consistent state
- Consistency in distributed systems with replication
  - Strong consistency: a schedule with read and write operations on a replicated object should give results and final state equivalent to some schedule on a single copy of the object, with order of operations from a single site preserved
  - Weak consistency (several forms)



# Availability

- Traditionally, availability of centralized server
- For distributed systems, availability of system to process requests
  - For large system, at almost any point in time there's a good chance that
    - ▶ a node is down or even
    - ▶ Network partitioning
- Distributed consensus algorithms will block during partitions to ensure consistency
  - Many applications require continued operation even during a network partition
    - ▶ Even at cost of consistency



# Brewer's CAP Theorem

- Three properties of a system
  - Consistency (all copies have same value)
  - Availability (system can run even if parts have failed)
    - Via replication
  - Partitions (network can break into two or more parts, each with active systems that can't talk to other parts)
- Brewer's CAP “Theorem”: You can have at most two of these three properties for any system
- Very large systems will partition at some point  
→ Choose one of consistency or availability
  - Traditional database choose consistency
  - Most Web applications choose availability
    - Except for specific parts such as order processing



# Eventual Consistency

- When no updates occur for a long period of time, eventually all updates will propagate through the system and all the nodes will be consistent
- For a given accepted update and a given node, eventually either the update reaches the node or the node is removed from service
- Known as **BASE** (**B**asically **A**vailable, **S**oft state, **E**ventual consistency), as opposed to ACID
  - **Soft state**: copies of a data item may be inconsistent
  - **Eventually Consistent** – copies becomes consistent at some later time if there are no more updates to that data item



# Availability vs Latency

- CAP theorem only matters when there is a partition
  - Even if partitions are rare, applications may trade off consistency for latency
    - ▶ E.g. PNUTS allows inconsistent reads to reduce latency
      - Critical for many applications
    - ▶ But update protocol (via master) ensures consistency over availability
  - Thus there are two questions :
    - ▶ If there is partitioning, how does system tradeoff *availability* for *consistency*
    - ▶ else how does system trade off *latency* for *consistency*



# Cloud Databases

**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Data Storage on the Cloud

- Need to store and retrieve massive amounts of data
- Traditional parallel databases not designed to scale to 1000's of nodes (and expensive)
- Initial needs did not include full database functionality
  - Store and retrieve data items by key value is minimum functionality
    - ▶ **Key-value stores**
- Several implementations
  - Bigtable from Google,
  - HBase, an open source clone of Bigtable
  - Dynamo, which is a key-value storage system from Amazon
  - Cassandra, from FaceBook
  - Sherpa/PNUTS from Yahoo!