

# Chapter Outline

## Functional Dependencies (FDs)

- Definition of FD
- Inference Rules for FDs
- Equivalence of Sets of FDs
- Minimal Sets of FDs

# Chapter Outline

## Normal Forms Based on Primary Keys

- Normalization of Relations
  - Practical Use of Normal Forms
  - Definitions of Keys and Attributes Participating in Keys
  - First Normal Form
  - Second Normal Form
  - Third Normal Form
- General Normal Form Definitions (For Multiple Keys)
  - BCNF (Boyce-Codd Normal Form)

# Functional Dependencies

- Functional dependencies (FDs)
  - Are used to specify *formal measures* of the "goodness" of relational designs
  - And keys are used to define **normal forms** for relations
  - Are **constraints** that are derived from the *meaning* and *interrelationships* of the data attributes
- A set of attributes *X functionally determines* a set of attributes *Y* if the value of *X* determines a unique value for *Y*

# Functional Dependencies

- $X \rightarrow Y$  holds if whenever two tuples have the same value for  $X$ , they *must have* the same value for  $Y$ 
  - For any two tuples  $t_1$  and  $t_2$  in any relation instance  $r(R)$ : If  $t_1[X]=t_2[X]$ , *then*  $t_1[Y]=t_2[Y]$
- $X \rightarrow Y$  in  $R$  specifies a *constraint* on all relation instances  $r(R)$
- Written as  $X \rightarrow Y$ ; can be displayed graphically on a relation schema as in Figures. ( denoted by the arrow: ).
- FDs are derived from the real-world constraints on the attributes

# Examples of FD constraints

- Social security number determines employee name
  - $SSN \rightarrow ENAME$
- Project number determines project name and location
  - $PNUMBER \rightarrow \{PNAME, PLOCATION\}$
- Employee ssn and project number determines the hours per week that the employee works on the project
  - $\{SSN, PNUMBER\} \rightarrow HOURS$

# Examples of FD constraints

- An FD is a property of the attributes in the schema  $R$
- The constraint must hold on *every* relation instance  $r(R)$
- If  $K$  is a key of  $R$ , then  $K$  functionally determines all attributes in  $R$ 
  - (since we never have two distinct tuples with  $t_1[K]=t_2[K]$ )

# Inference Rules for FDs

- Given a set of FDs  $F$ , we can **infer** additional FDs that hold whenever the FDs in  $F$  hold
- Armstrong's inference rules:
  - IR1. (**Reflexive**) If  $Y \text{ subset-of } X$ , then  $X \rightarrow Y$
  - IR2. (**Augmentation**) If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$ 
    - (Notation:  $XZ$  stands for  $X \cup Z$ )
  - IR3. (**Transitive**) If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$
- IR1, IR2, IR3 form a **sound** and **complete** set of inference rules
  - These are rules hold and all other rules that hold can be deduced from these

# Inference Rules for FDs

- Some additional inference rules that are useful:
  - **Decomposition:** If  $X \rightarrow YZ$ , then  $X \rightarrow Y$  and  $X \rightarrow Z$
  - **Union:** If  $X \rightarrow Y$  and  $X \rightarrow Z$ , then  $X \rightarrow YZ$
  - **Pseudotransitivity:** If  $X \rightarrow Y$  and  $WY \rightarrow Z$ , then  $WX \rightarrow Z$
- The last three inference rules, as well as any other inference rules, can be deduced from IR1, IR2, and IR3 (completeness property)



# Inference Rules for FDs

- **Closure** of a set  $F$  of FDs is the set  $F^+$  of all FDs that can be inferred from  $F$
- **Closure** of a set of attributes  $X$  with respect to  $F$  is the set  $X^+$  of all attributes that are functionally determined by  $X$
- $X^+$  can be calculated by repeatedly applying IR1, IR2, IR3 using the FDs in  $F$

