

# Introduction to Data Management

## CSE 344

Unit 6: Conceptual Design  
E/R Diagrams  
Integrity Constraints  
BCNF

(3 lectures)

# Introduction to Data Management

## CSE 344

### E/R Diagrams

# Class Overview

- Unit 1: Intro
- Unit 2: Relational Data Models and Query Languages
- Unit 3: Non-relational data
- Unit 4: RDMBS internals and query optimization
- Unit 5: Parallel query processing
- Unit 6: DBMS usability, conceptual design
  - E/R diagrams
  - Schema normalization
- Unit 7: Transactions
- Unit 8: Advanced topics (time permitting)

# Database Design

What it is:

- Starting from scratch, design the database schema: relation, attributes, keys, foreign keys, constraints etc

Why it's hard

- The database will be in operation for a very long time (years). Updating the schema while in production is very expensive (why?)



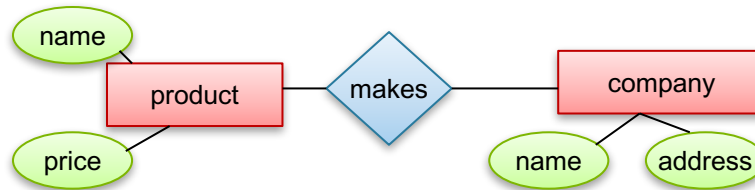
# Database Design

- Consider issues such as:
  - What entities to model
  - How entities are related
  - What constraints exist in the domain
- Several formalisms exists
  - We discuss E/R diagrams
  - UML, model-driven architecture
- Reading: Sec. 4.1-4.6



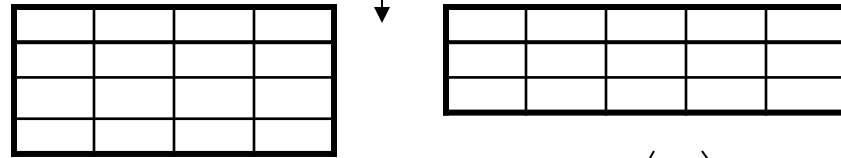
# Database Design Process

Conceptual Model:



Relational Model:

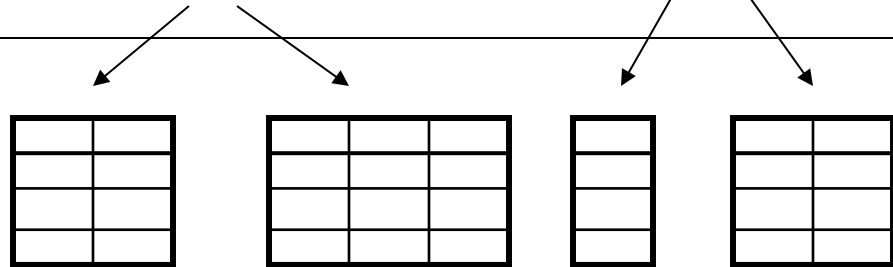
Tables + constraints  
And also functional dep.



Normalization:

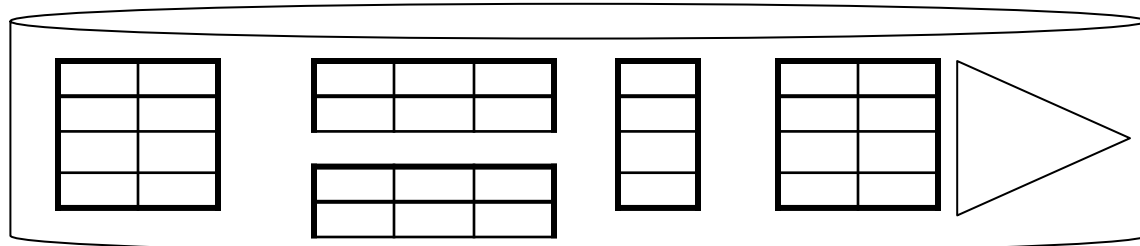
Eliminates anomalies

Conceptual Schema



Physical storage details

Physical Schema



# Entity / Relationship Diagrams

- Entity set = a class
  - An entity = an object

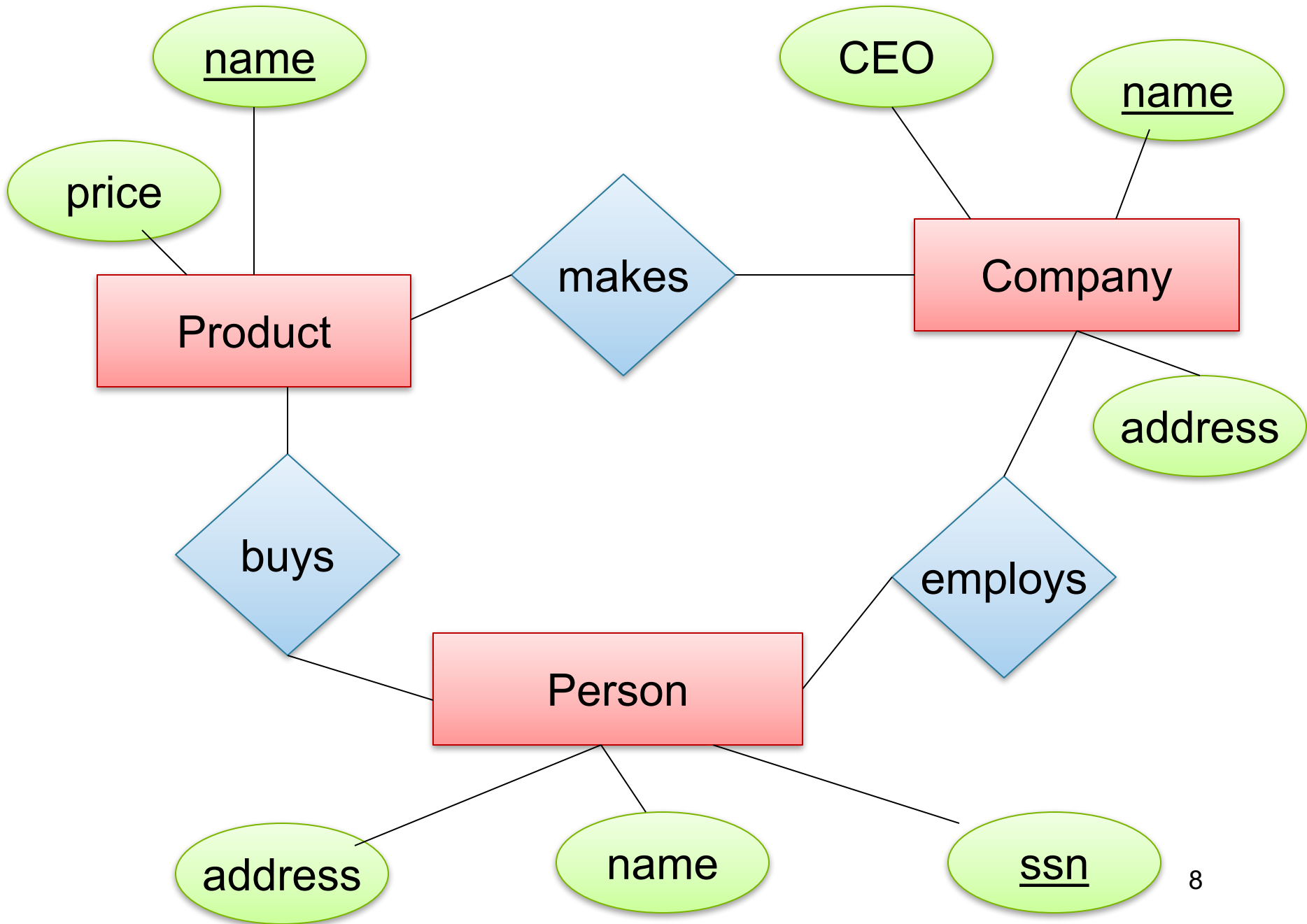


- Attribute



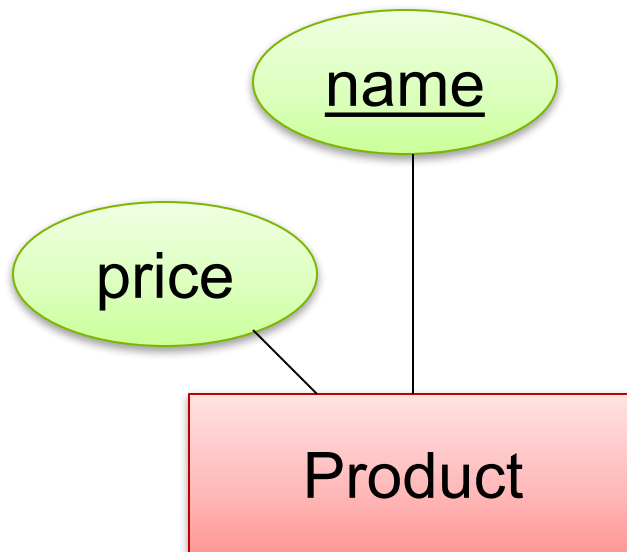
- Relationship





# Keys in E/R Diagrams

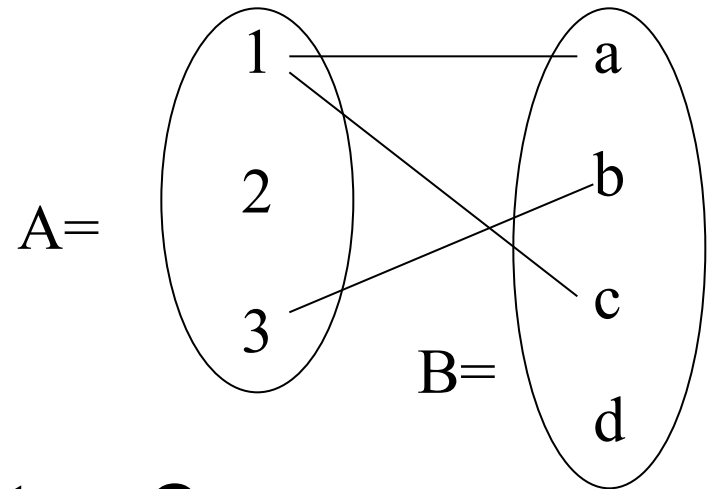
- Every entity set must have a key



# What is a Relation ?

- A mathematical definition:
  - if A, B are sets, then a relation R is a subset of  $A \times B$

- $A = \{1, 2, 3\}$ ,  $B = \{a, b, c, d\}$ ,  
 $A \times B = \{(1, a), (1, b), \dots, (3, d)\}$   
 $R = \{(1, a), (1, c), (3, b)\}$

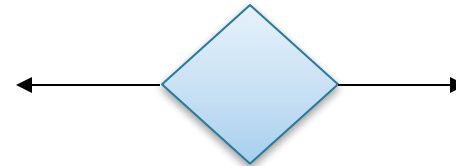
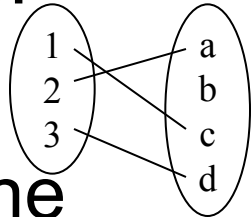


- **makes** is a subset of **Product**  $\times$  **Company**:

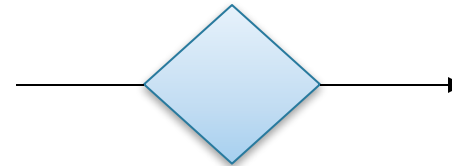
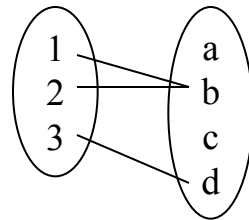


# Multiplicity of E/R Relations

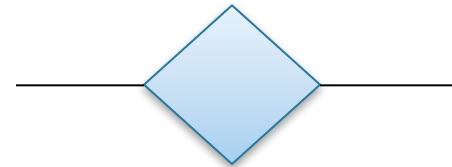
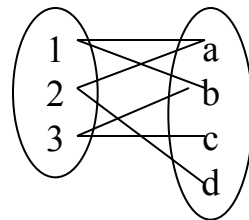
- one-one:

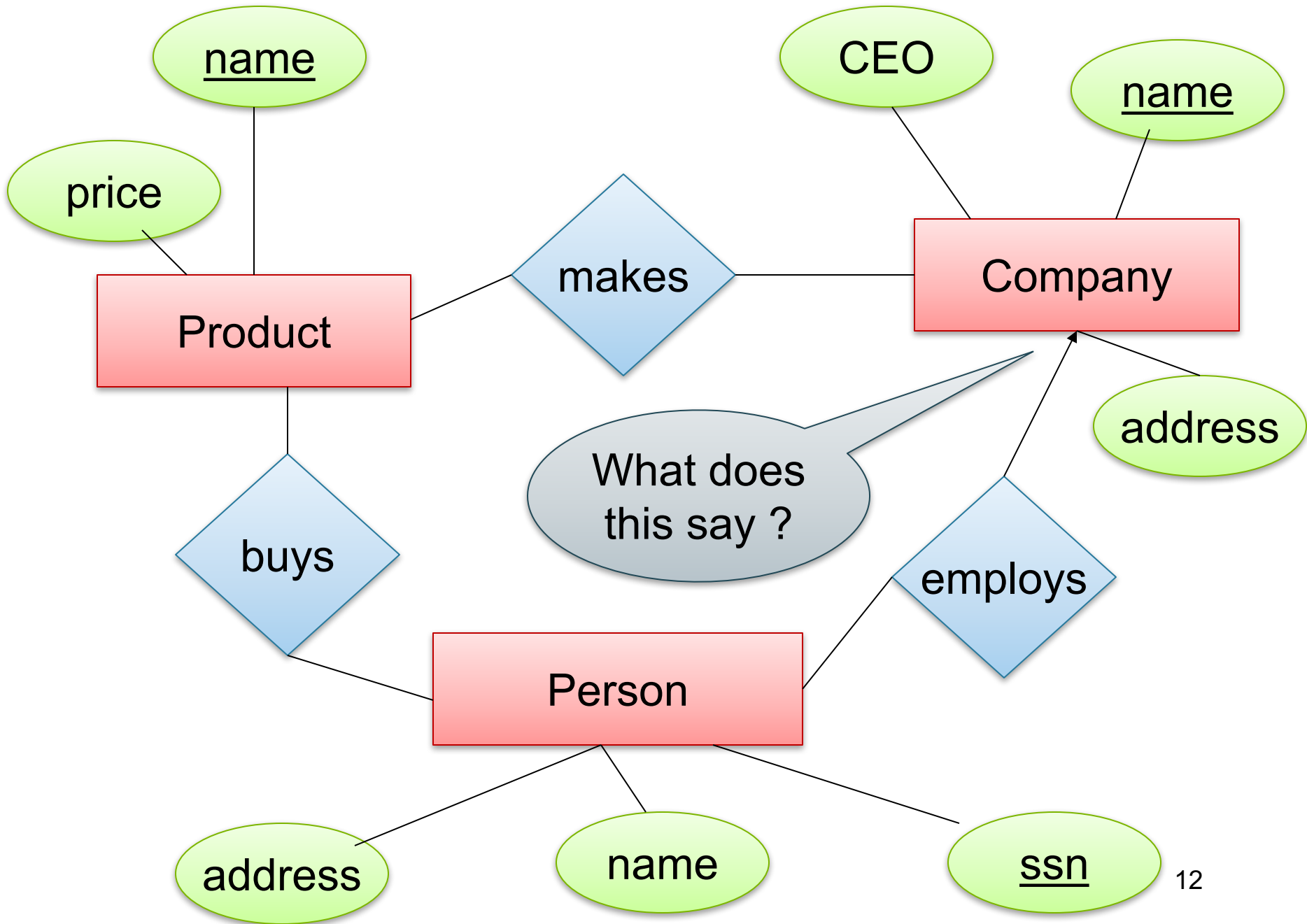


- many-one



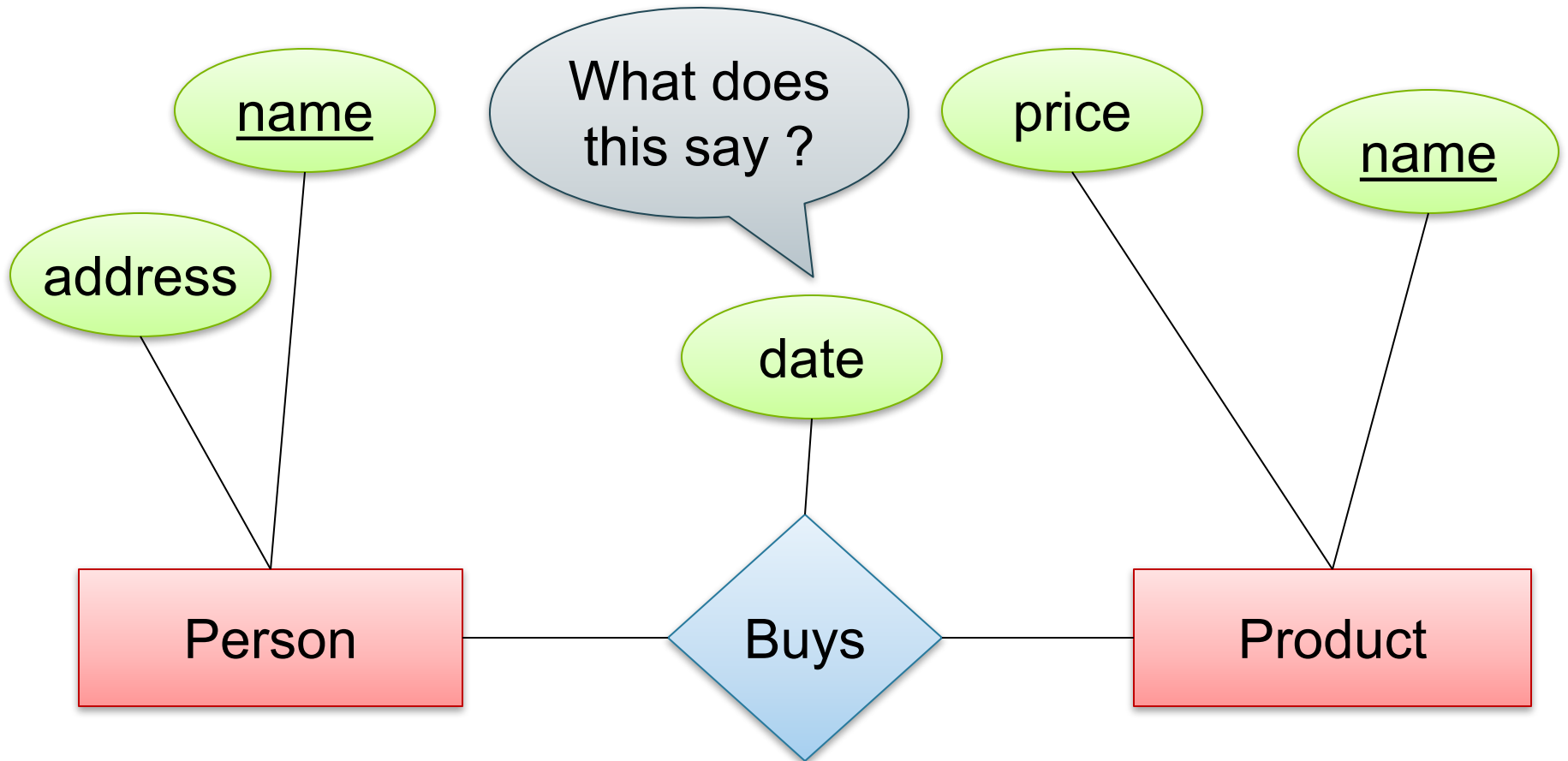
- many-many





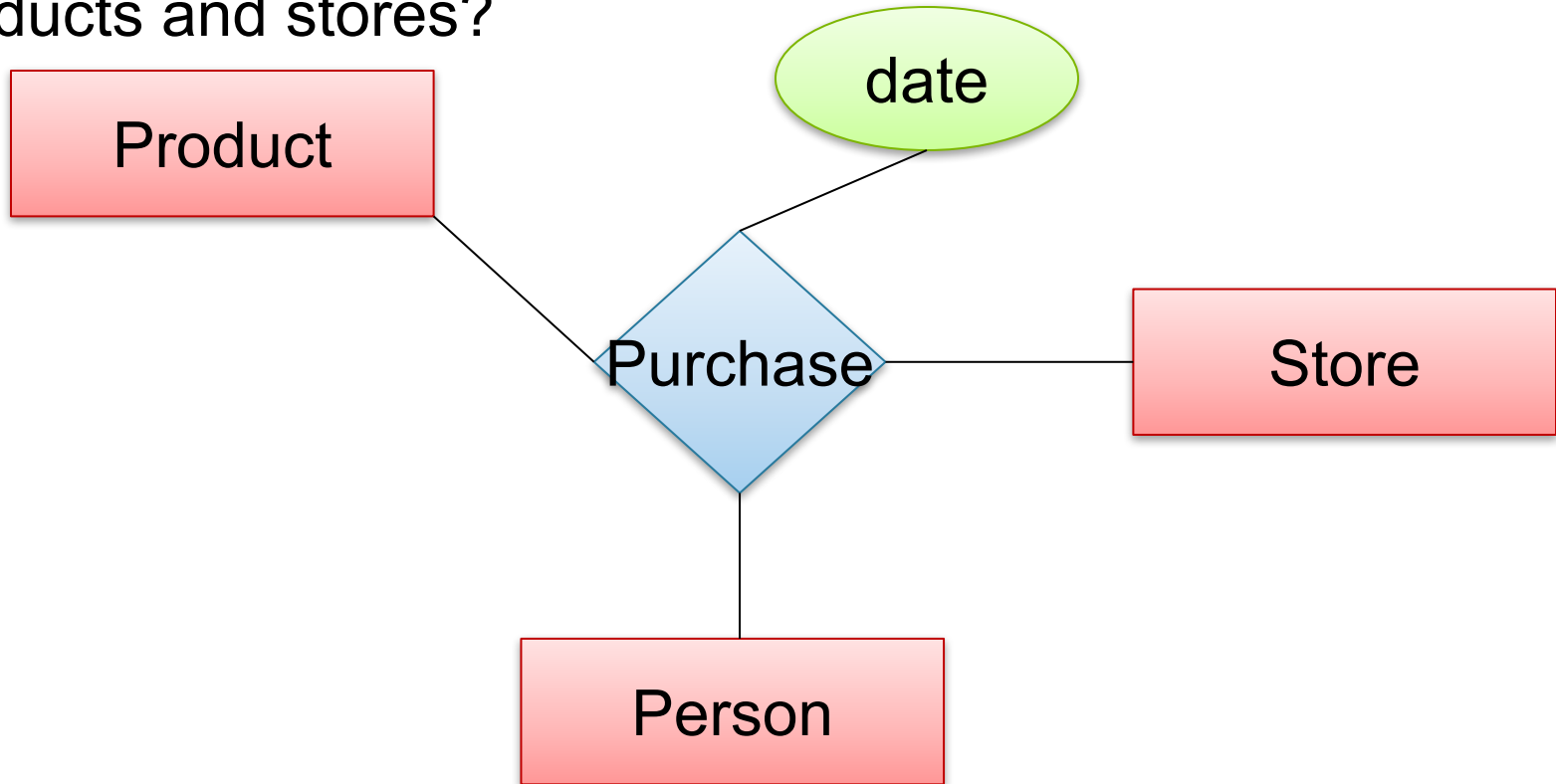


# Attributes on Relationships



# Multi-way Relationships

How do we model a purchase relationship between buyers, products and stores?

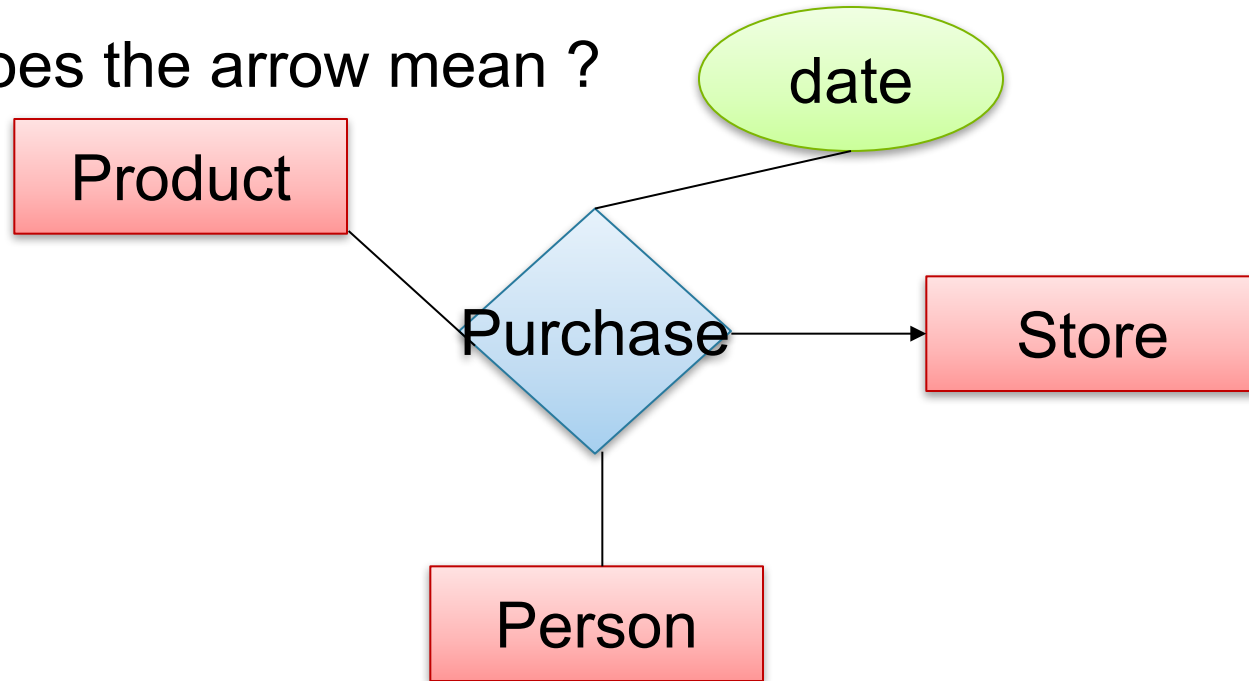


Can still model as a mathematical set (How?)

