

# **Chapter 7: Normalization**

# Overview of Normalization

# Features of Good Relational Designs

- 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)

# A Combined Schema Without Repetition

Not all combined schemas result in repetition of information

- Consider combining relations
  - *sec\_class(sec\_id, building, room\_number)* and
  - *section(course\_id, sec\_id, semester, year)*into one relation
  - *section(course\_id, sec\_id, semester, year, building, room\_number)*
- No repetition in this case

# 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

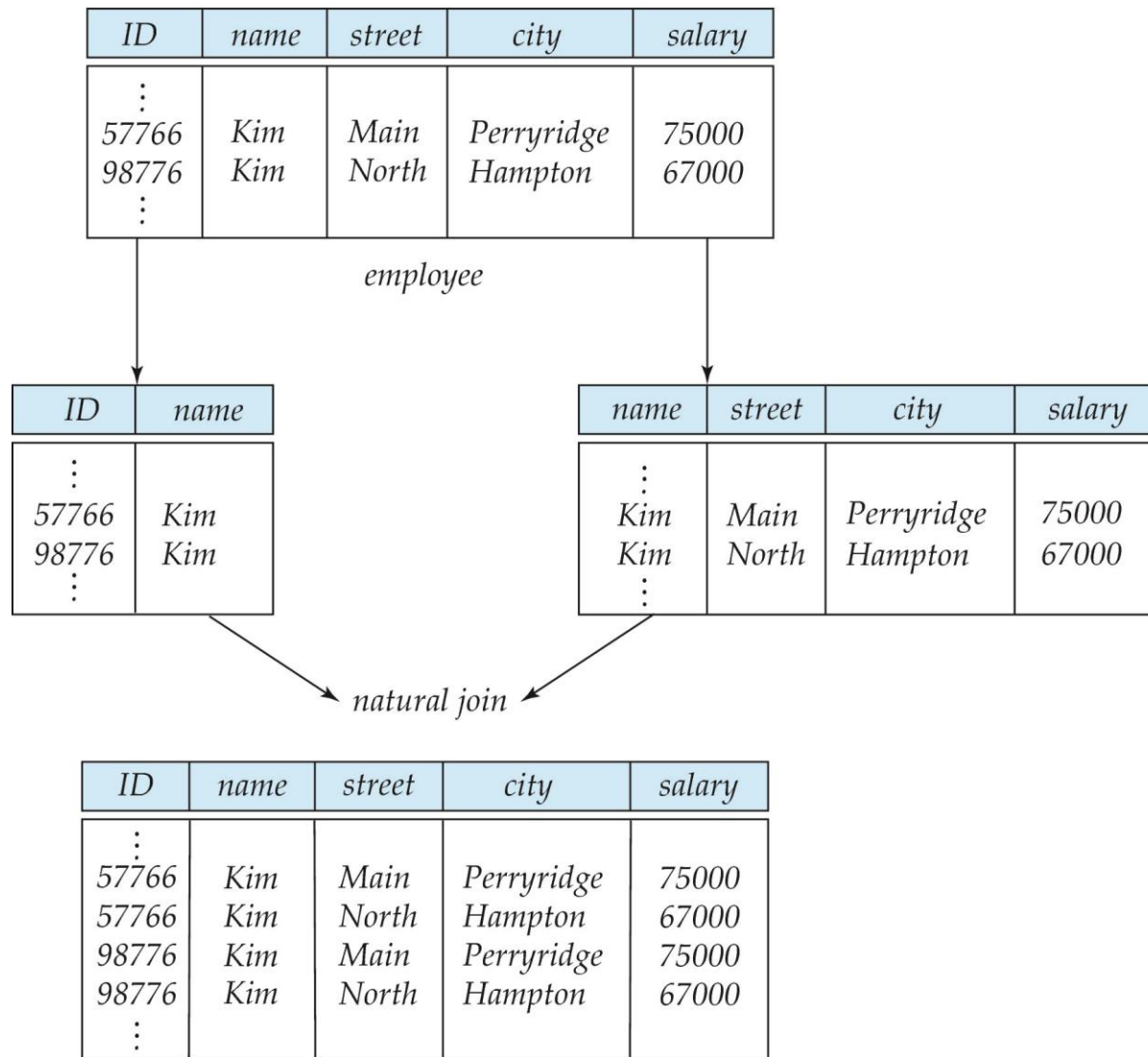
*employee1* (*ID*, *name*)

*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**.