# Equivalence of Sets of FDs

- Two sets of FDs  $F$  and  $G$  are **equivalent** if:
  - Every FD in  $F$  can be inferred from  $G$ , and
  - Every FD in  $G$  can be inferred from  $F$
  - Hence,  $F$  and  $G$  are equivalent if  $F^+ = G^+$
- Definition (**Covers**):
  - $F$  **covers**  $G$  if every FD in  $G$  can be inferred from  $F$ 
    - (i.e., if  $G^+ \text{ subset-of } F^+$ )
- $F$  and  $G$  are equivalent if  $F$  covers  $G$  and  $G$  covers  $F$
- There is an algorithm for checking equivalence of sets of FDs

# Equivalence of Sets of FDs

- Two sets of FDs  $F$  and  $G$  are **equivalent** if:
  - Every FD in  $F$  can be inferred from  $G$ , and
  - Every FD in  $G$  can be inferred from  $F$
  - Hence,  $F$  and  $G$  are equivalent if  $F^+ = G^+$
- Definition (**Covers**):
  - $F$  **covers**  $G$  if every FD in  $G$  can be inferred from  $F$ 
    - (i.e., if  $G^+ \text{ subset-of } F^+$ )
- $F$  and  $G$  are equivalent if  $F$  covers  $G$  and  $G$  covers  $F$
- There is an algorithm for checking equivalence of sets of FDs

# Equivalence of Sets of FDs

**$R(A,B,C,D)$**

**$F = \{A \rightarrow B, AB \rightarrow C, D \rightarrow AC, D \rightarrow E\}$  and  $G = \{A \rightarrow BC, D \rightarrow AE\}$**

**Is F and G are Equivalent?**

# Minimal Sets of FDs

- A set of FDs is **minimal** if it satisfies the following conditions:
  1. Every dependency in  $F$  has a single attribute for its RHS.
  2. We cannot remove any dependency from  $F$  and have a set of dependencies that is equivalent to  $F$ .
  3. We cannot replace any dependency  $X \rightarrow A$  in  $F$  with a dependency  $Y \rightarrow A$ , where  $Y$  proper-subset-of  $X$  ( $Y$  subset-of  $X$ ) and still have a set of dependencies that is equivalent to  $F$ .

# Minimal Sets of FDs

- Every set of FDs has an equivalent minimal set
- There can be several equivalent minimal sets
- There is no simple algorithm for computing a minimal set of FDs that is equivalent to a set  $F$  of FDs
- To synthesize a set of relations, we assume that we start with a set of dependencies that is a minimal set

# Minimal Sets of FDs

$F = (A \rightarrow B, C \rightarrow B, D \rightarrow ABC, AC \rightarrow D)$  Find Minimal set.

# Normalization of Relations

- **Normalization:**
  - The process of decomposing unsatisfactory "bad" relations by breaking up their attributes into smaller relations
- **Normal form:**
  - Condition using keys and FDs of a relation to certify whether a relation schema is in a particular normal form



# Normalization of Relations

- 2NF, 3NF, BCNF
  - based on keys and FDs of a relation schema
- 4NF
  - based on keys, multi-valued dependencies :  
MVDs; 5NF based on keys, join dependencies :  
JDs (Chapter 11)
- Additional properties may be needed to ensure a good relational design (lossless join, dependency preservation; Chapter 11)

# Practical Use of Normal Forms

- **Normalization** is carried out in practice so that the resulting designs are of high quality and meet the desirable properties
- The practical utility of these normal forms becomes questionable when the constraints on which they are based are *hard to understand* or to *detect*
- The database designers *need not* normalize to the highest possible normal form
  - (usually up to 3NF, BCNF or 4NF)
- **Denormalization:**
  - The process of storing the join of higher normal form relations as a base relation—which is in a lower normal form

# Definitions of Keys and Attributes Participating in Keys

- A **superkey** of a relation schema  $R = \{A_1, A_2, \dots, A_n\}$  is a set of attributes  $S$  *subset-of*  $R$  with the property that no two tuples  $t_1$  and  $t_2$  in any legal relation state  $r$  of  $R$  will have  $t_1[S] = t_2[S]$
- A **key**  $K$  is a **superkey** with the *additional property* that removal of any attribute from  $K$  will cause  $K$  not to be a superkey any more.