As a set of triples  $\subseteq \text{Person} \times \text{Product} \times \text{Store}$

# Arrows in Multiway Relationships

**Q:** What does the arrow mean ?

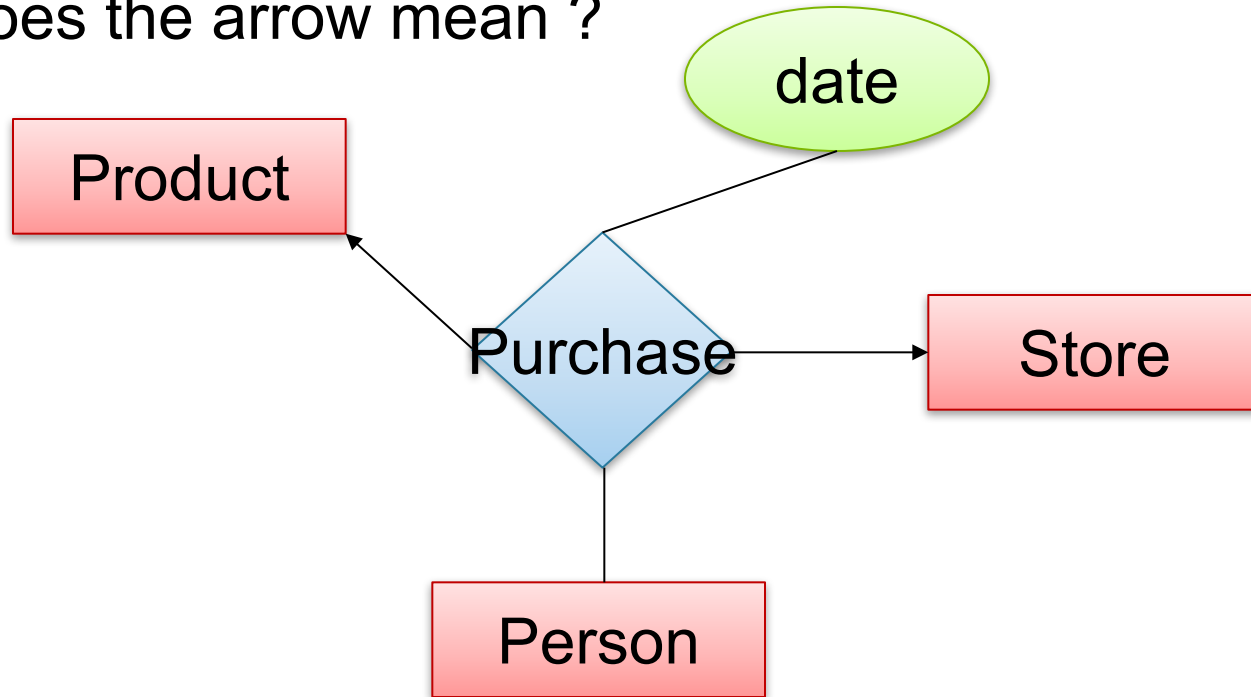


**A:** Any person buys a given product from at most one store

[Fine print: Arrow pointing to E means that if we select one entity from each of the other entity sets in the relationship, those entities are related to at most one entity in E]

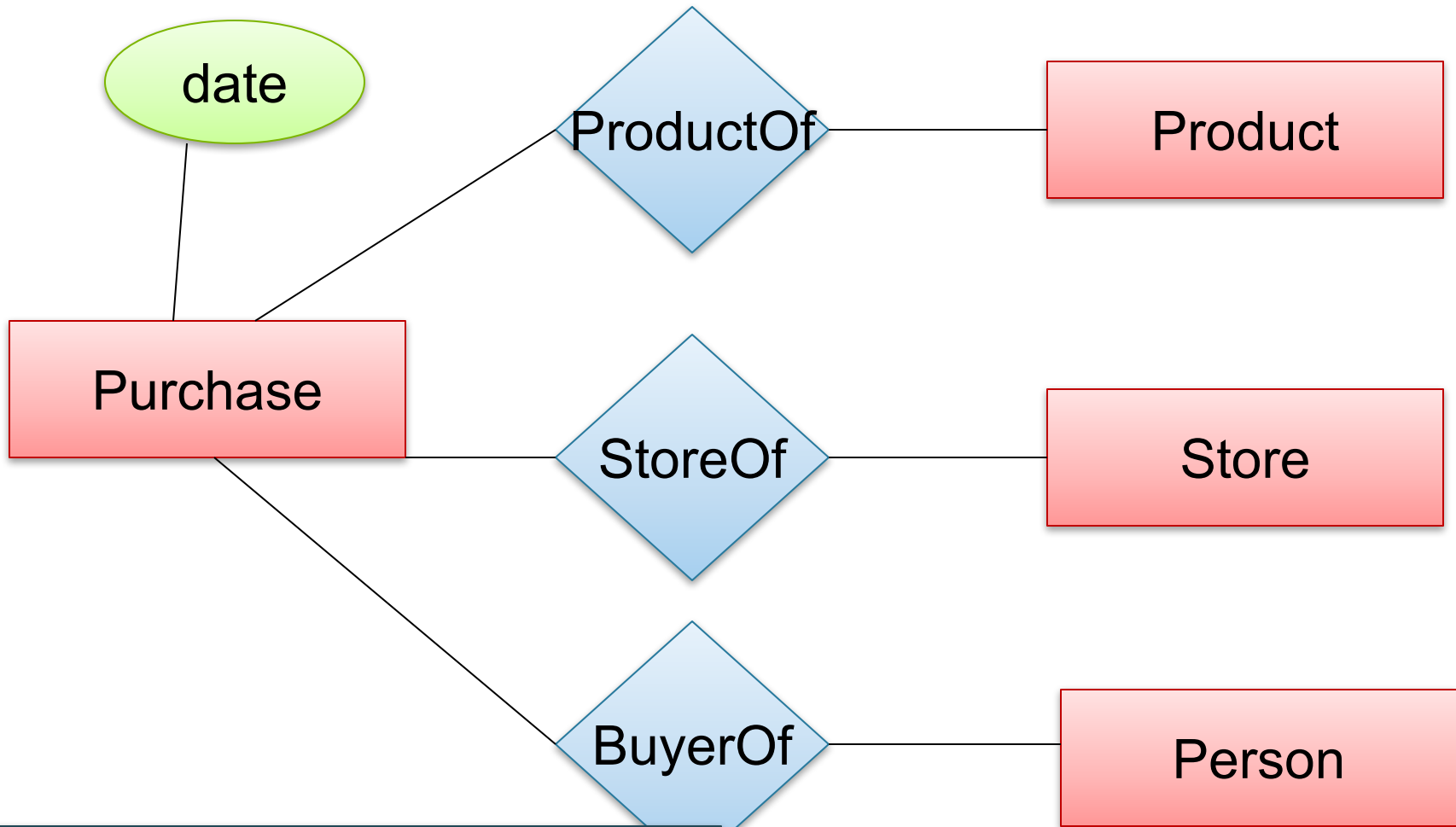
# Arrows in Multiway Relationships

**Q:** What does the arrow mean ?



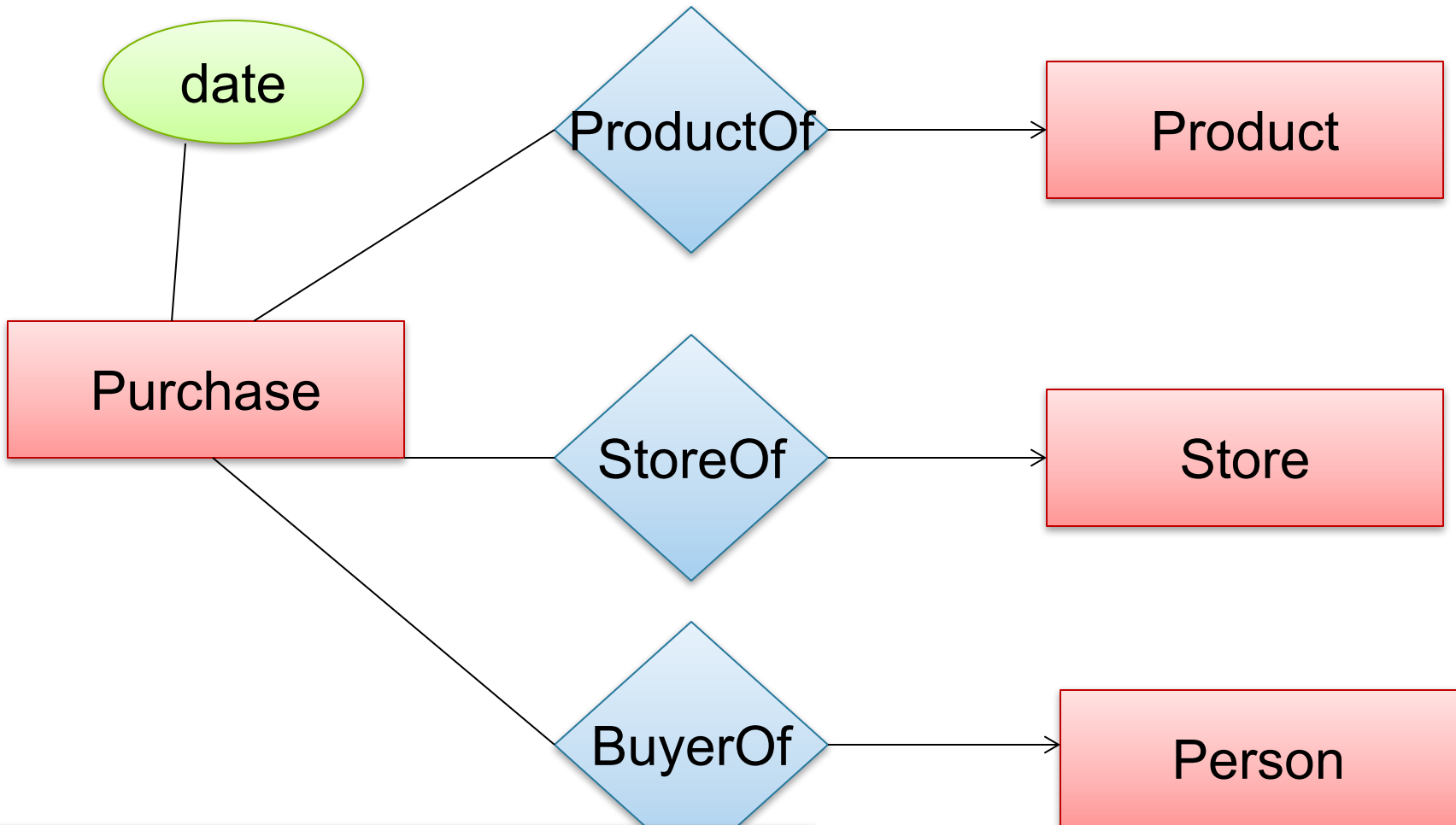
**A:** Any person buys a given product from at most one store  
AND every store sells to every person at most one product

# Converting Multi-way Relationships to Binary



Arrows go in which direction?

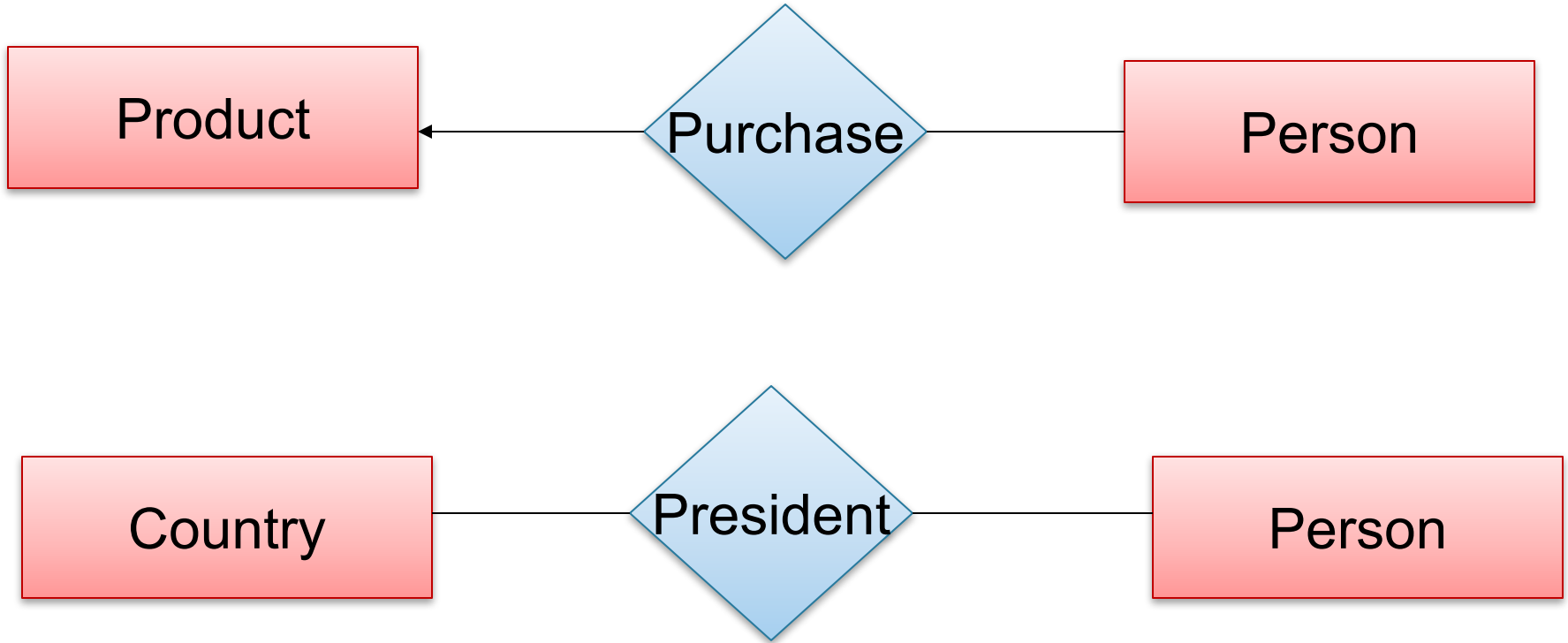
# Converting Multi-way Relationships to Binary



Make sure you understand why!

