

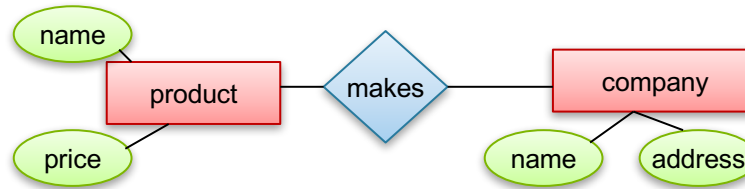
Database Systems

CSE 414

Lectures 18-19: Design Theory
(Ch. 3.1, 3.3-4)

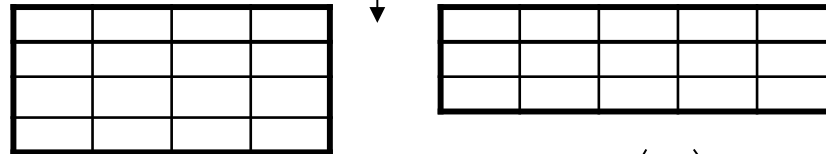
Database Design Process

Conceptual Model:



Relational Model:

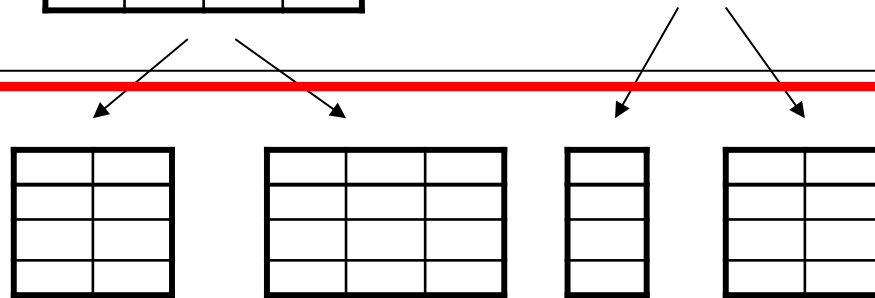
Tables + constraints
And also functional dep.



Normalization:

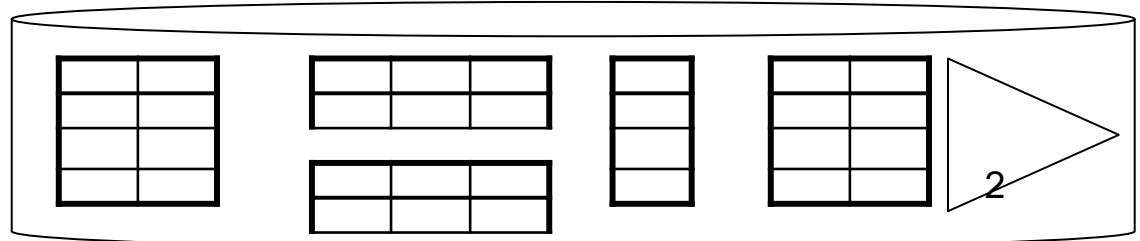
Eliminates anomalies

Conceptual Schema

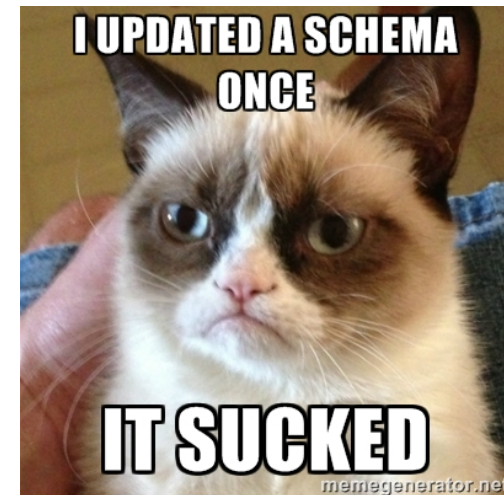
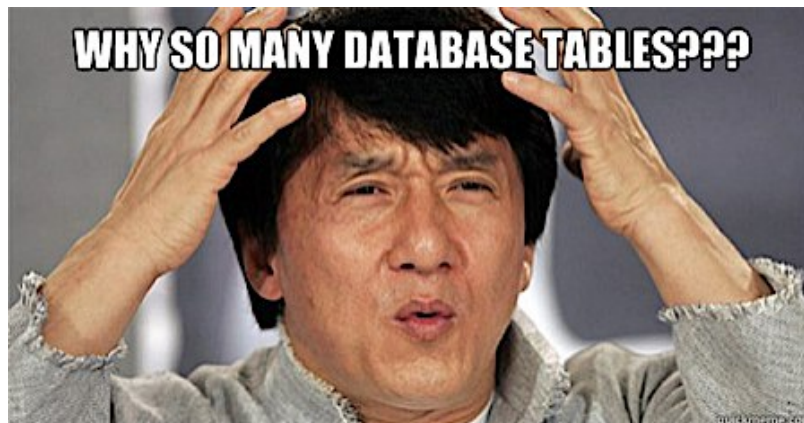