# Definitions of Keys and Attributes Participating in Keys

- If a relation schema has more than one key, each is called a **candidate** key.
  - One of the candidate keys is *arbitrarily* designated to be the **primary key**, and the others are called **secondary keys**.
- A **Prime attribute** must be a member of *some* candidate key
- A **Nonprime attribute** is not a prime attribute—that is, it is not a member of any candidate key.


# First Normal Form

- Disallows
  - composite attributes
  - multivalued attributes
  - **nested relations**; attributes whose values for an *individual tuple* are non-atomic
- Considered to be part of the definition of relation

# Normalization into 1NF

(a)

DEPARTMENT

Dname	<u>Dnumber</u>	Dmgr_ssn	Dlocations
			

(b)

DEPARTMENT

Dname	<u>Dnumber</u>	Dmgr_ssn	Dlocations
Research	5	333445555	{Bellaire, Sugarland, Houston}
Administration	4	987654321	{Stafford}
Headquarters	1	888665555	{Houston}

# Normalization into 1NF

**DEPARTMENT**

Dname	<u>Dnumber</u>	Dmgr_ssn	<u>Dlocation</u>
Research	5	333445555	Bellaire
Research	5	333445555	Sugarland
Research	5	333445555	Houston
Administration	4	987654321	Stafford
Headquarters	1	888665555	Houston

# Normalization nested relations into 1NF

(a)

EMP\_PROJ

Ssn	Ename	Projs	
		Pnumber	Hours

(b)

EMP\_PROJ

Ssn	Ename	Pnumber	Hours
123456789	Smith, John B.	1	32.5
		2	7.5
666884444	Narayan, Ramesh K.	3	40.0
453453453	English, Joyce A.	1	20.0
		2	20.0
333445555	Wong, Franklin T.	2	10.0
		3	10.0
		10	10.0
		20	10.0
999887777	Zelaya, AliciaJ.	30	30.0
		10	10.0
987987987	Jabbar, Ahmad V.	10	35.0
		30	5.0
987654321	Wallace, Jennifer S.	30	20.0
		20	15.0
888665555	Borg, James E.	20	NULL



# Normalization nested relations into 1NF

EMP\_PROJ1

<u>Ssn</u>	Ename
------------	-------

EMP\_PROJ2

<u>Ssn</u>	<u>Pnumber</u>	Hours
------------	----------------	-------

# Second Normal Form

- Uses the concepts of **FDs, primary key**
- Definitions
  - **Prime attribute:** An attribute that is member of the primary key K
  - **Full functional dependency:** a FD  $Y \rightarrow Z$  where removal of any attribute from Y means the FD does not hold any more
- Examples:
  - $\{SSN, PNUMBER\} \rightarrow HOURS$  is a full FD since neither  $SSN \rightarrow HOURS$  nor  $PNUMBER \rightarrow HOURS$  hold
  - $\{SSN, PNUMBER\} \rightarrow ENAME$  is not a full FD (it is called a partial dependency ) since  $SSN \rightarrow ENAME$  also holds

# Second Normal Form

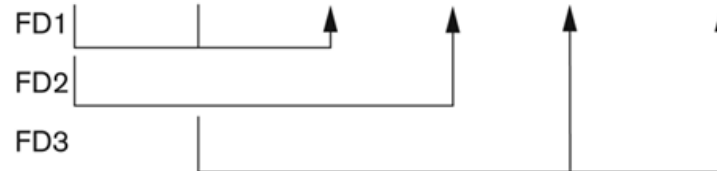
- A relation schema  $R$  is in **second normal form (2NF)** if every non-prime attribute  $A$  in  $R$  is fully functionally dependent on the primary key
- $R$  can be decomposed into 2NF relations via the process of 2NF normalization