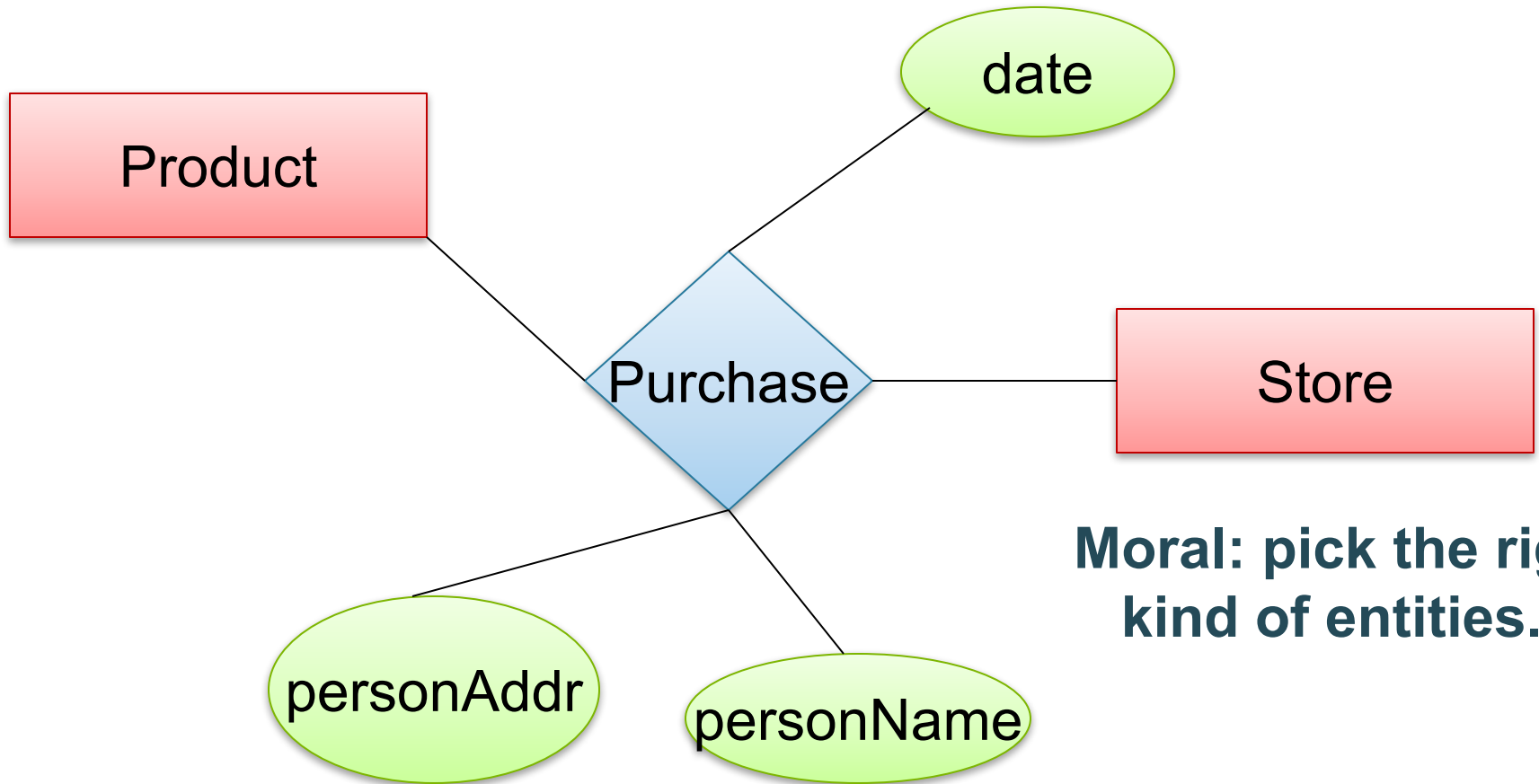
# 3. Design Principles

What's wrong?



**Moral: Be faithful to the specifications of the application!**

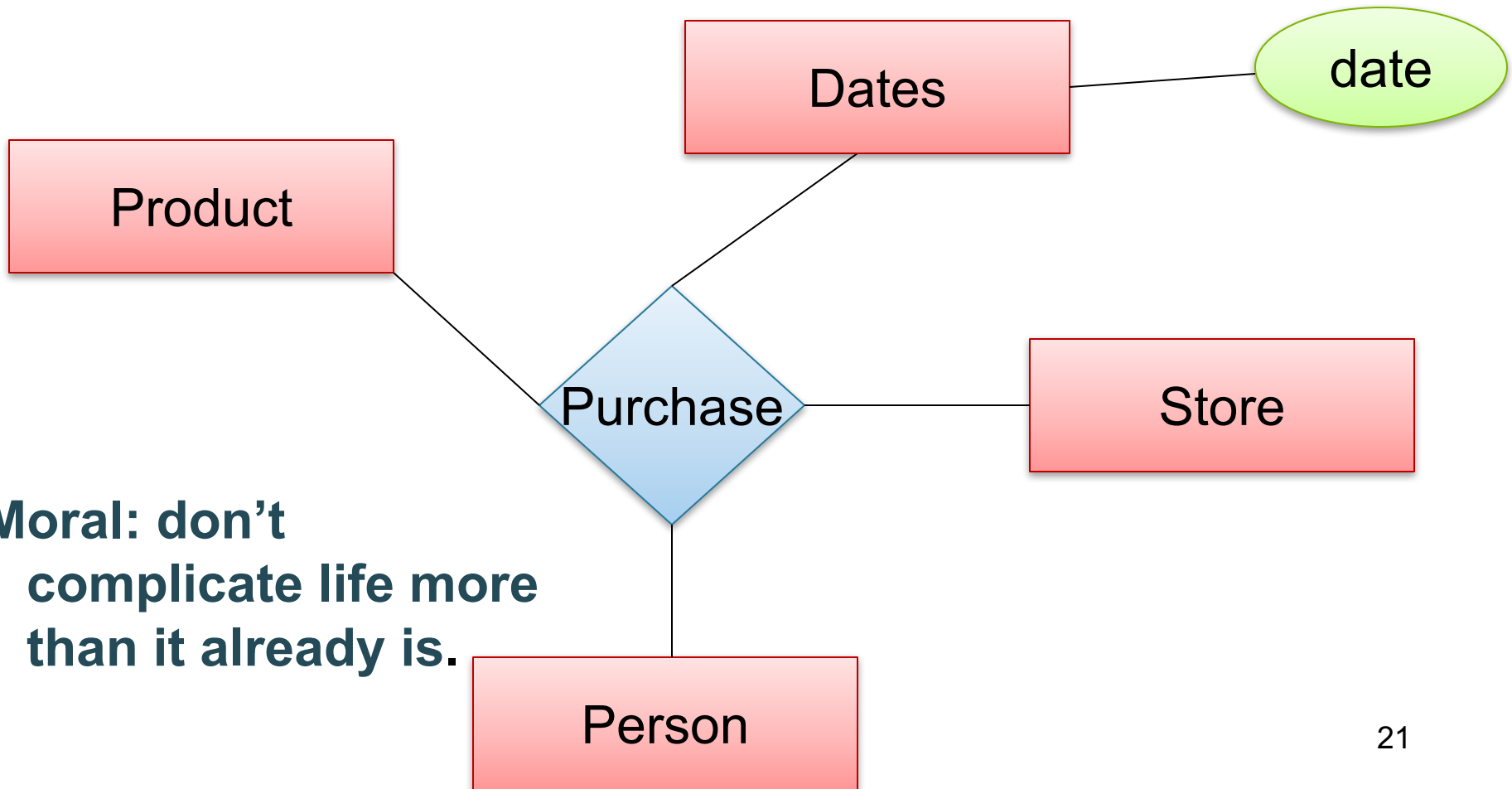
# Design Principles: What's Wrong?