Physical storage details

Physical Schema



What makes good schemas?



Relational Schema Design

Name	<u>SSN</u>	<u>PhoneNumber</u>	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

Primary key is thus (SSN, PhoneNumber)

What is the problem with this schema?

Relational Schema Design

Name	<u>SSN</u>	<u>PhoneNumber</u>	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

These can cause bugs!
Worry most about later two.

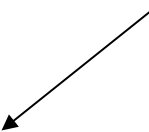
Anomalies:

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


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 numbers (how ?)

Relational Schema Design (or Logical Design)

How do we do this systematically?

- Start with some relational schema
- Find out its **functional dependencies** (FDs)
- Use FDs to **normalize** the relational schema

Functional Dependencies (FDs)

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:

$A_1 \dots A_n$ determines $B_1 \dots B_m$

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

Functional Dependencies (FDs)

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

for every pair of tuples $t, t' \in R$,

$(t.A_1 = t'.A_1 \text{ and } \dots \text{ and } t.A_m = t'.A_m \rightarrow t.B_1 = t'.B_1 \text{ and } \dots \text{ and } t.B_n = t'.B_n)$

R	A_1	...	A_m		B_1	...	B_n		
t									
t'									

Never have equal As
but different Bs!

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

Example

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 → Name, Phone, Position

Position → Phone

but not Phone → 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

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

Do all the FDs hold on this instance?

Example

name → color
category → department
color, category → price

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

What about this one ?

Terminology

- FD **holds** or **does not hold** on an *instance*
- If we can be sure that *every instance of R* will be one in which a given FD is true, then we say that **R satisfies the FD**
- If we say that R satisfies an FD F, we are **stating a constraint on R** (part of schema)

An Interesting Observation

If all these FDs are true:

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

Then this FD also holds:

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

If we find out from application domain that a relation satisfies some FDs, it doesn't mean that we found all the FDs that it satisfies!
There could be more FDs implied by the ones we have.

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:

1. name \rightarrow color
2. category \rightarrow department
3. color, category \rightarrow 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:

1. name \rightarrow color
2. category \rightarrow department
3. color, category \rightarrow 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	→	C
A, D	→	E
B	→	D
A, F	→	B

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

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

Example

In class:

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

A, B	→	C
A, D	→	E
B	→	D
A, F	→	B

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

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

Example

In class:

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

A, B	→	C
A, D	→	E
B	→	D
A, F	→	B

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

Compute $\{A, F\}^+$ $X = \{A, F, B, C, D, E\}$

Practice at Home

Find all FD's implied by:

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

Practice at Home

Find all FD's implied by:

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

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 $X \cap Y = \emptyset$:

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

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
 - superkey and for which no subset is a superkey

Computing (Super)Keys

- For all sets X , compute X^+
- If $X^+ = [\text{all attributes}]$, then X is a superkey
- Try only the minimal X 's to get the key

Example

Product(name, price, category, color)

name, category \rightarrow price
category \rightarrow color

What is the key?

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

Hence $\{\text{name, category}\}$ is a (super)key

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

$A \rightarrow B$
$B \rightarrow C$
$C \rightarrow A$

or

$AB \rightarrow C$
$BC \rightarrow A$

or

$A \rightarrow BC$
$B \rightarrow AC$

what are the keys here ?