# Normalizing into 2NF and 3NF

(a)

**EMP\_PROJ**

<u>Ssn</u>	<u>Pnumber</u>	Hours	Ename	Pname	Plocation
------------	----------------	-------	-------	-------	-----------



2NF Normalization

**EP1**

<u>Ssn</u>	<u>Pnumber</u>	Hours
------------	----------------	-------



**EP2**

<u>Ssn</u>	Ename
------------	-------



**EP3**

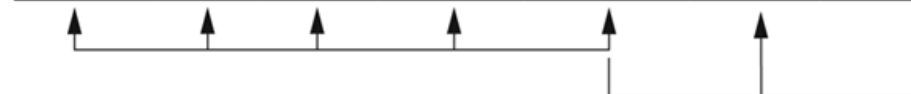
<u>Pnumber</u>	Pname	Plocation
----------------	-------	-----------



(b)

**EMP\_DEPT**

Ename	<u>Ssn</u>	Bdate	Address	Dnumber	Dname	Dmgr_ssn
-------	------------	-------	---------	---------	-------	----------



3NF Normalization

**ED1**

Ename	<u>Ssn</u>	Bdate	Address	Dnumber
-------	------------	-------	---------	---------

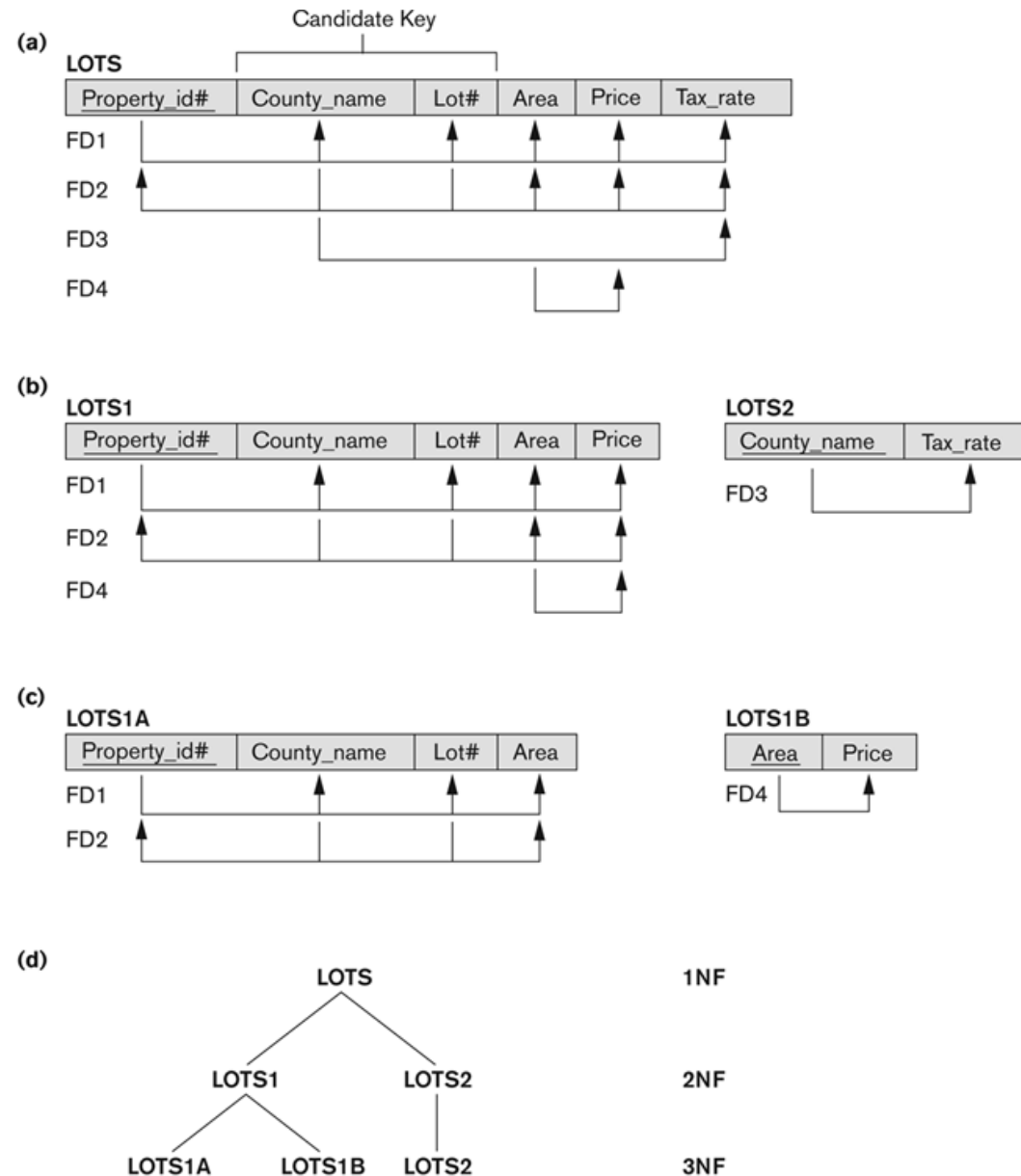


**ED2**

<u>Dnumber</u>	Dname	Dmgr_ssn
----------------	-------	----------



# Normalization into 2NF and 3NF



# Third Normal Form

- Definition:
  - **Transitive functional dependency:** a FD  $X \rightarrow Z$  that can be derived from two FDs  $X \rightarrow Y$  and  $Y \rightarrow Z$
- Examples:
  - $SSN \rightarrow DMGRSSN$  is a **transitive** FD
    - Since  $SSN \rightarrow DNUMBER$  and  $DNUMBER \rightarrow DMGRSSN$  hold
  - $SSN \rightarrow ENAME$  is **non-transitive**
    - Since there is no set of attributes  $X$  where  $SSN \rightarrow X$  and  $X \rightarrow ENAME$