**Moral: pick the right  
kind of entities.**



# Design Principles: What's Wrong?

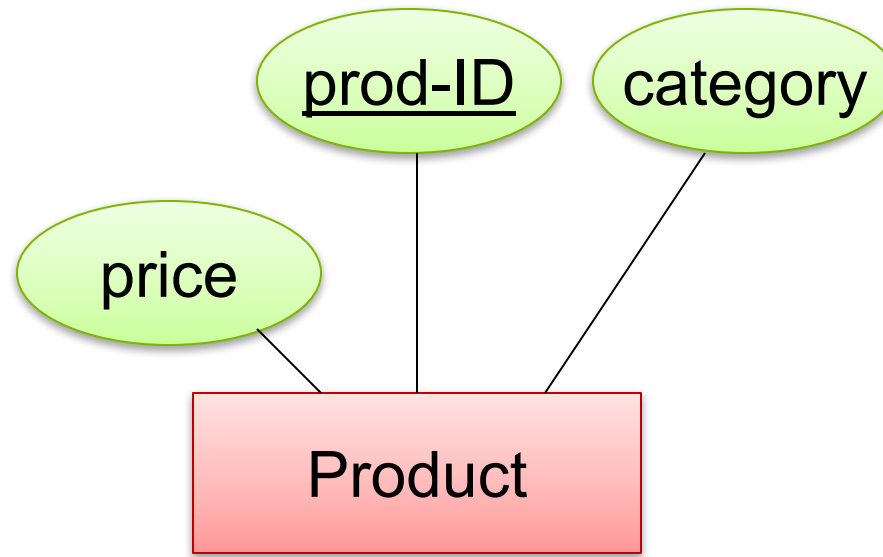


**Moral: don't  
complicate life more  
than it already is.**

# From E/R Diagrams to Relational Schema

- Entity set  $\rightarrow$  relation
- Relationship  $\rightarrow$  relation

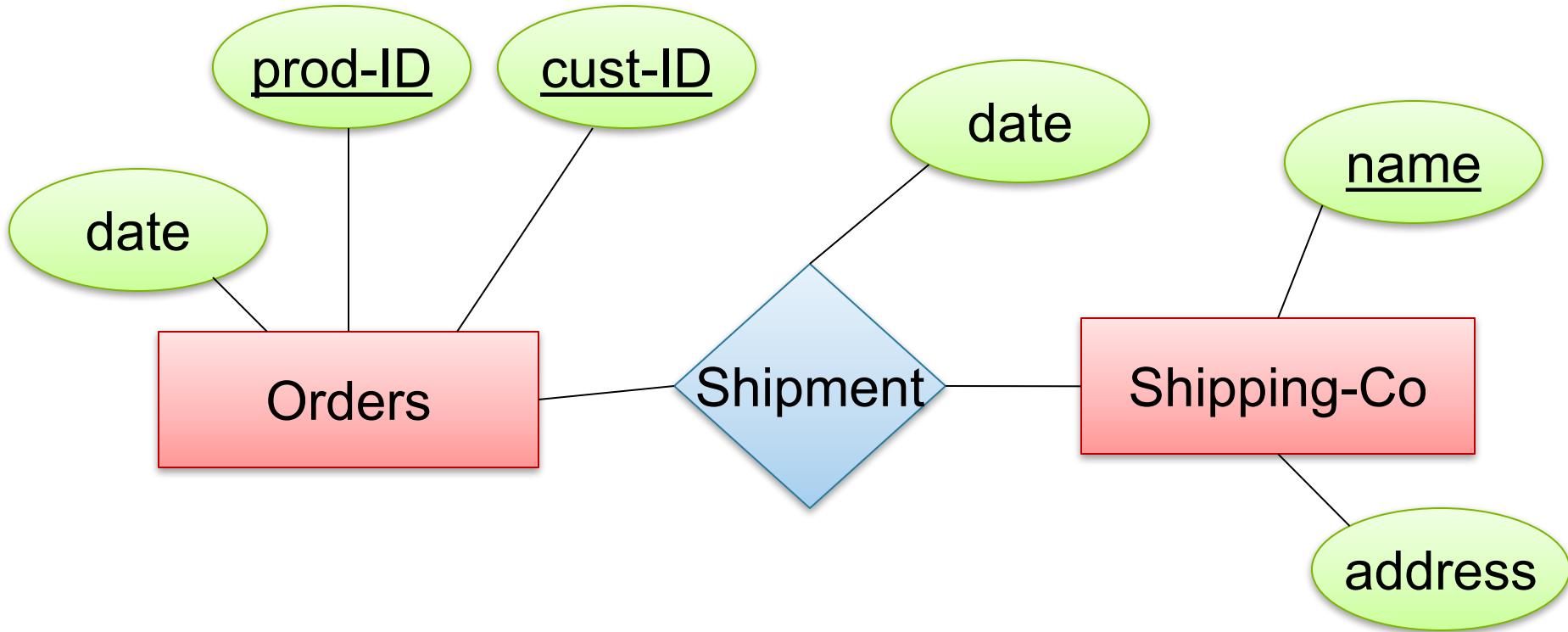
# Entity Set to Relation



**Product**(prod-ID, category, price)

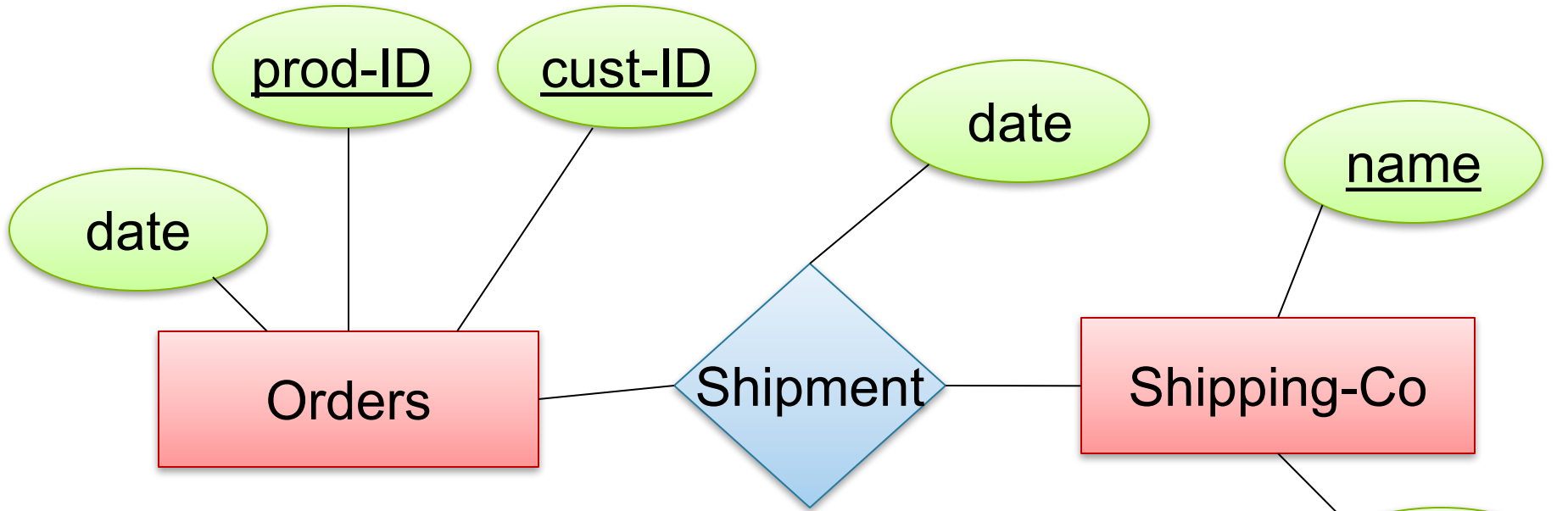
<u>prod-ID</u>	category	price
Gizmo55	Camera	99.99
Pokemn19	Toy	29.99

# N-N Relationships to Relations



Represent this in relations

# N-N Relationships to Relations



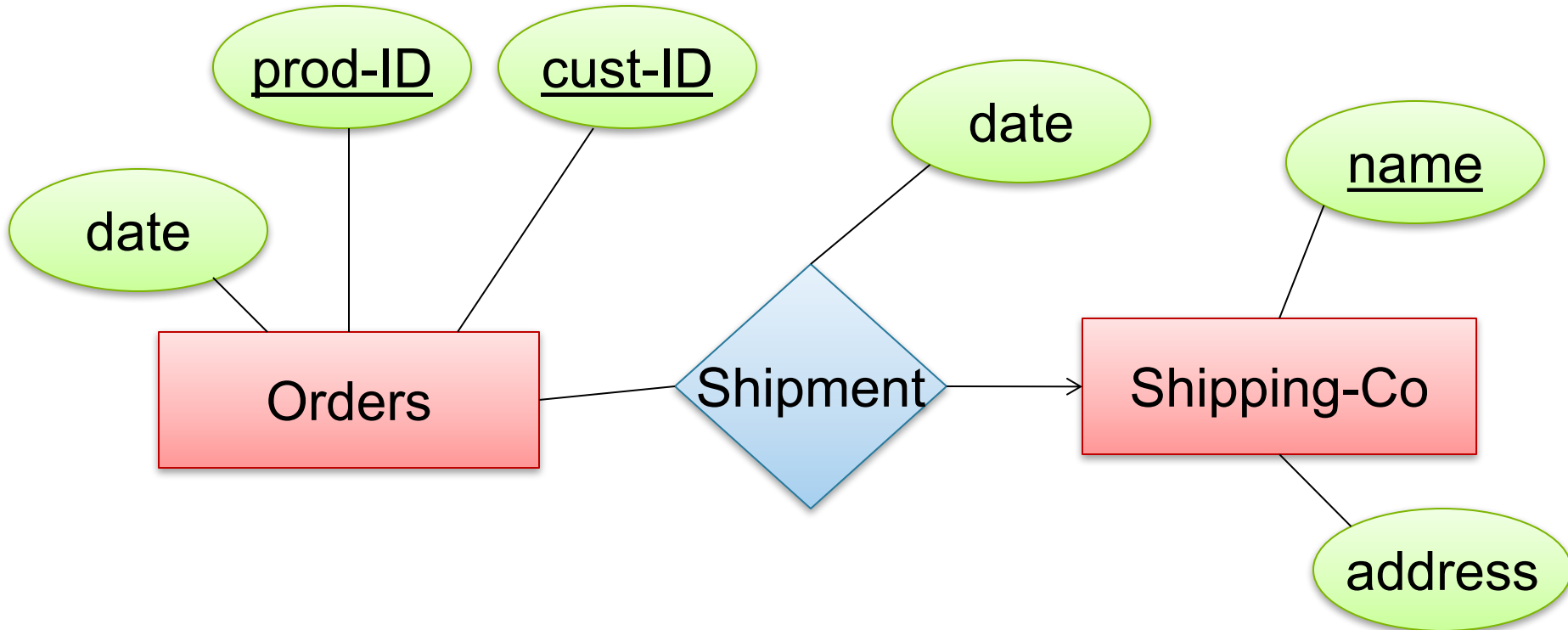
**Orders**(prod-ID, cust-ID, date)

**Shipment**(prod-ID, cust-ID, name, date)

**Shipping-Co**(name, address)

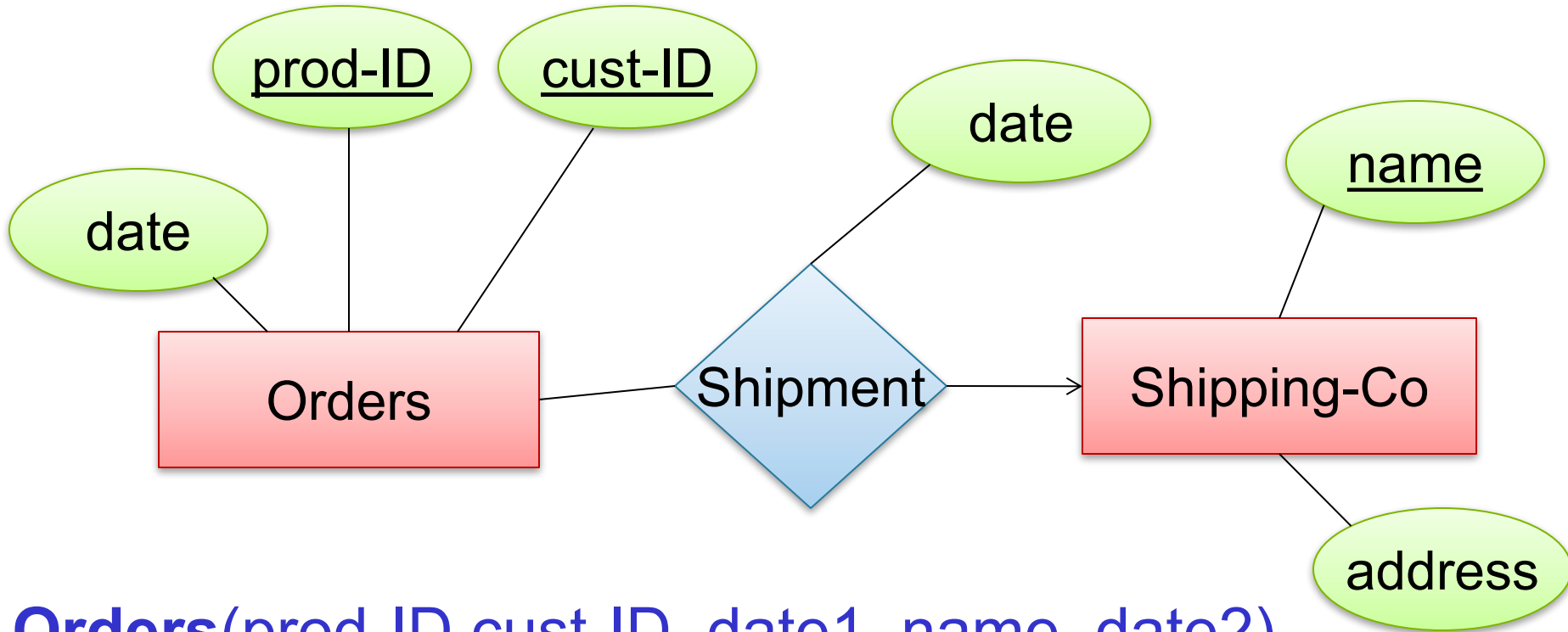
<u>prod-ID</u>	<u>cust-ID</u>	<u>name</u>	date
Gizmo55	Joe12	UPS	4/10/2011
Gizmo55	Joe12	FEDEX	4/9/2011

# N-1 Relationships to Relations



Represent this in relations

# N-1 Relationships to Relations

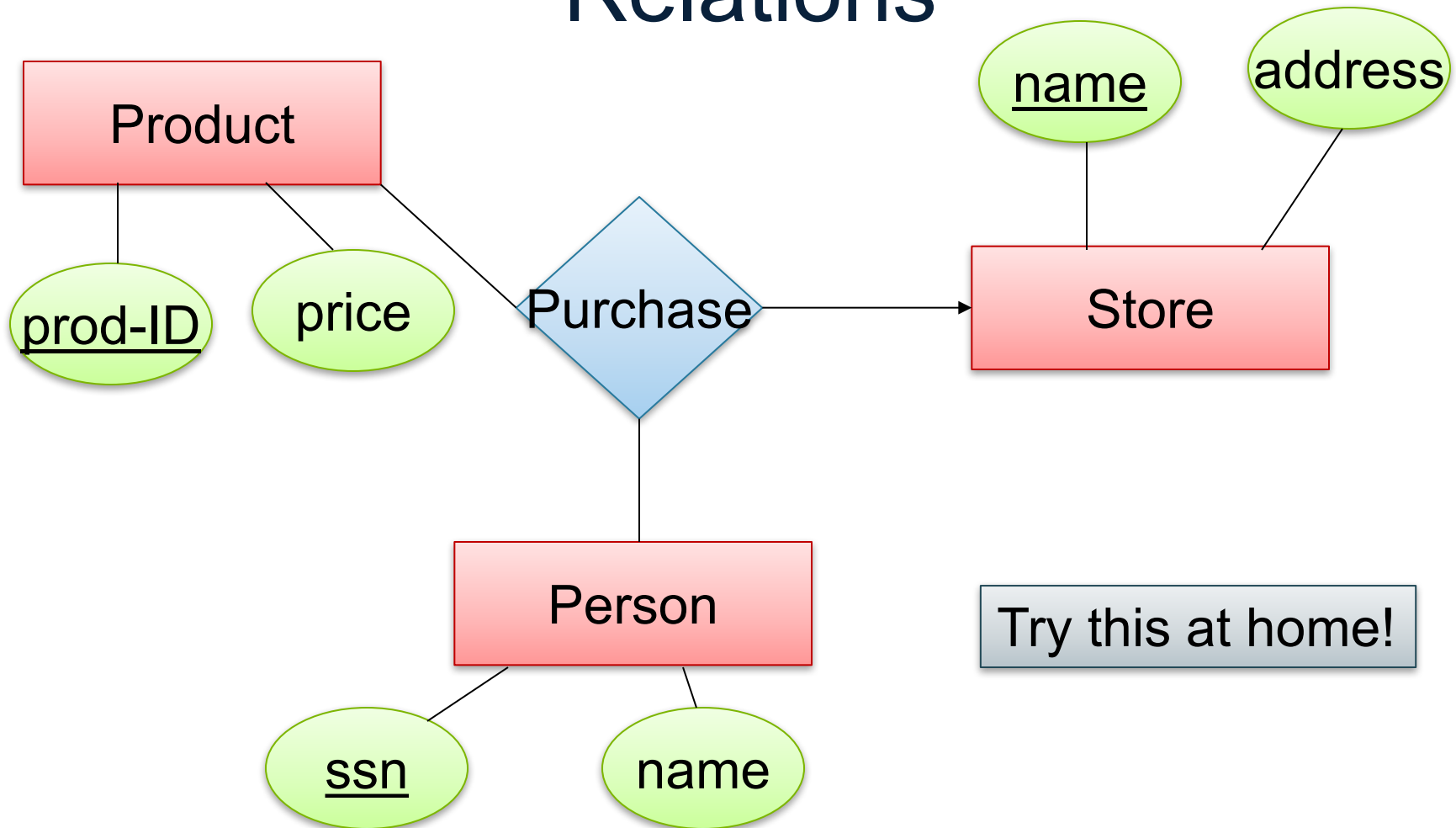


**Orders**(prod-ID, cust-ID, date1, name, date2)

**Shipping-Co**(name, address)

Remember: no separate relations for many-one relationship

# Multi-way Relationships to Relations



Try this at home!

