

# Lectures 8 and 9: Database Design

Friday, October 12  
and Monday, October 15, 2006

# Announcements/Reminders

- Homework 1: solutions are posted
- Homework 2: posted (due Wed. Oct. 24)
- Project Phase 1 due Wednesday

# Outline

- The relational data model: 3.1
- Functional dependencies: 3.4

# The Relational Data Model

- Main idea: store data in relations (= tables)
- What kind of tables ?
  - Flat tables = First Normal Form
  - No anomalies = Boyce Codd Normal Form

# First Normal Form (1NF)

- A database schema is in First Normal Form if all tables are flat

Student

Name	GPA	Courses			
Alice	3.8	<table border="1"> <tr><td>Math</td></tr> <tr><td>DB</td></tr> <tr><td>OS</td></tr> </table>	Math	DB	OS
Math					
DB					
OS					
Bob	3.7	<table border="1"> <tr><td>DB</td></tr> <tr><td>OS</td></tr> </table>	DB	OS	
DB					
OS					
Carol	3.9	<table border="1"> <tr><td>Math</td></tr> <tr><td>OS</td></tr> </table>	Math	OS	
Math					
OS					

Student

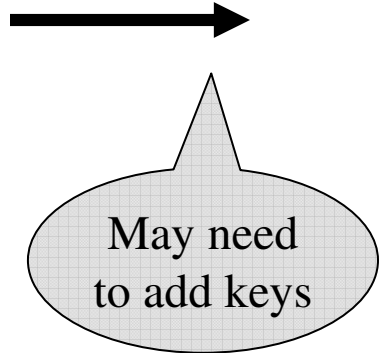
Name	GPA
Alice	3.8
Bob	3.7
Carol	3.9

Takes

Student	Course
Alice	Math
Carol	Math
Alice	DB
Bob	DB
Alice	OS
Carol	OS

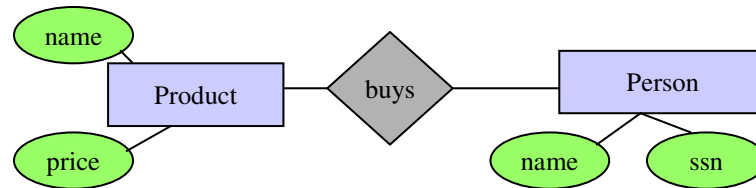
Course

Course
Math
DB
OS

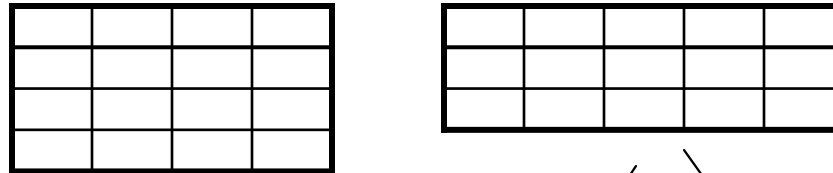


# Relational Schema Design

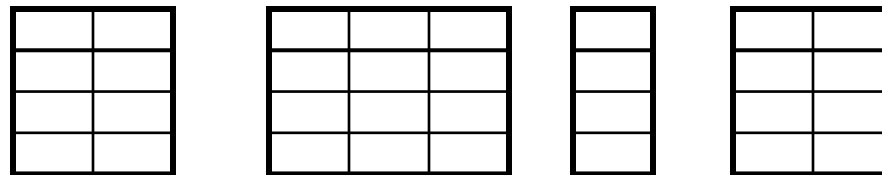
Conceptual Model:



Relational Model:  
plus FD's



Normalization:  
Eliminates anomalies



# Data Anomalies

When a database is poorly designed we get anomalies:

**Redundancy**: data is repeated

**Updated anomalies**: need to change in several places

**Delete anomalies**: may lose data when we don't want

# Relational Schema Design

Recall set attributes (persons with several phones):

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield

One person may have multiple phones, but lives in only one city

## Anomalies:

- Redundancy = repeat data
- Update anomalies = Fred moves to “Bellevue”
- Deletion anomalies = Joe deletes his phone number:  
what is his city ?