# 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, if we project  $r$  onto  $R_1$  and  $R_2$ , and then compute the natural join of the projection results, we get back exactly  $r$ .

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

- And conversely a decomposition is lossy if the natural join of the projection results a proper superset of  $r$ .

$$r \subset \Pi_{R_1}(r) \bowtie \Pi_{R_2}(r)$$

# Example of Lossless Decomposition

- Decomposition of  $R = (A, B, C)$

$$R_1 = (A, B) \quad R_2 = (B, C)$$

A	B	C
$\alpha$	1	A
$\beta$	2	B

$r$

A	B
$\alpha$	1
$\beta$	2

$\Pi_{A,B}(r)$

B	C
1	A
2	B

$\Pi_{B,C}(r)$

$\Pi_A(r) \bowtie \Pi_B(r)$

A	B	C
$\alpha$	1	A
$\beta$	2	B



# 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  $r$ .

A	B
$a_1$	$b_1$
$a_1$	$b_1$
$a_2$	$b_2$
$a_2$	$b_2$
$a_3$	$b_2$

On this instance,  $A \rightarrow B$  hold;  $B \rightarrow A$  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*

# 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 decomposition are lossless.
- For the case of  $R = (R_1, R_2)$ , we require that for all possible relations  $r$  on schema  $R$

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

- A decomposition of  $R$  into  $R_1$  and  $R_2$  is lossless decomposition if at least one of the following dependencies is in  $F^+$ :
  - $R_1 \cap R_2 \rightarrow R_1$
  - $R_1 \cap R_2 \rightarrow R_2$
- Example

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

*instructor (ID, name, dept\_name, salary)*

*department (dept\_name, building, budget)*

***dept\_name → dept\_name, building, budget***

# Example

- $R = (A, B, C)$   
 $F = \{A \rightarrow B, B \rightarrow C\}$
- $R_1 = (A, B), \quad R_2 = (B, C)$ 
  - Lossless decomposition:  
 $R_1 \cap R_2 = \{B\}$  and  $B \rightarrow BC$
- $R_1 = (A, B), \quad R_2 = (A, C)$ 
  - Lossless decomposition:  
 $R_1 \cap R_2 = \{A\}$  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

- An instructor can be associated with only a single department, and a student may have more than one advisor, but no more than one from a given department.

- Consider a schema:

*dept\_advisor(s\_ID, i\_ID, department\_name)*

- With function dependencies:

$i\_ID \rightarrow dept\_name$

$s\_ID, dept\_name \rightarrow i\_ID$

- In the above design we are forced to repeat the department name once for each time an instructor participates in a *dept\_advisor* relationship.
- To fix this, we need to decompose *dept\_advisor*
- Any decomposition will not include all the attributes in

$s\_ID, dept\_name \rightarrow i\_ID$

- Thus, the composition NOT be dependency preserving

# 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  $R$

# Boyce-Codd Normal Form (Cont.)

- Example schema that is **not** in BCNF:

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

because :

- *dept\_name* → *building, budget*
    - holds on *in\_dep*
    - but
  - *dept\_name* is not a superkey
- When decompose *in\_dept* into *instructor* and *department*
    - *instructor* is in BCNF
    - *department* is in BCNF

# 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
  - $(\alpha \cup \beta) = (dept\_name, building, budget)$
  - $(R - (\beta - \alpha)) = (ID, name, dept\_name, salary)$



# Example

- $R = (A, B, C)$   
 $F = \{A \rightarrow B, B \rightarrow C\}$
- $R_1 = (A, B), R_2 = (B, C)$ 
  - Lossless-join decomposition:  
$$R_1 \cap R_2 = \{B\} \text{ and } B \rightarrow BC$$
  - Dependency preserving
- $R_1 = (A, B), R_2 = (A, C)$ 
  - Lossless-join decomposition:  
$$R_1 \cap R_2 = \{A\} \text{ and } A \rightarrow AB$$
  - Not dependency preserving  
(cannot check  $B \rightarrow C$  without computing  $R_1 \bowtie R_2$ )

# BCNF and Dependency Preservation

- It is not always possible to achieve both BCNF and dependency preservation

- Consider a schema:

*dept\_advisor(s\_ID, i\_ID, department\_name)*

- With function dependencies:

$i\_ID \rightarrow dept\_name$

$s\_ID, dept\_name \rightarrow i\_ID$

- *dept\_advisor* is not in BCNF

- $i\_ID$  is not a superkey.