Eliminating Anomalies

Main idea:

- $X \rightarrow A$ is OK if X is a (super)key
- $X \rightarrow A$ is not OK otherwise
 - Need to decompose the table, but how?

Boyce-Codd Normal Form

Boyce-Codd Normal Form

Dr. Raymond F. Boyce

Boyce-Codd Normal Form

There are no
“bad” FDs:

Definition. A relation R is in BCNF if:

Whenever $X \rightarrow A$ is a non-trivial dependency,
then X is a superkey.

Equivalently:

Definition. A relation R is in BCNF if:

$\forall X$, either $X^+ = X$ or $X^+ = [\text{all attributes}]$

BCNF Decomposition Algorithm

Normalize(R)

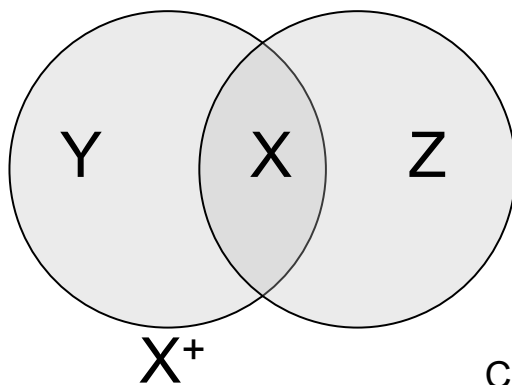
find X s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

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

let $Y = X^+ - X$; $Z = [\text{all attributes}] - X^+$

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

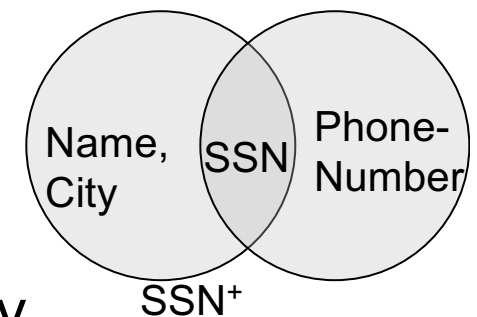
Normalize(R_1); Normalize(R_2);



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

$SSN \rightarrow Name, City$



The only key is: $\{SSN, PhoneNumber\}$

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

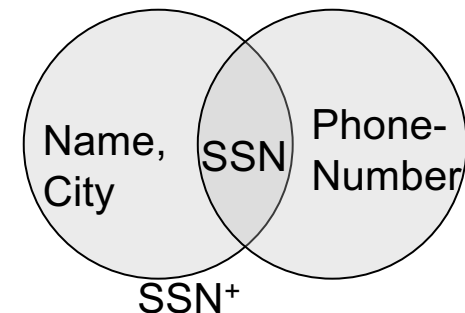
In other words:

$SSN^+ = SSN, Name, City$ and is neither SSN nor $All\ Attributes$

Example BCNF Decomposition

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

SSN \rightarrow 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 ?