# Relation Decomposition

**Break the relation into two:**

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield

Name	<u>SSN</u>	City
Fred	123-45-6789	Seattle
Joe	987-65-4321	Westfield

<u>SSN</u>	<u>PhoneNumber</u>
123-45-6789	206-555-1234
123-45-6789	206-555-6543
987-65-4321	908-555-2121

**Anomalies have gone:**

- No more repeated data
- Easy to move Fred to “Bellevue” (how ?)
- Easy to delete all Joe’s phone number (how ?)

# Relational Schema Design (or Logical Design)

Main idea:

- Start with some relational schema
- Find out its *functional dependencies*
- Use them to design a better relational schema

# Functional Dependencies

- A form of constraint
  - hence, part of the schema
- Finding them is part of the database design
- Also used in normalizing the relations

# Functional Dependencies

## Definition:

If two tuples agree on the attributes

$$A_1, A_2, \dots, A_n$$

then they must also agree on the attributes

$$B_1, B_2, \dots, B_m$$

## Formally:

$$A_1, A_2, \dots, A_n \rightarrow B_1, B_2, \dots, B_m$$

# When Does an FD Hold

Definition:  $A_1, \dots, A_m \rightarrow B_1, \dots, B_n$  holds in R if:

$$\forall t, t' \in R, (t.A_1=t'.A_1 \wedge \dots \wedge t.A_m=t'.A_m \Rightarrow t.B_1=t'.B_1 \wedge \dots \wedge t.B_n=t'.B_n)$$

R

	$A_1$	...	$A_m$		$B_1$	...	$B_m$		
t									
t'									

if t, t' agree here      then t, t' agree here

# Examples

An FD holds, or does not hold on an instance:

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	Lawyer

EmpID  $\rightarrow$  Name, Phone, Position

Position  $\rightarrow$  Phone

but not Phone  $\rightarrow$  Position

# Example

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876 ←	Salesrep
E1111	Smith	9876 ←	Salesrep
E9999	Mary	1234	Lawyer

Position → Phone

# Example

EmpID	Name	Phone	Position
E0045	Smith	1234 →	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234 →	Lawyer

but not Phone → Position



# Example

FD's are constraints:

- On some instances they hold
- On others they don't

name  $\rightarrow$  color  
category  $\rightarrow$  department  
color, category  $\rightarrow$  price

name	category	color	department	price
Gizmo	Gadget	Green	Toys	49
Tweaker	Gadget	Green	Toys	99

Does this instance satisfy all the FDs ?

# Example

name → color  
category → department  
color, category → price

name	category	color	department	price
Gizmo	Gadget	Green	Toys	49
Tweaker	Gadget	Black	Toys	99
Gizmo	Stationary	Green	Office-supp.	59

What about this one ?

# An Interesting Observation

If all these FDs are true:

name  $\rightarrow$  color  
category  $\rightarrow$  department  
color, category  $\rightarrow$  price

Then this FD also holds:

name, category  $\rightarrow$  price

Why ??

# Goal: Find ALL Functional Dependencies

- Anomalies occur when certain “bad” FDs hold
- We know some of the FDs
- Need to find *all* FDs, then look for the bad ones

# Armstrong's Rules (1/3)

$$A_1, A_2, \dots, A_n \rightarrow B_1, B_2, \dots, B_m$$

Is equivalent to

$$\begin{aligned} A_1, A_2, \dots, A_n &\rightarrow B_1 \\ A_1, A_2, \dots, A_n &\rightarrow B_2 \\ &\dots \\ A_1, A_2, \dots, A_n &\rightarrow B_m \end{aligned}$$

**Splitting rule  
and  
Combing rule**

	A1	...	Am		B1	...	Bm	

# Armstrong's Rules (1/3)

$$A_1, A_2, \dots, A_n \rightarrow A_i$$

**Trivial Rule**

where  $i = 1, 2, \dots, n$

Why ?

	$A_1$	...	$A_m$	

# Armstrong's Rules (1/3)

## Transitive Closure Rule

If

$$A_1, A_2, \dots, A_n \rightarrow B_1, B_2, \dots, B_m$$

and

$$B_1, B_2, \dots, B_m \rightarrow C_1, C_2, \dots, C_p$$

then

$$A_1, A_2, \dots, A_n \rightarrow C_1, C_2, \dots, C_p$$

Why ?