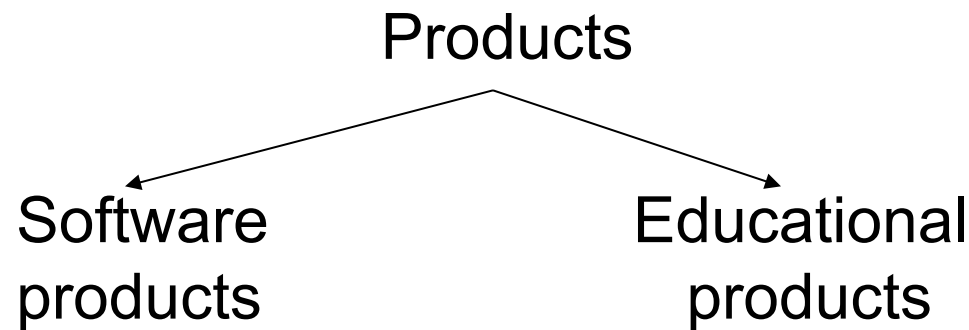
**Purchase(prod-ID, ssn, name)**



# Modeling Subclasses

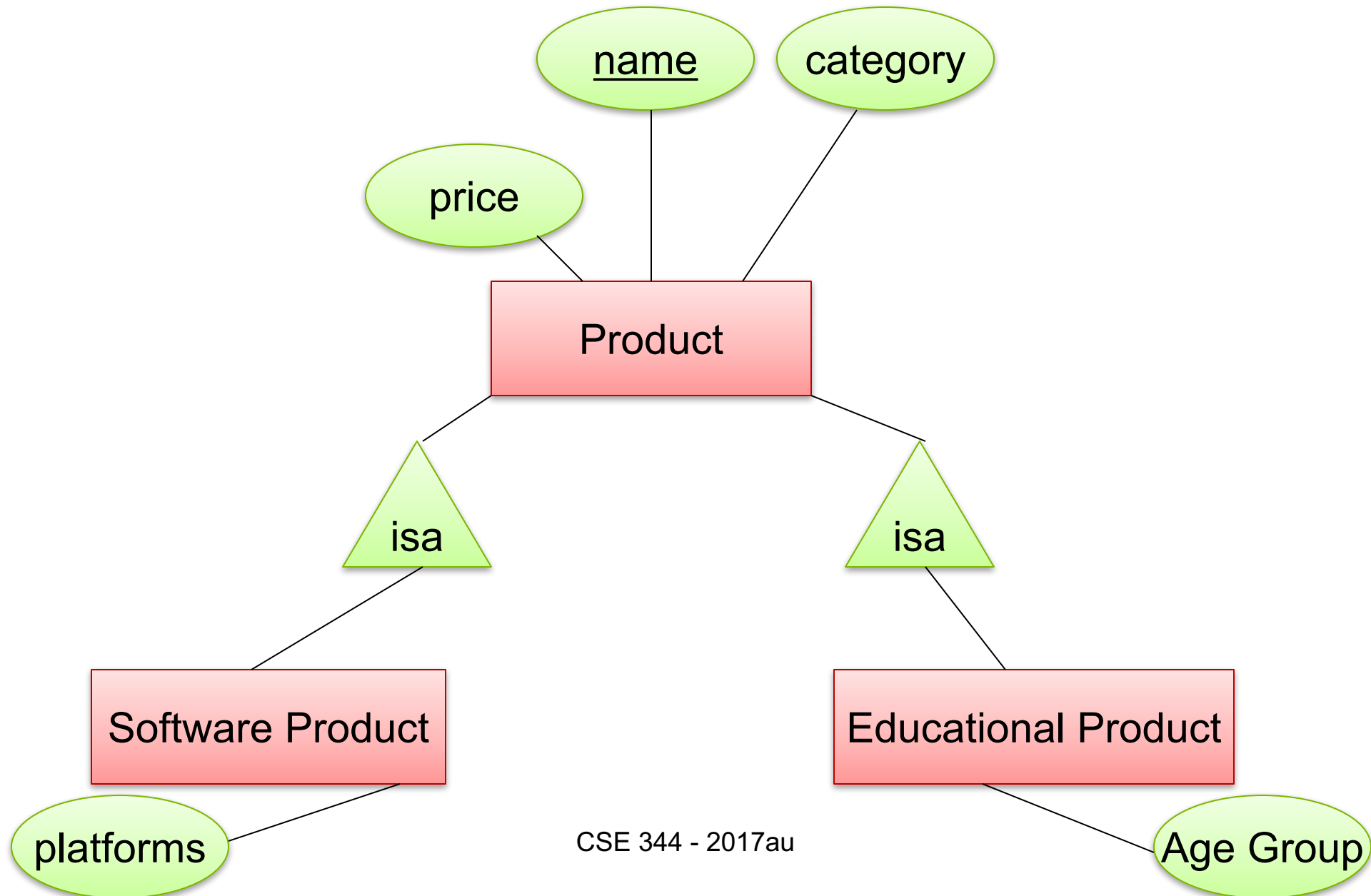
Some objects in a class may be special

- define a new class
- better: define a *subclass*



So --- we define subclasses in E/R

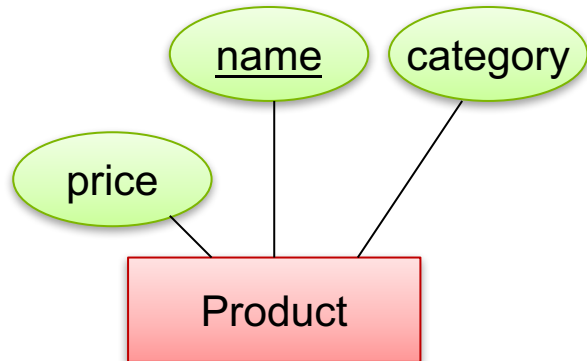
# Subclasses



# Subclasses to Relations

## Product

<u>Name</u>	Price	Category
Gizmo	99	gadget
Camera	49	photo
Toy	39	gadget



## Sw.Product

<u>Name</u>	platforms
Gizmo	unix

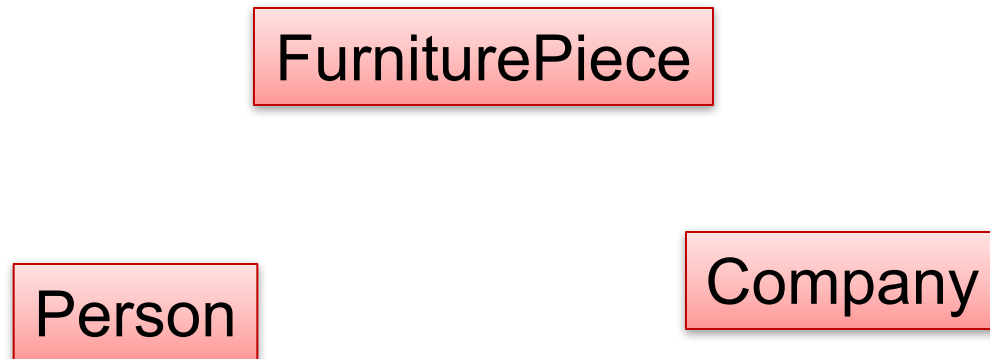


## Ed.Product

<u>Name</u>	Age Group
Gizmo	toddler
Toy	retired

Other ways to convert are possible

# Modeling Union Types with Subclasses

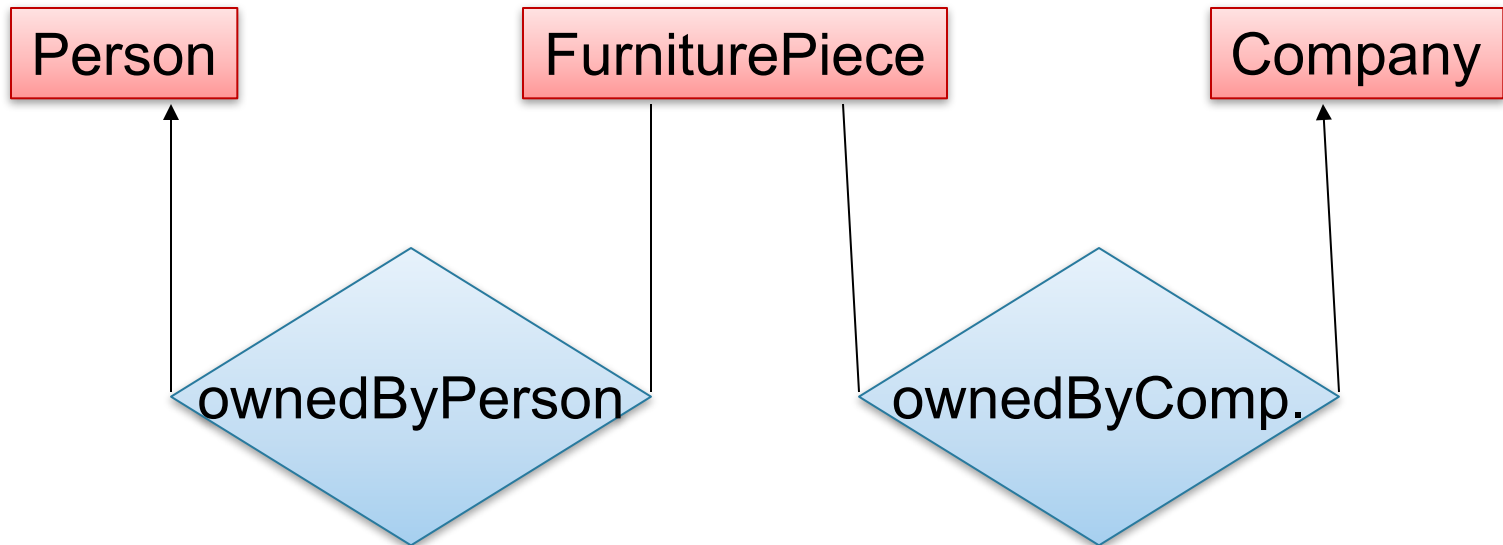


Say: each piece of furniture is owned either by a person or by a company

# Modeling Union Types with Subclasses

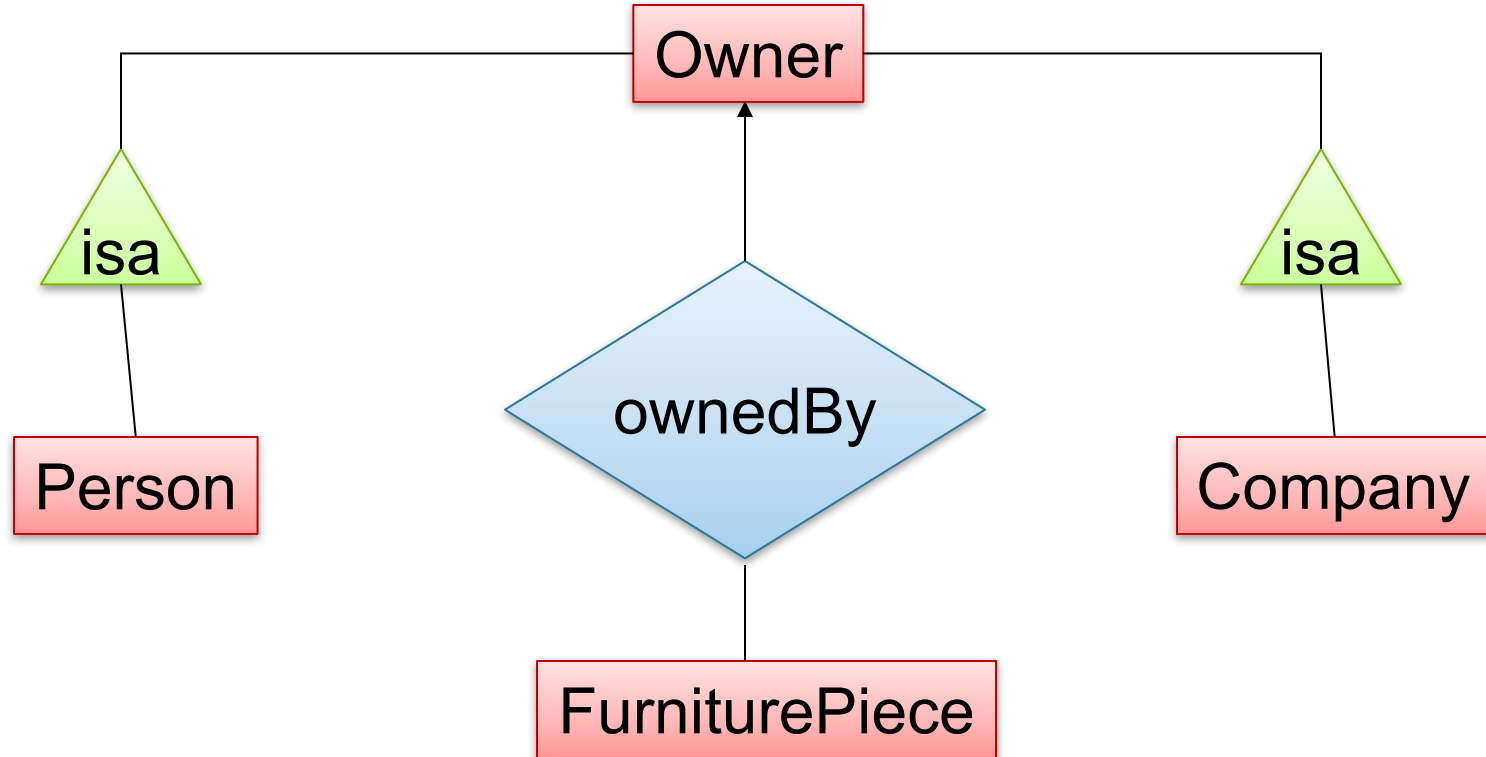
Say: each piece of furniture is owned either by a person or by a company

Solution 1. Acceptable but imperfect (What's wrong ?)



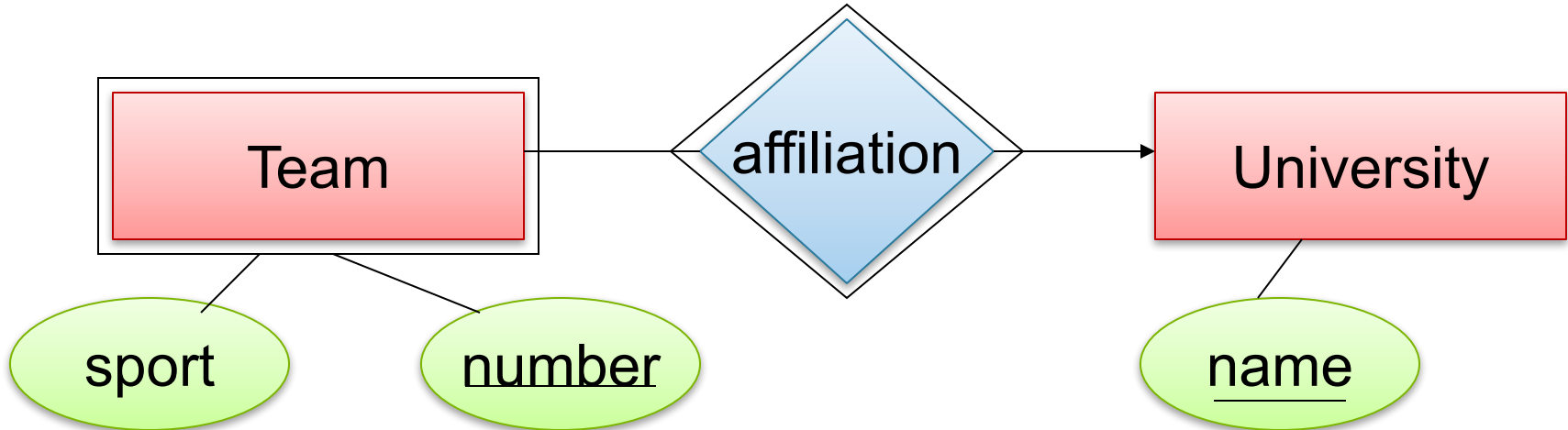
# Modeling Union Types with Subclasses

Solution 2: better, more laborious



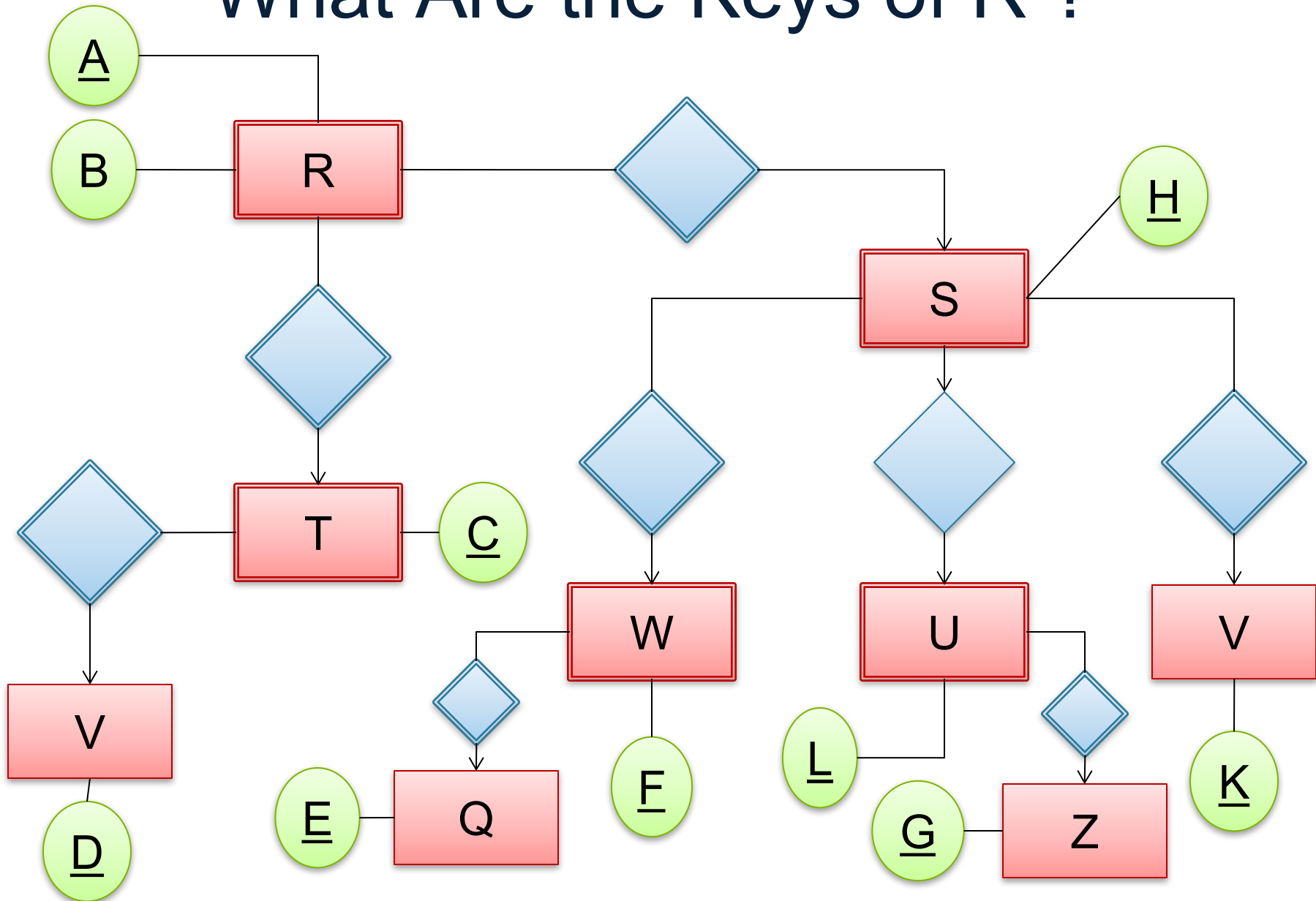