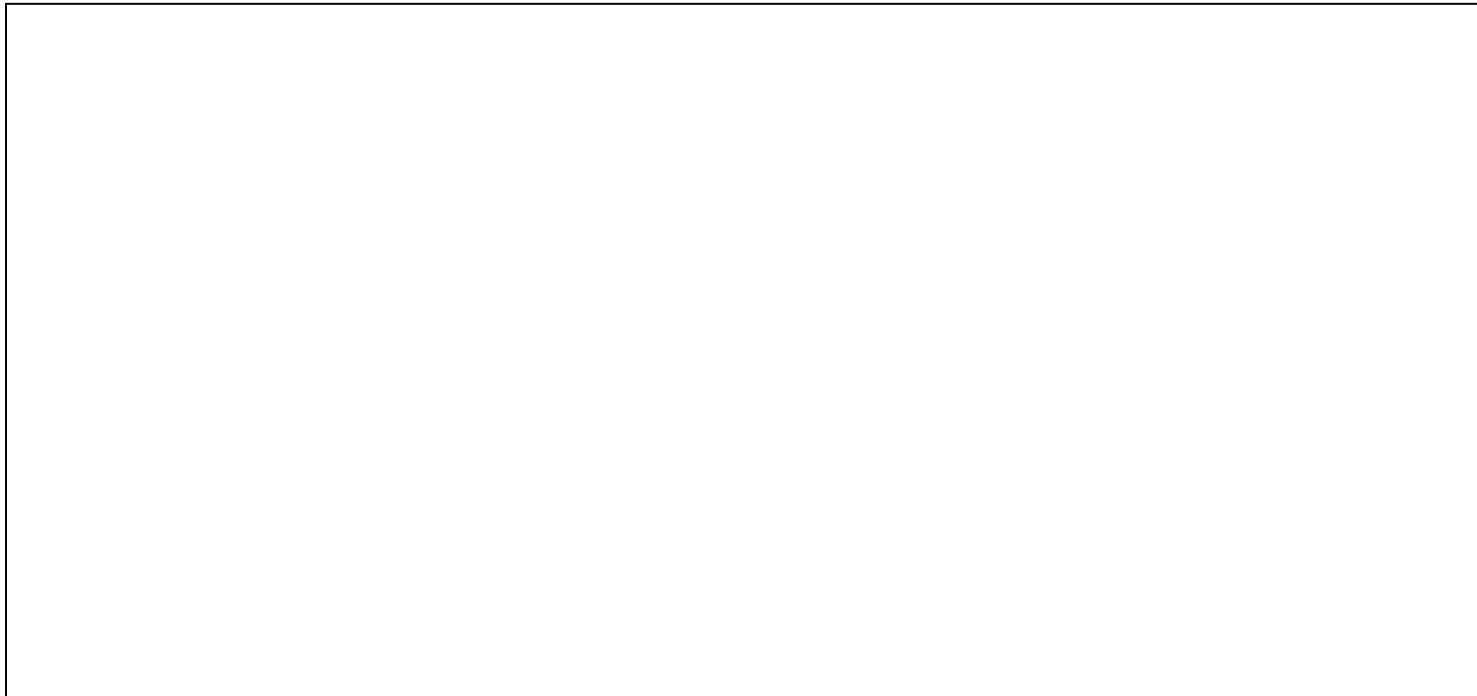
Find X s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

Example BCNF Decomposition

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

SSN \rightarrow name, age

age \rightarrow hairColor



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

Example BCNF Decomposition

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

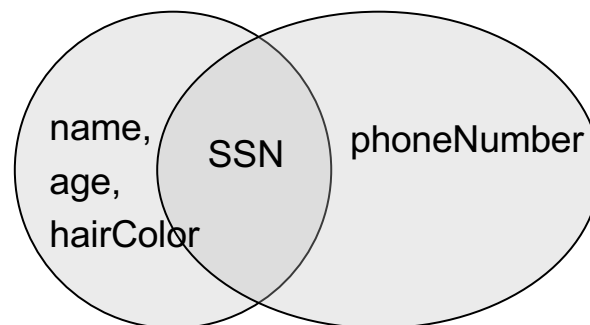
SSN \rightarrow name, age

age \rightarrow hairColor

Iteration 1: **Person**: SSN⁺ = SSN, name, age, hairColor

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

Phone(SSN, phoneNumber)



Find X s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

Example BCNF Decomposition

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

SSN \rightarrow name, age

age \rightarrow hairColor

What are
the keys ?

Iteration 1: **Person**: SSN⁺ = SSN, name, age, hairColor

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

Phone(SSN, phoneNumber)

Iteration 2: **P**: age⁺ = age, hairColor

Decompose: **People**(SSN, name, age)

Hair(age, hairColor)

Phone(SSN, phoneNumber)

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

Example BCNF Decomposition

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

SSN \rightarrow name, age

age \rightarrow hairColor

Note the keys!