	$A_1$	...	$A_m$		$B_1$	...	$B_m$		$C_1$	...	$C_p$	



# Example (continued)

Start from the following FDs:

1. name  $\rightarrow$  color
2. category  $\rightarrow$  department
3. color, category  $\rightarrow$  price

Infer the following FDs:

Inferred FD	Which Rule did we apply ?
4. name, category $\rightarrow$ name	
5. name, category $\rightarrow$ color	
6. name, category $\rightarrow$ category	
7. name, category $\rightarrow$ color, category	
8. name, category $\rightarrow$ price	

# Example (continued)

Answers:

1. name  $\rightarrow$  color
2. category  $\rightarrow$  department
3. color, category  $\rightarrow$  price

Inferred FD	Which Rule did we apply ?
4. name, category $\rightarrow$ name	Trivial rule
5. name, category $\rightarrow$ color	Transitivity on 4, 1
6. name, category $\rightarrow$ category	Trivial rule
7. name, category $\rightarrow$ color, category	Split/combine on 5, 6
8. name, category $\rightarrow$ price	Transitivity on 3, 7

THIS IS TOO HARD ! Let's see an easier way.

# Closure of a set of Attributes

**Given** a set of attributes  $A_1, \dots, A_n$

The **closure**,  $\{A_1, \dots, A_n\}^+$  = the set of attributes  $B$   
s.t.  $A_1, \dots, A_n \rightarrow B$

Example:

$\text{name} \rightarrow \text{color}$   
 $\text{category} \rightarrow \text{department}$   
 $\text{color, category} \rightarrow \text{price}$

Closures:

$\text{name}^+ = \{\text{name}, \text{color}\}$

$\{\text{name}, \text{category}\}^+ = \{\text{name}, \text{category}, \text{color}, \text{department}, \text{price}\}$

$\text{color}^+ = \{\text{color}\}$

# Closure Algorithm

$X = \{A_1, \dots, A_n\}$ .

**Repeat until**  $X$  doesn't change **do:**

**if**  $B_1, \dots, B_n \rightarrow C$  is a FD **and**  
 $B_1, \dots, B_n$  are all in  $X$   
**then** add  $C$  to  $X$ .

Example:

$\text{name} \rightarrow \text{color}$   
 $\text{category} \rightarrow \text{department}$   
 $\text{color, category} \rightarrow \text{price}$

$\{\text{name, category}\}^+ =$

$\{\text{name, category, color, department, price}\}$

Hence:  $\text{name, category} \rightarrow \text{color, department, price}$

# Example

In class:

$R(A,B,C,D,E,F)$

$A, B$	$\rightarrow$	$C$
$A, D$	$\rightarrow$	$E$
$B$	$\rightarrow$	$D$
$A, F$	$\rightarrow$	$B$

Compute  $\{A,B\}^+$   $X = \{A, B,$   $\}$

Compute  $\{A, F\}^+$   $X = \{A, F,$   $\}$

# Why Do We Need Closure

- With closure we can find all FD's easily
- To check if  $X \rightarrow A$ 
  - Compute  $X^+$
  - Check if  $A \in X^+$

# Using Closure to Infer ALL FDs

Example:

$$\begin{array}{l} A, B \rightarrow C \\ A, D \rightarrow B \\ B \rightarrow D \end{array}$$

Step 1: Compute  $X^+$ , for every  $X$ :

$$A^+ = A, \quad B^+ = BD, \quad C^+ = C, \quad D^+ = D$$
$$AB^+ = ABCD, \quad AC^+ = AC, \quad AD^+ = ABCD,$$
$$BC^+ = BCD, \quad BD^+ = BD, \quad CD^+ = CD$$
$$ABC^+ = ABD^+ = ACD^+ = ABCD \text{ (no need to compute— why ?)}$$
$$BCD^+ = BCD, \quad ABCD^+ = ABCD$$