# Weak Entity Sets

Entity sets are weak when their key comes from other classes to which they are related.



Team(sport, number, universityName)  
University(name)

# What Are the Keys of R ?





# Introduction to Data Management

## CSE 344

### Integrity Constraints

# Integrity Constraints Motivation

An integrity constraint is a condition specified on a database schema that restricts the data that can be stored in an instance of the database.

- ICs help prevent entry of incorrect information
- How? DBMS enforces integrity constraints
  - Allows only legal database instances (i.e., those that satisfy all constraints) to exist
  - Ensures that all necessary checks are always performed and avoids duplicating the verification logic in each application

# Constraints in E/R Diagrams

Finding constraints is part of the modeling process.  
Commonly used constraints:

**Keys:** social security number uniquely identifies a person.

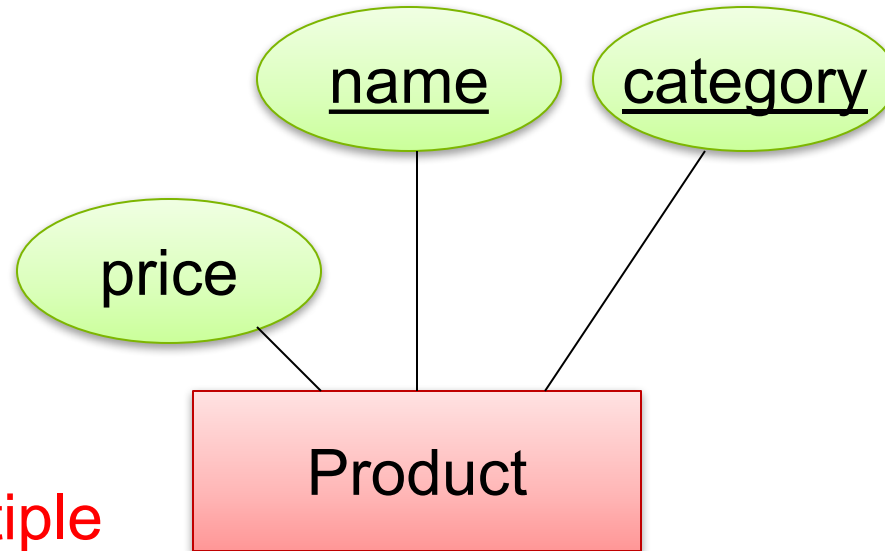
**Single-value constraints:** a person can have only one father.

**Referential integrity constraints:** if you work for a company, it must exist in the database.

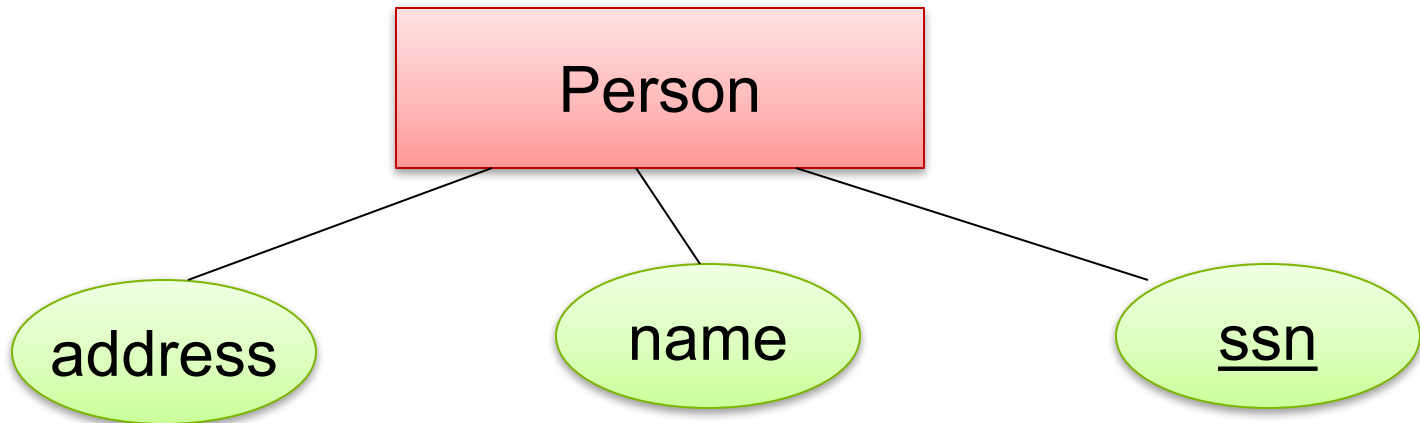
**Other constraints:** peoples' ages are between 0 and 150.