Iteration 1: **Person**: SSN⁺ = SSN, name, age, hairColor

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

Phone(SSN, phoneNumber)

Iteration 2: **P**: age⁺ = age, hairColor

Decompose: **People**(SSN, name, age)

Hair(age, hairColor)

Phone(SSN, phoneNumber)

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

Example: BCNF

$A \rightarrow B$
$B \rightarrow C$

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

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

$A \rightarrow B$
 $B \rightarrow C$

Example: BCNF

Recall: find X s.t.
 $X \subsetneq X^+ \subsetneq [\text{all-attrs}]$

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

R(A,B,C,D)

A	→	B
B	→	C

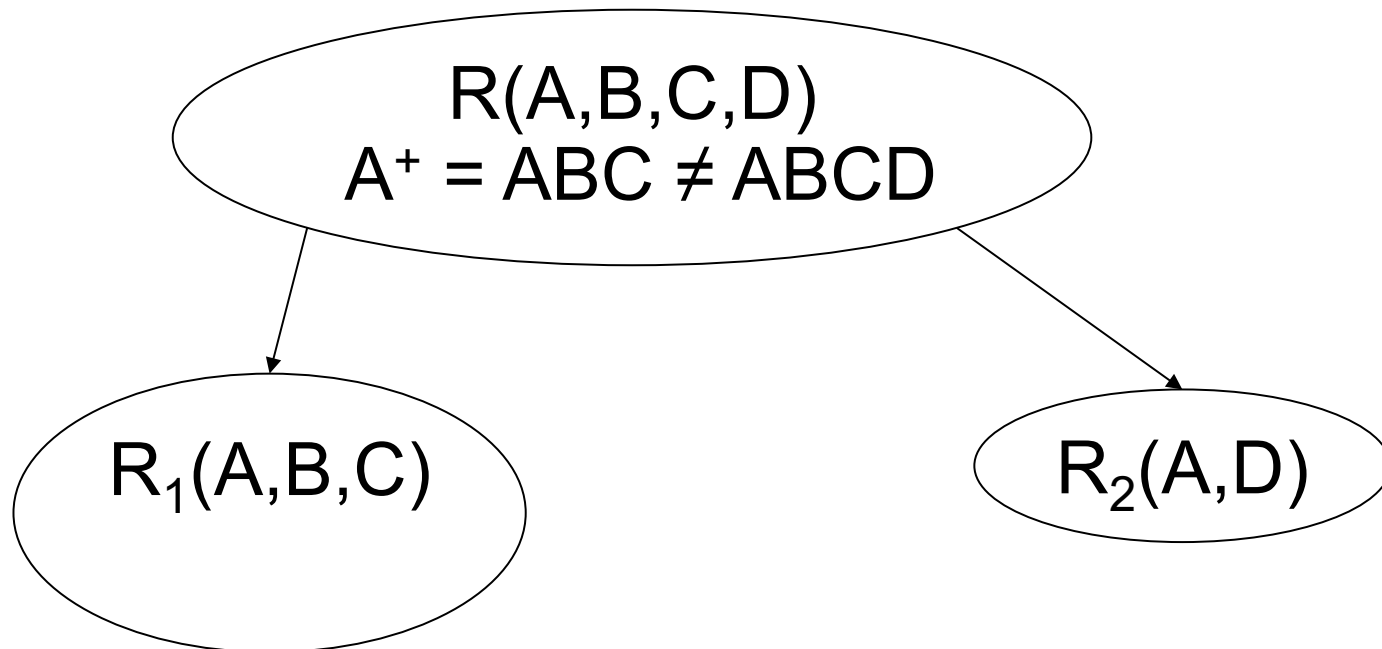
Example: BCNF

R(A,B,C,D)
 $A^+ = ABC \neq ABCD$

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

$A \rightarrow B$