Step 2: Enumerate all FD's  $X \rightarrow Y$ , s.t.  $Y \subseteq X^+$  and  $X \cap Y = \emptyset$ :

$$AB \rightarrow CD, \quad AD \rightarrow BC, \quad ABC \rightarrow D, \quad ABD \rightarrow C, \quad ACD \rightarrow B$$

# Another Example

- Enrollment(student, major, course, room, time)  
student → major  
major, course → room  
course → time

What else can we infer ? [in class, or at home]



# Keys

- A **superkey** is a set of attributes  $A_1, \dots, A_n$  s.t. for any other attribute  $B$ , we have  $A_1, \dots, A_n \rightarrow B$
- A **key** is a minimal superkey
  - I.e. set of attributes which is a superkey and for which no subset is a superkey

# Computing (Super)Keys

- Compute  $X^+$  for all sets  $X$
- If  $X^+ =$  all attributes, then  $X$  is a key
- List only the minimal  $X$ 's

# Example

Product(name, price, category, color)

name, category  $\rightarrow$  price  
category  $\rightarrow$  color

What is the key ?

# Example

Product(name, price, category, color)

name, category  $\rightarrow$  price  
category  $\rightarrow$  color

What is the key ?

(name, category) + = name, category, price, color

Hence (name, category) is a key

# Examples of Keys

Enrollment(student, address, course, room, time)

student  $\rightarrow$  address

room, time  $\rightarrow$  course

student, course  $\rightarrow$  room, time

(find keys at home)

# Eliminating Anomalies

Main idea:

- $X \rightarrow A$  is OK if  $X$  is a (super)key
- $X \rightarrow A$  is not OK otherwise

# Example

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield
Joe	987-65-4321	908-555-1234	Westfield

SSN  $\rightarrow$  Name, City

What the key?

{SSN, PhoneNumber}

Hence SSN  $\rightarrow$  Name, City  
is a “bad” dependency

# Key or Keys ?

Can we have more than one key ?

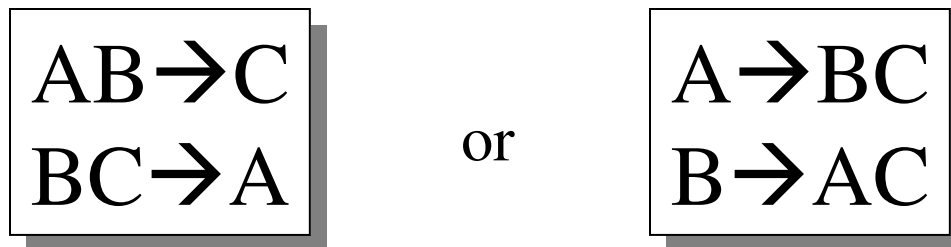
Given  $R(A,B,C)$  define FD's s.t. there are two or more keys



# Key or Keys ?

Can we have more than one key ?

Given  $R(A,B,C)$  define FD's s.t. there are two or more keys



what are the keys here ?

Can you design FDs such that there are *three* keys ?

# Boyce-Codd Normal Form

A simple condition for removing anomalies from relations:

A relation  $R$  is in BCNF if:

If  $A_1, \dots, A_n \rightarrow B$  is a non-trivial dependency in  $R$ , then  $\{A_1, \dots, A_n\}$  is a superkey for  $R$

In other words: there are no “bad” FDs

Equivalently:

$\forall X$ , either  $(X^+ = X)$  or  $(X^+ = \text{all attributes})$

# BCNF Decomposition Algorithm

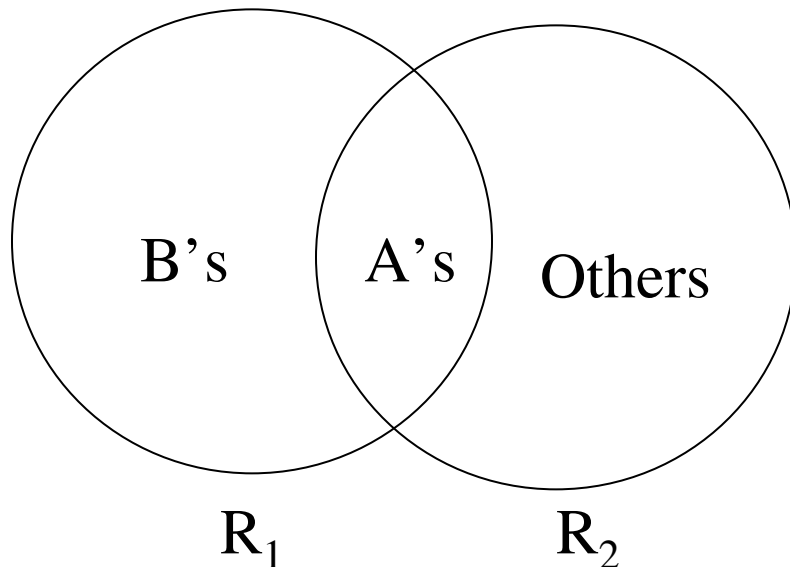
**repeat**

choose  $A_1, \dots, A_m \rightarrow B_1, \dots, B_n$  that violates BCNF

split  $R$  into  $R_1(A_1, \dots, A_m, B_1, \dots, B_n)$  and  $R_2(A_1, \dots, A_m, [\text{others}])$

continue with both  $R_1$  and  $R_2$

**until** no more violations



Is there a  
2-attribute  
relation that is  
not in BCNF ?

In practice, we have  
a better algorithm (coming<sup>43</sup> up)

# Example

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield
Joe	987-65-4321	908-555-1234	Westfield

SSN → Name, City

What the key?

{SSN, PhoneNumber}

use SSN → Name, City  
to split

# Example

<u>Name</u>	<u>SSN</u>	<u>City</u>
Fred	123-45-6789	Seattle
Joe	987-65-4321	Westfield

SSN → Name, City

<u>SSN</u>	<u>PhoneNumber</u>
123-45-6789	206-555-1234
123-45-6789	206-555-6543
987-65-4321	908-555-2121
987-65-4321	908-555-1234

Let's check anomalies:

- Redundancy ?
- Update ?
- Delete ?

# Example Decomposition

Person(name, SSN, age, hairColor, phoneNumber)

SSN  $\rightarrow$  name, age

age  $\rightarrow$  hairColor

Decompose in BCNF (in class):

# BCNF Decomposition Algorithm

BCNF\_Decompose(R)

find  $X$  s.t.:  $X \neq X^+ \neq$  [all attributes]

**if** (not found) **then** “R is in BCNF”

**let**  $Y = X^+ - X$

**let**  $Z =$  [all attributes]  $- X^+$

decompose R into  $R_1(X \cup Y)$  and  $R_2(X \cup Z)$

continue to decompose recursively  $R_1$  and  $R_2$

Find  $X$  s.t.:  $X \neq X^+ \neq$  [all attributes]

# Example BCNF Decomposition

Person(name, SSN, age, hairColor, phoneNumber)

SSN  $\rightarrow$  name, age

age  $\rightarrow$  hairColor

Iteration 1: Person

SSN<sup>+</sup> = SSN, name, age, hairColor

Decompose into: P(SSN, name, age, hairColor)

Phone(SSN, phoneNumber)

Iteration 2: P

age<sup>+</sup> = age, hairColor

Decompose: People(SSN, name, age)

Hair(age, hairColor)

Phone(SSN, phoneNumber)