# Third Normal Form

- A relation schema  $R$  is in **third normal form (3NF)** if it is in 2NF *and* no non-prime attribute  $A$  in  $R$  is transitively dependent on the primary key
- $R$  can be decomposed into 3NF relations via the process of 3NF normalization
- NOTE:
  - In  $X \rightarrow Y$  and  $Y \rightarrow Z$ , with  $X$  as the primary key, we consider this a problem only if  $Y$  is not a candidate key.
  - When  $Y$  is a candidate key, there is no problem with the transitive dependency .
  - E.g., Consider EMP (SSN, Emp#, Salary ).
    - Here,  $SSN \rightarrow Emp\# \rightarrow Salary$  and  $Emp\#$  is a candidate key.

# Normal Forms Defined Informally

- 1<sup>st</sup> normal form
  - All attributes depend on **the key**
- 2<sup>nd</sup> normal form
  - All attributes depend on **the whole key**
- 3<sup>rd</sup> normal form
  - All attributes depend on **nothing but the key**



## 4 General Normal Form Definitions (For Multiple Keys)

- The above definitions consider the primary key only
- The following more general definitions take into account relations with multiple candidate keys
- A relation schema  $R$  is in **second normal form (2NF)** if every non-prime attribute  $A$  in  $R$  is fully functionally dependent on *every* key of  $R$