$B \rightarrow C$

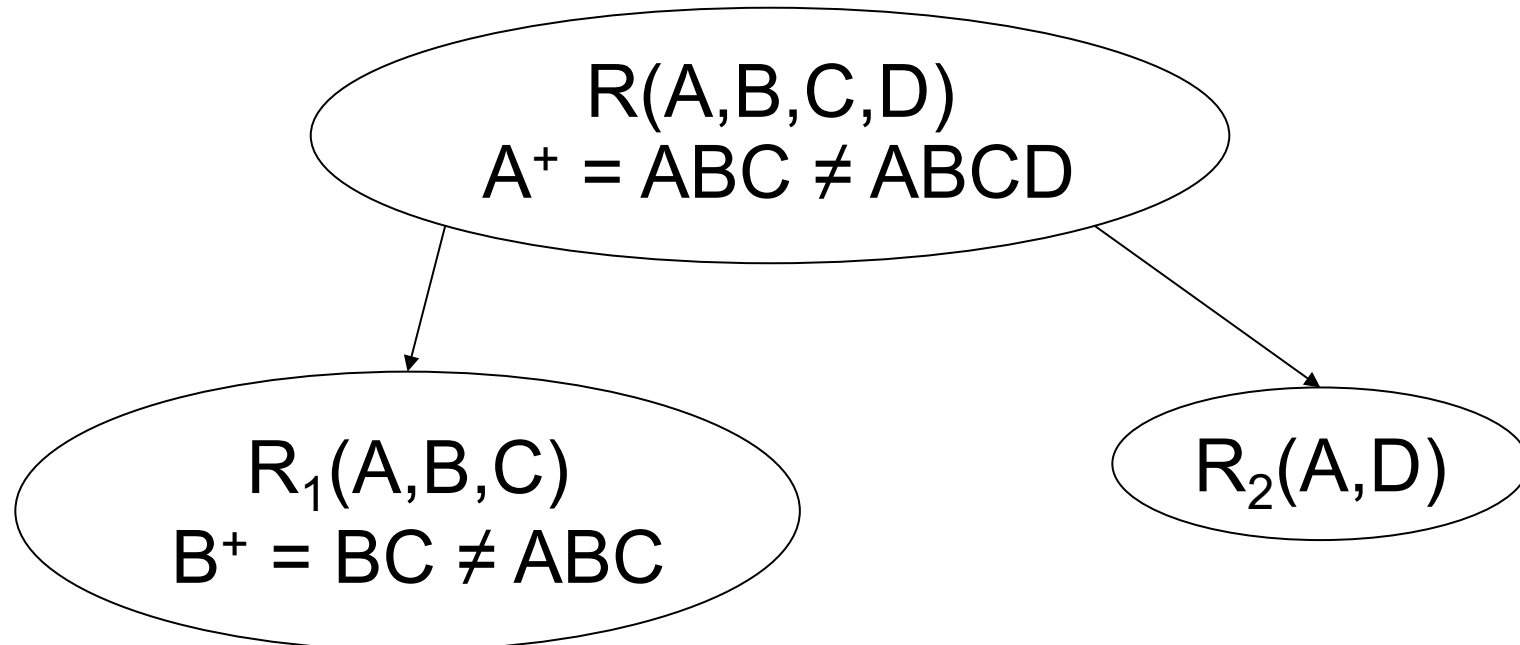
Example: BCNF



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

$A \rightarrow B$
$B \rightarrow C$

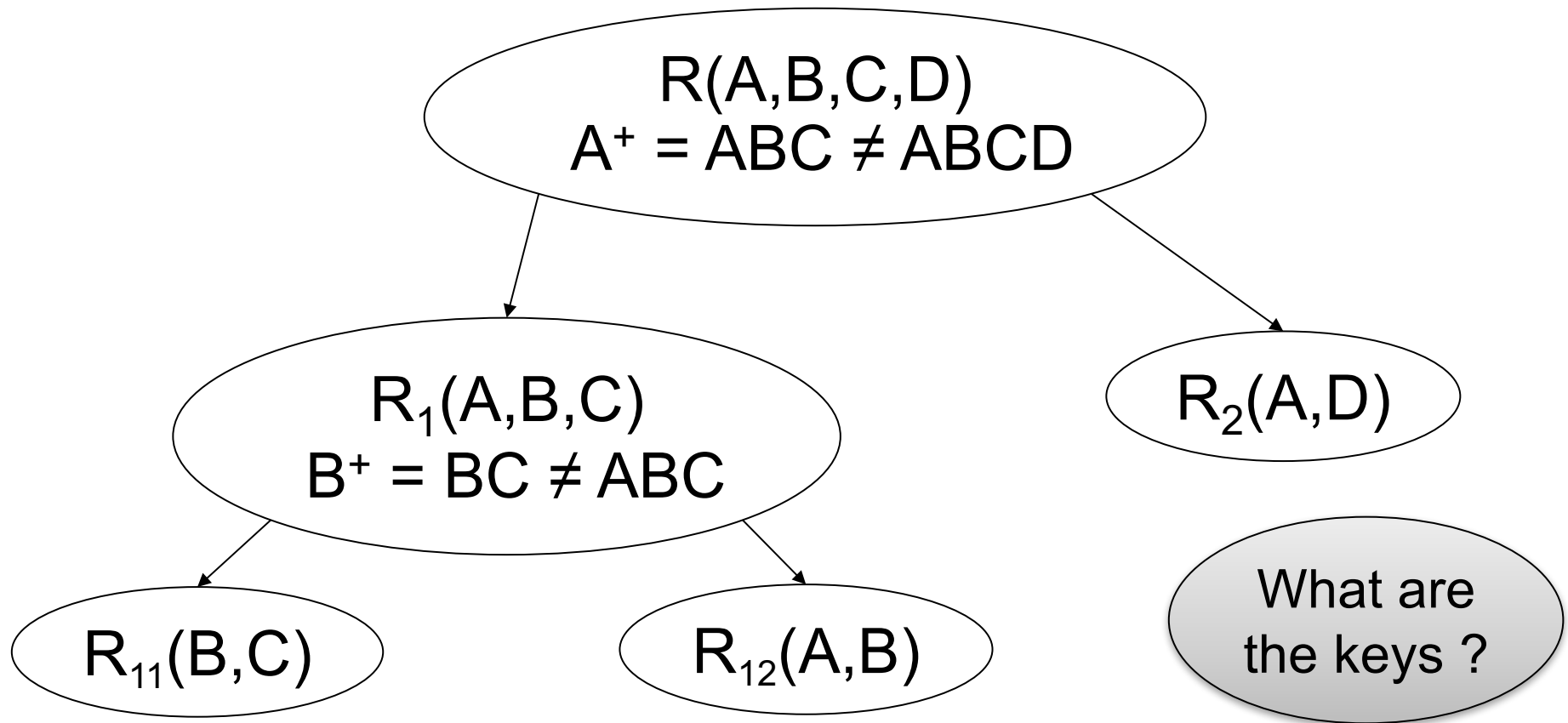
Example: BCNF



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

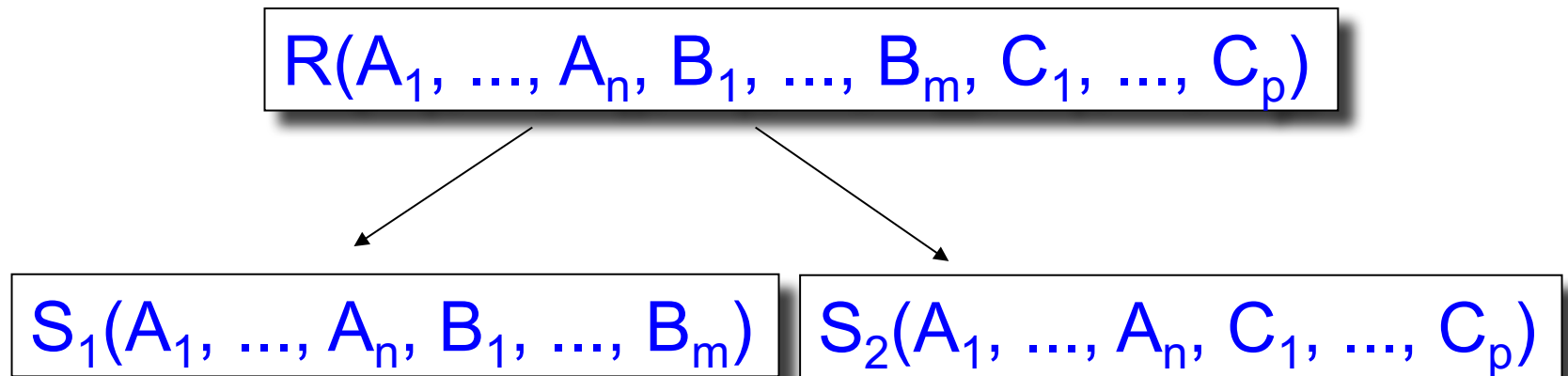
$A \rightarrow B$
 $B \rightarrow C$

Example: BCNF



What happens if in R we first pick B^+ ? Or AB^+ ?