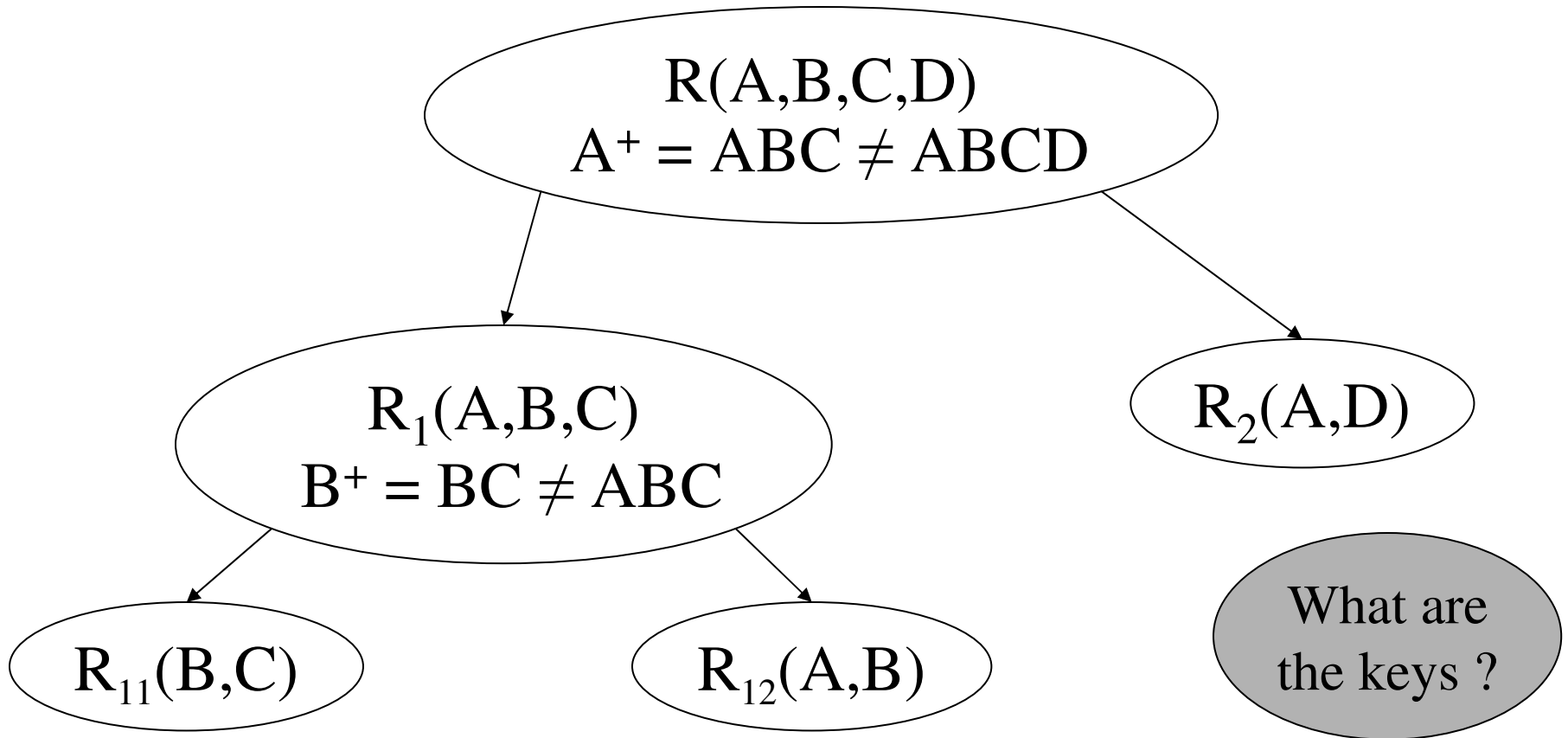
What are  
the keys ?



R(A,B,C,D)