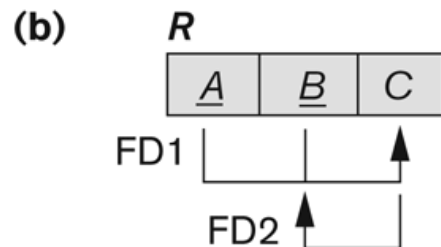
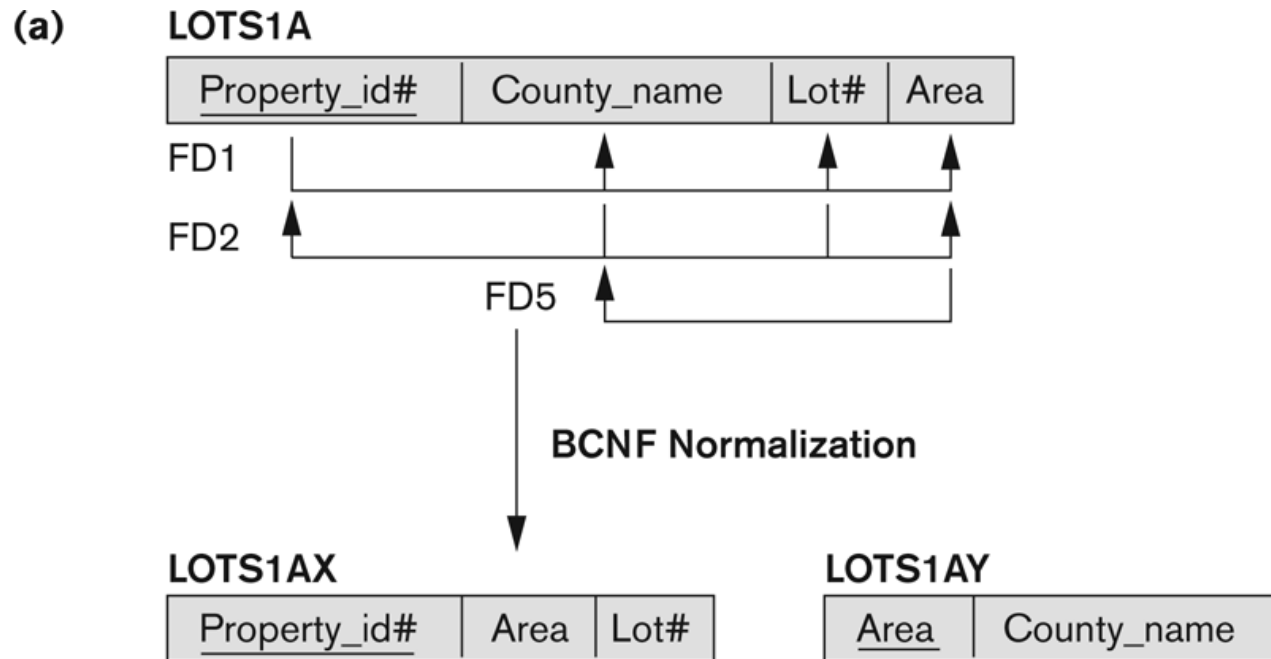
# General Normal Form Definitions

- Definition:
  - **Superkey** of relation schema R - a set of attributes S of R that contains a key of R
  - A relation schema R is in **third normal form (3NF)** if whenever a FD  $X \rightarrow A$  holds in R, then either:
    - (a) X is a superkey of R, or
    - (b) A is a prime attribute of R
- NOTE: Boyce-Codd normal form disallows condition (b) above

# BCNF (Boyce-Codd Normal Form)

- A relation schema  $R$  is in **Boyce-Codd Normal Form (BCNF)** if whenever an **FD  $X \rightarrow A$**  holds in  $R$ , then  **$X$  is a superkey** of  $R$
- Each normal form is strictly stronger than the previous one
  - Every 2NF relation is in 1NF
  - Every 3NF relation is in 2NF
  - Every BCNF relation is in 3NF
- There exist relations that are in 3NF but not in BCNF
- The goal is to have each relation in BCNF (or 3NF)

# Boyce-Codd normal form



**Figure 10.12**

Boyce-Codd normal form. (a) BCNF normalization of LOTSA with the functional dependency FD2 being lost in the decomposition. (b) A schematic relation with FDs; it is in 3NF, but not in BCNF.