- Any decomposition of *dept\_advisor* will not include all the attributes in

$s\_ID, dept\_name \rightarrow i\_ID$

- 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:

*dept\_advisor(s\_ID, i\_ID, dept\_name)*

- With function dependencies:

$i\_ID \rightarrow dept\_name$

$s\_ID, dept\_name \rightarrow i\_ID$

- Two candidate keys =  $\{s\_ID, dept\_name\}, \{s\_ID, i\_ID\}$
- We have seen before that *dept\_advisor* is **not** in BCNF
- *R*, however, is in 3NF
  - $s\_ID, dept\_name$  is a superkey
  - $i\_ID \rightarrow dept\_name$  and  $i\_ID$  is NOT a superkey, but:
    - $\{dept\_name\} - \{i\_ID\} = \{dept\_name\}$  and
    - $dept\_name$  is contained in a candidate key

# 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

## 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)

# 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



# **Functional-Dependency Theory**

# 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  $F^+$ .

# Closure of a Set of Functional Dependencies

- We can compute  $F^+$ , the closure of  $F$ , by repeatedly applying **Armstrong's Axioms**:
  - **Reflexive rule**: if  $\beta \subseteq \alpha$ , then  $\alpha \rightarrow \beta$
  - **Augmentation rule**: if  $\alpha \rightarrow \beta$ , then  $\gamma \alpha \rightarrow \gamma \beta$
  - **Transitivity rule**: if  $\alpha \rightarrow \beta$ , and  $\beta \rightarrow \gamma$ , then  $\alpha \rightarrow \gamma$
- These rules are
  - **Sound** -- generate only functional dependencies that actually hold, and
  - **Complete** -- generate all functional dependencies that hold.

# Example of $F^+$

- $R = (A, B, C, G, H, I)$

$F = \{$   
     $A \rightarrow B$   
     $A \rightarrow C$   
     $CG \rightarrow H$   
     $CG \rightarrow I$   
     $B \rightarrow H\}$

- Some members of  $F^+$

- $A \rightarrow H$

- by transitivity from  $A \rightarrow B$  and  $B \rightarrow H$

- $AG \rightarrow I$

- by augmenting  $A \rightarrow C$  with  $G$ , to get  $AG \rightarrow CG$   
and then transitivity with  $CG \rightarrow I$

- $CG \rightarrow HI$

- by augmenting  $CG \rightarrow I$  to infer  $CG \rightarrow CGI$ ,  
and augmenting of  $CG \rightarrow H$  to infer  $CGI \rightarrow HI$ ,  
and then transitivity

# Closure of Functional Dependencies (Cont.)

- Additional rules:
  - **Union rule:** If  $\alpha \rightarrow \beta$  holds and  $\alpha \rightarrow \gamma$  holds, then  $\alpha \rightarrow \beta\gamma$  holds.
  - **Decomposition rule:** If  $\alpha \rightarrow \beta\gamma$  holds, then  $\alpha \rightarrow \beta$  holds and  $\alpha \rightarrow \gamma$  holds.
  - **Pseudotransitivity rule:** If  $\alpha \rightarrow \beta$  holds and  $\gamma\beta \rightarrow \delta$  holds, then  $\alpha\gamma \rightarrow \delta$  holds.
- The above rules can be inferred from Armstrong's axioms.

# Procedure for Computing $F^+$

- To compute the closure of a set of functional dependencies  $F$ :

$F^+ = F$

**repeat**

**for each** functional dependency  $f$  in  $F^+$

        apply reflexivity and augmentation rules on  $f$

        add the resulting functional dependencies to  $F^+$

**for each** pair of functional dependencies  $f_1$  and  $f_2$  in  $F^+$

**if**  $f_1$  and  $f_2$  can be combined using transitivity

**then** add the resulting functional dependency to  $F^+$

**until**  $F^+$  does not change any further

- **NOTE:** We shall see an alternative procedure for this task later