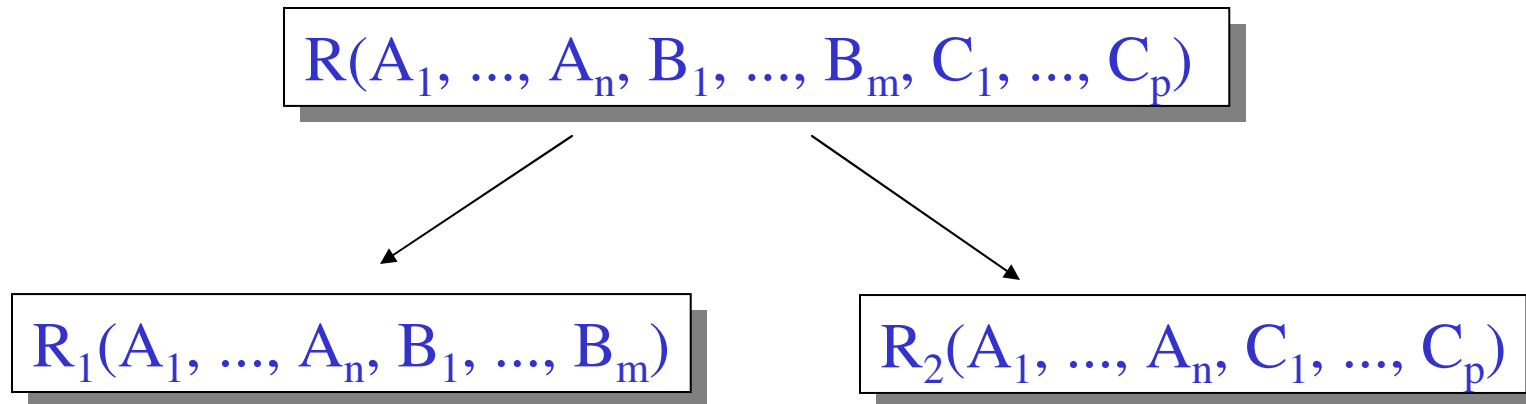
A → B  
B → C

# Example



What happens if in R we first pick B<sup>+</sup> ? Or AB<sub>49</sub><sup>+</sup> ?

# Decompositions in General




$R_1$  = projection of  $R$  on  $A_1, \dots, A_n, B_1, \dots, B_m$

$R_2$  = projection of  $R$  on  $A_1, \dots, A_n, C_1, \dots, C_p$

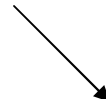
# Theory of Decomposition

- Sometimes it is correct:

Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera



Name	Price
Gizmo	19.99
OneClick	24.99
<del>Gizmo</del>	<del>19.99</del>



Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Lossless decomposition

# Incorrect Decomposition

- Sometimes it is not:

Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera

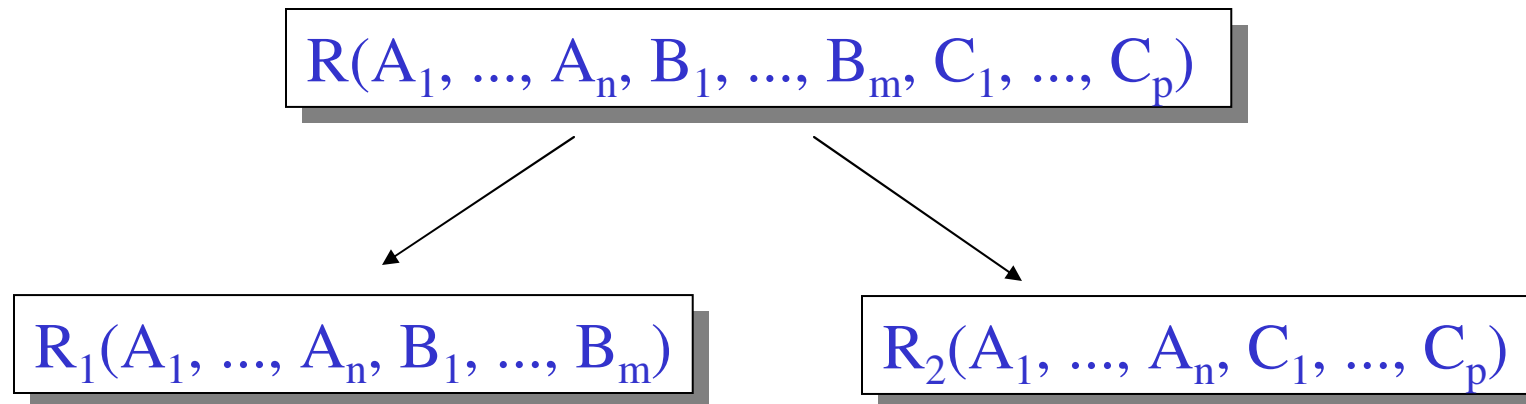
What's incorrect ??

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Price	Category
19.99	Gadget
24.99	Camera
19.99	Camera

Lossy decomposition

# Decompositions in General



If  $A_1, \dots, A_n \rightarrow B_1, \dots, B_m$   
Then the decomposition is lossless

Note: don't need  $A_1, \dots, A_n \rightarrow C_1, \dots, C_p$

BCNF decomposition is always lossless. WHY ?