a relation TEACH that is in 3NF but not in BCNF

**TEACH**

Student	Course	Instructor
Narayan	Database	Mark
Smith	Database	Navathe
Smith	Operating Systems	Ammar
Smith	Theory	Schulman
Wallace	Database	Mark
Wallace	Operating Systems	Ahamad
Wong	Database	Omiecinski
Zelaya	Database	Navathe
Narayan	Operating Systems	Ammar

**Figure 10.13**  
A relation TEACH that  
is in 3NF but not  
BCNF.

# Achieving the BCNF by Decomposition (1)

- Two FDs exist in the relation TEACH:
  - fd1: { student, course} -> instructor
  - fd2: instructor -> course
- {student, course} is a candidate key for this relation and that the dependencies shown follow the pattern in Figure 10.12 (b).
  - So this relation is in 3NF *but not in* BCNF
- A relation **NOT** in BCNF should be decomposed so as to meet this property, while possibly forgoing the preservation of all functional dependencies in the decomposed relations.
  - (See Algorithm 11.3)

# Achieving the BCNF by Decomposition (2)

- Three possible decompositions for relation TEACH
  - {student, instructor} and {student, course}
  - {course, instructor} and {course, student}
  - {instructor, course} and {instructor, student}
- All three decompositions will lose fd1.
  - We have to settle for sacrificing the functional dependency preservation. But we cannot sacrifice the non-additivity property after decomposition.
- Out of the above three, only the 3rd decomposition will not generate spurious tuples after join.(and hence has the non-additivity property).
- A test to determine whether a binary decomposition (decomposition into two relations) is non-additive (lossless) is discussed in section 11.1.4 under Property LJ1. Verify that the third decomposition above meets the property.

# Properties of Relational Decompositions

**Universal relation schema**  $R = \{A_1, A_2, \dots, A_n\}$  that includes *all* the attributes of the database.

Using the functional dependencies, the algorithms decompose the universal relation schema  $R$  into a set of relation schemas  $D = \{R_1, R_2, \dots, R_m\}$  that will become the relational database schema;  $D$  is called a **decomposition** of  $R$ .

This is called the **attribute preservation** condition of a decomposition.



# Properties of Relational Decompositions

**Definition.** Formally, a decomposition  $D = \{R_1, R_2, \dots, R_m\}$  of  $R$  has the **lossless (non additive) join property** with respect to the set of dependencies  $F$  on  $R$  if, for *every* relation state  $r$  of  $R$  that satisfies  $F$ , the following holds, where  $*$  is the NATURAL JOIN of all the relations in  $D$ :  $*(\pi R_1(r), \dots, \pi R_m(r)) = r$ .

# Properties of Relational Decompositions

## **Algorithm 16.3.** Testing for Nonadditive Join Property

**Input:** A universal relation  $R$ , a decomposition  $D = \{R_1, R_2, \dots, R_m\}$  of  $R$ , and a set  $F$  of functional dependencies.

*Note:* Explanatory comments are given at the end of some of the steps. They follow the format: (*\* comment \**).

1. Create an initial matrix  $S$  with one row  $i$  for each relation  $R_i$  in  $D$ , and one column  $j$  for each attribute  $A_j$  in  $R$ .
2. Set  $S(i, j) := b_{ij}$  for all matrix entries. (*\* each  $b_{ij}$  is a distinct symbol associated with indices  $(i, j)$  \**).
3. For each row  $i$  representing relation schema  $R_i$   
{for each column  $j$  representing attribute  $A_j$   
{if (relation  $R_i$  includes attribute  $A_j$ ) then set  $S(i, j) := a_j$ ;}}; (*\* each  $a_j$  is a distinct symbol associated with index  $(j)$  \**).

# Properties of Relational Decompositions

4. Repeat the following loop until a *complete loop execution* results in no changes to  $S$

{for each functional dependency  $X \rightarrow Y$  in  $F$

{for all rows in  $S$  that have the same symbols in the columns corresponding

to attributes in  $X$

{make the symbols in each column that correspond to an attribute in  $Y$  be the same in all these rows as follows: If any of the rows has an  $a$  symbol for the column, set the other rows to that *same a* symbol in the column.