Decompositions in General




S_1 = projection of R on $A_1, \dots, A_n, B_1, \dots, B_m$

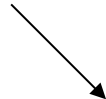
S_2 = projection of R on $A_1, \dots, A_n, C_1, \dots, C_p$

Lossless Decomposition

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



Name	Price
Gizmo	19.99
OneClick	24.99
Gizmo	19.99



Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Lossy Decomposition

What is lossy here?

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

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Price	Category
19.99	Gadget
24.99	Camera
19.99	Camera

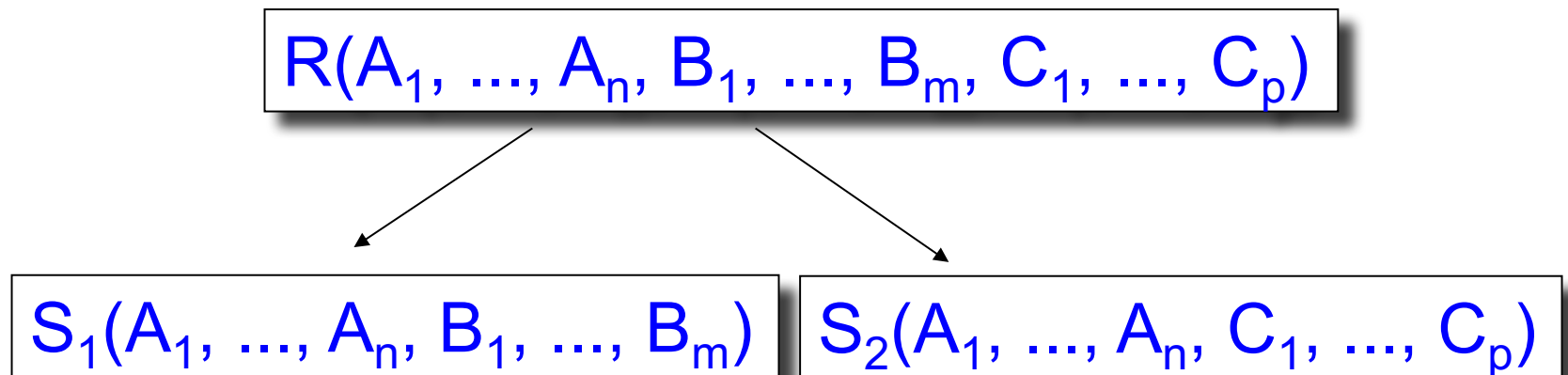
Lossy Decomposition

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

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Price	Category
19.99	Gadget
24.99	Camera
19.99	Camera

Decomposition in General



Let: S_1 = projection of R on $A_1, \dots, A_n, B_1, \dots, B_m$

S_2 = projection of R on $A_1, \dots, A_n, C_1, \dots, C_p$

The decomposition is called lossless if $R = S_1 \bowtie S_2$

Fact: If $A_1, \dots, A_n \rightarrow B_1, \dots, B_m$ then the decomposition is lossless

It follows that every BCNF decomposition is lossless 48

DBs and QM (off topic)

Lossless joins are related to quantum mechanics:

- Tables of measurement outcomes from one object
 - each measures two properties at once
 - e.g., {height, hair color}, {hair color, weight}, etc.
- Each 2-property table should be a projection of a table with all properties
- Somehow this does not happen for QM systems

DBs and QM (off topic)

measurement

A	B
0	0
0	1
1	0
1	1

A'	B
0	1
1	0
1	1

A'	B'
0	0
0	1
1	0

A	B'
0	1
1	0
1	1

(measuring a population of objects, say, so we can get different outcomes each time)

natural join all =

A	A'	B	B'
0	0	1	1
1	0	1	0
1	1	1	0
1	0	1	1
1	1	0	0

projection onto (A,B)
does not include (0,0)!

Schema Refinements = Normal Forms

- 1st Normal Form = all tables are flat (no list values)
- 2nd Normal Form = obsolete
- Boyce Codd Normal Form = no bad FDs
- 3rd Normal Form = see book
 - BCNF is lossless but can cause lose ability to check some FDs without a join (see book 3.4.4)
 - 3NF fixes that (is lossless and dependency-preserving), but some tables might not be in BCNF – i.e., they may have redundancy anomalies