If no  $a$  symbol exists for the attribute in any of the rows, choose one of the  $b$  symbols that appears in one of the rows for the attribute and set the other rows to that same  $b$  symbol in the column ;} ; } ;};

5. If a row is made up entirely of  $a$  symbols, then the decomposition has the non-additive join property; otherwise, it does not.

**(continued)**  
 Lossless  
 (nonadditive) join  
 test for  $n$ -ary  
 decompositions.  
 (c) Case 2:  
 Decomposition of  
 EMP\_PROJ into EMP,  
 PROJECT, and  
 WORKS\_ON satisfies  
 test.

(c)  $R = \{SSN, ENAME, PNUMBER, PNAME, PLOCATION, HOURS\}$   $D = \{R_1, R_2, R_3\}$   
 $R_1 = EMP = \{SSN, ENAME\}$   
 $R_2 = PROJ = \{PNUMBER, PNAME, PLOCATION\}$   
 $R_3 = WORKS\_ON = \{SSN, PNUMBER, HOURS\}$

$F = \{SSN \rightarrow \{ENAME, PNUMBER\}; PNUMBER \rightarrow \{PNAME, PLOCATION\}; \{SSN, PNUMBER\} \rightarrow HOURS\}$

	SSN	ENAME	PNUMBER	PNAME	PLOCATION	HOURS
$R_1$	$a_1$	$a_2$	$b_{13}$	$b_{14}$	$b_{15}$	$b_{16}$
$R_2$	$b_{21}$	$b_{22}$	$a_3$	$a_4$	$a_5$	$b_{26}$
$R_3$	$a_1$	$b_{32}$	$a_3$	$b_{34}$	$b_{35}$	$a_6$

(original matrix S at start of algorithm)

	SSN	ENAME	PNUMBER	PNAME	PLOCATION	HOURS
$R_1$	$a_1$	$a_2$	$b_{13}$	$b_{14}$	$b_{15}$	$b_{16}$
$R_2$	$b_{21}$	$b_{22}$	$a_3$	$a_4$	$a_5$	$b_{26}$
$R_3$	$a_1$	<del><math>b_{32}</math></del> $a_2$	$a_3$	<del><math>b_{34}</math></del> $a_4$	<del><math>b_{35}</math></del> $a_5$	$a_6$

(matrix S after applying the first two functional dependencies -  
 last row is all "a" symbols, so we stop)

Lossless (nonadditive) join test for  $n$ -ary decompositions.

(a) Case 1: Decomposition of EMP\_PROJ into EMP\_PROJ1 and EMP\_LOCS fails test. (b) A decomposition of EMP\_PROJ that has the lossless join property.

(a)  $R = \{SSN, ENAME, PNUMBER, PNAME, PLOCATION, HOURS\}$   $D = \{R_1, R_2\}$   
 $R_1 = \text{EMP\_LOCS} = \{ENAME, PLOCATION\}$   
 $R_2 = \text{EMP\_PROJ1} = \{SSN, PNUMBER, HOURS, PNAME, PLOCATION\}$   
 $F = \{SSN \rightarrow ENAME; PNUMBER \rightarrow \{PNAME, PLOCATION\}; \{SSN, PNUMBER\} \rightarrow HOURS\}$

	SSN	ENAME	PNUMBER	PNAME	PLOCATION	HOURS
$R_1$	$b_{11}$	$a_2$	$b_{13}$	$b_{14}$	$a_5$	$b_{16}$
$R_2$	$a_1$	$b_{22}$	$a_3$	$a_4$	$a_5$	$a_6$

(no changes to matrix after applying functional dependencies)

(b)

EMP		PROJECT			WORKS_ON		
SSN	ENAME	PNUMBER	PNAME	PLOCATION	SSN	PNUMBER	HOURS