# Keys in E/R Diagrams

Underline:



No formal way  
to specify multiple  
keys in E/R diagrams



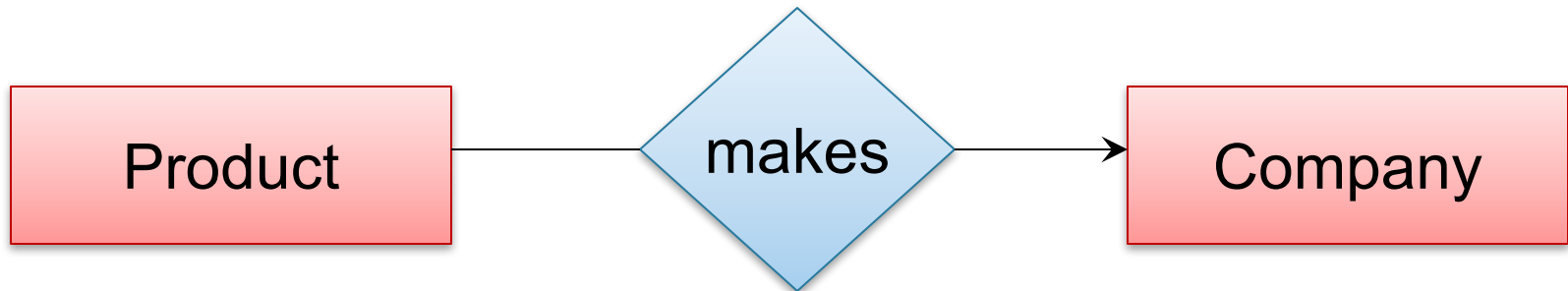
# Single Value Constraints



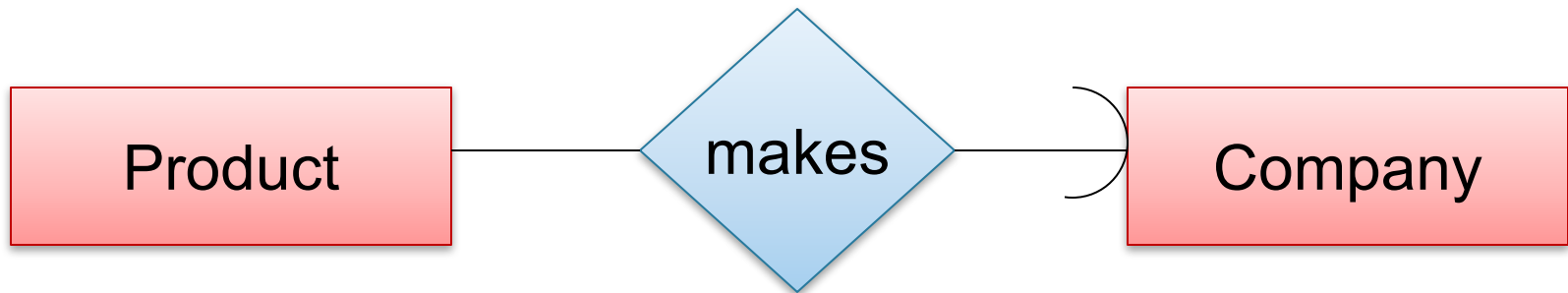
vs.



# Referential Integrity Constraints

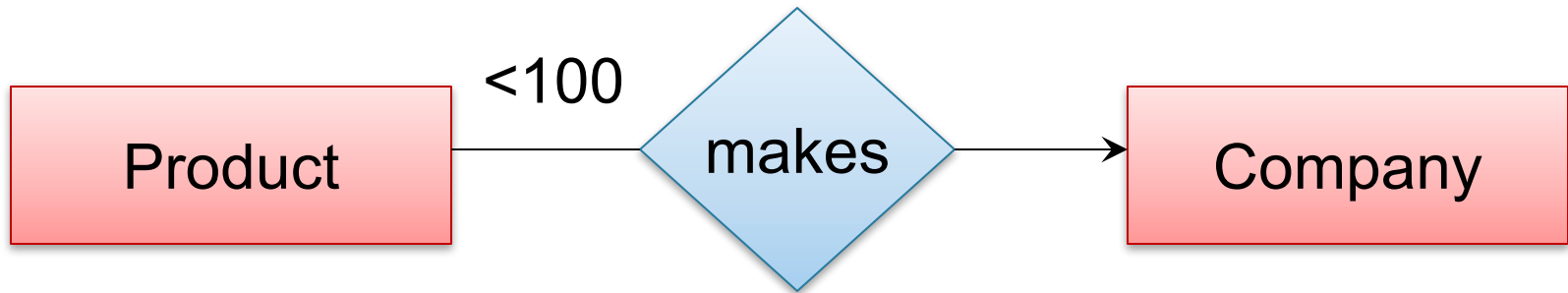


Each product made by at most one company.  
Some products made by no company



Each product made by exactly one company.

# Other Constraints



Q: What does this mean ?

A: A Company entity cannot be connected by relationship to more than 99 Product entities

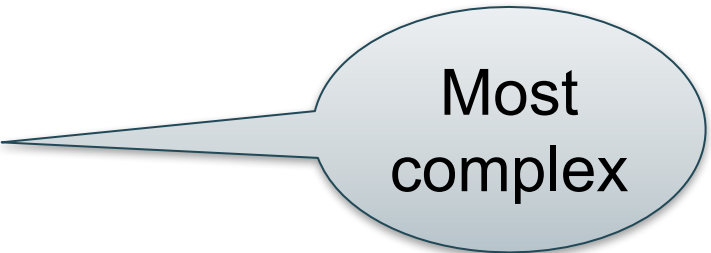
# Constraints in SQL

Constraints in SQL:

- **Keys, foreign keys**
- **Attribute-level** constraints
- **Tuple-level** constraints
- **Global** constraints: assertions



simplest



Most  
complex

- The more complex the constraint, the harder it is to check and to enforce



# Key Constraints

Product(name, category)

```
CREATE TABLE Product (  
    name CHAR(30) PRIMARY KEY,  
    category VARCHAR(20))
```

OR:

```
CREATE TABLE Product (  
    name CHAR(30),  
    category VARCHAR(20),  
    PRIMARY KEY (name))
```

# Keys with Multiple Attributes

Product(name, category, price)

```
CREATE TABLE Product (  
    name CHAR(30),  
    category VARCHAR(20),  
    price INT,  
    PRIMARY KEY (name, category))
```

Name	Category	Price
Gizmo	Gadget	10
Camera	Photo	20
Gizmo	Photo	30
<del>Gizmo</del>	<del>Gadget</del>	<del>40</del>

# Other Keys

```
CREATE TABLE Product (  
    productID CHAR(10),  
    name CHAR(30),  
    category VARCHAR(20),  
    price INT,  
    PRIMARY KEY (productID),  
    UNIQUE (name, category))
```

There is at most one **PRIMARY KEY**;  
there can be many **UNIQUE**

# Foreign Key Constraints

```
CREATE TABLE Purchase (  
  prodName CHAR(30)  
  REFERENCES Product(name),  
  date DATETIME)
```

Referential  
integrity  
constraints

prodName is a **foreign key** to Product(name)  
name must be a **key** in Product

May write  
just Product  
if name is PK

# Foreign Key Constraints

- Example with multi-attribute primary key

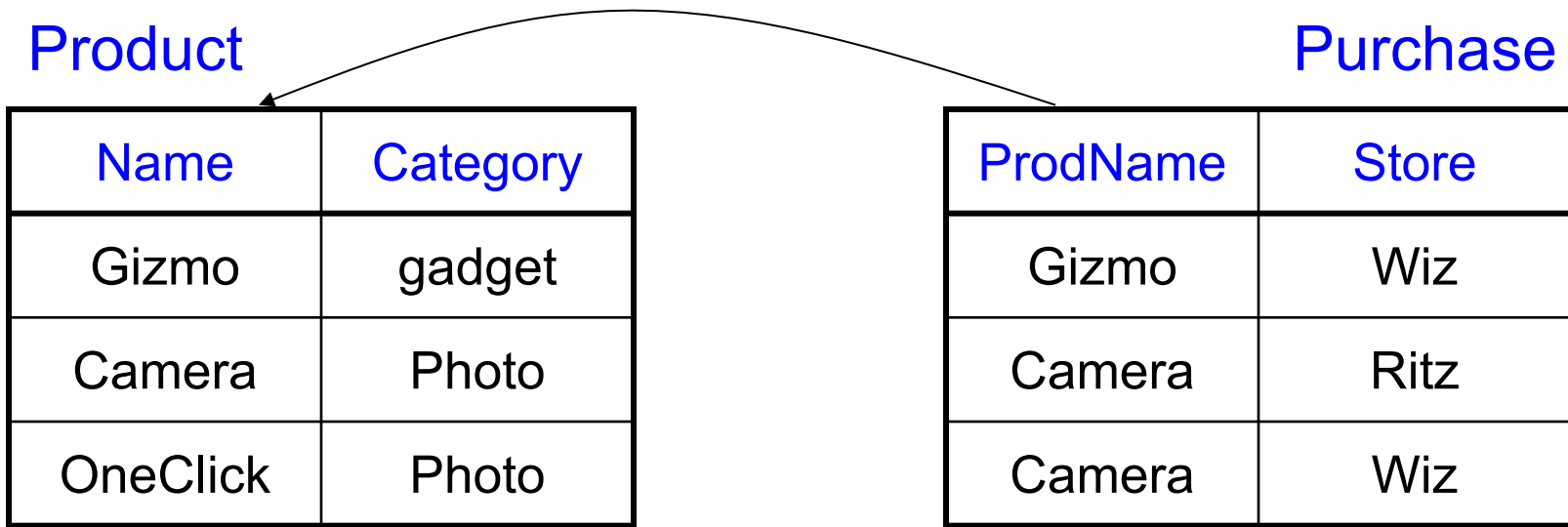
```
CREATE TABLE Purchase (  
    prodName CHAR(30),  
    category VARCHAR(20),  
    date DATETIME,  
    FOREIGN KEY (prodName, category)  
    REFERENCES Product(name, category)
```

- (name, category) must be a KEY in Product

# What happens when data changes?

Types of updates:

- In Purchase: insert/update
- In Product: delete/update



# What happens when data changes?

- SQL has three policies for maintaining referential integrity:
- NO ACTION reject violating modifications (default)
- CASCADE after delete/update do delete/update
- SET NULL set foreign-key field to NULL
- SET DEFAULT set foreign-key field to default value
  - need to be declared with column, e.g.,  
`CREATE TABLE Product (pid INT DEFAULT 42)`

# Maintaining Referential Integrity

```
CREATE TABLE Purchase (  
    prodName CHAR(30),  
    category VARCHAR(20),  
    date DATETIME,  
    FOREIGN KEY (prodName, category)  
    REFERENCES Product(name, category)  
    ON UPDATE CASCADE  
    ON DELETE SET NULL )
```

Product

Name	Category
Gizmo	gadget
Camera	Photo
OneClick	Photo

Purchase

ProdName	Category
Gizmo	Gizmo
Snap	Camera
EasyShoot	Camera



# Constraints on Attributes and Tuples

- Constraints on attributes:
  - NOT NULL** -- obvious meaning...
  - CHECK** condition -- any condition !
- Constraints on tuples
  - CHECK** condition

# Constraints on Attributes and Tuples

```
CREATE TABLE R (  
  A int NOT NULL,  
  B int CHECK (B > 50 and B < 100),  
  C varchar(20),  
  D int,  
  CHECK (C >= 'd' or D > 0))
```

# Constraints on Attributes and Tuples

```
CREATE TABLE Product (  
    productID CHAR(10),  
    name CHAR(30),  
    category VARCHAR(20),  
    price INT CHECK (price > 0),  
    PRIMARY KEY (productID),  
    UNIQUE (name, category))
```

# Constraints on Attributes and Tuples

What does this constraint do?

```
CREATE TABLE Purchase (  
  prodName CHAR(30)  
  CHECK (prodName IN  
    (SELECT Product.name  
     FROM Product),  
  date DATETIME NOT NULL)
```

What  
is the difference from  
Foreign-Key ?

# General Assertions

```
CREATE ASSERTION myAssert CHECK  
(NOT EXISTS(  
    SELECT Product.name  
    FROM Product, Purchase  
    WHERE Product.name = Purchase.prodName  
    GROUP BY Product.name  
    HAVING count(*) > 200) )
```

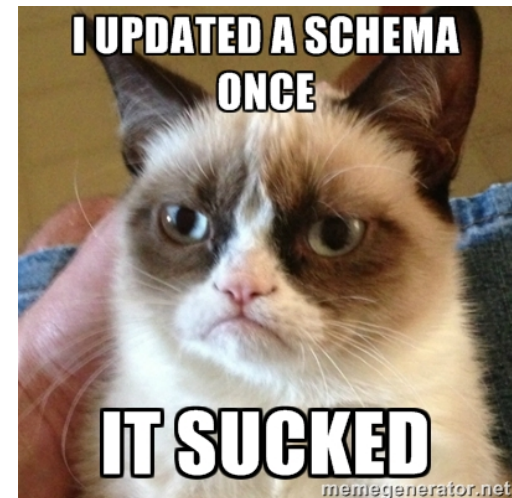
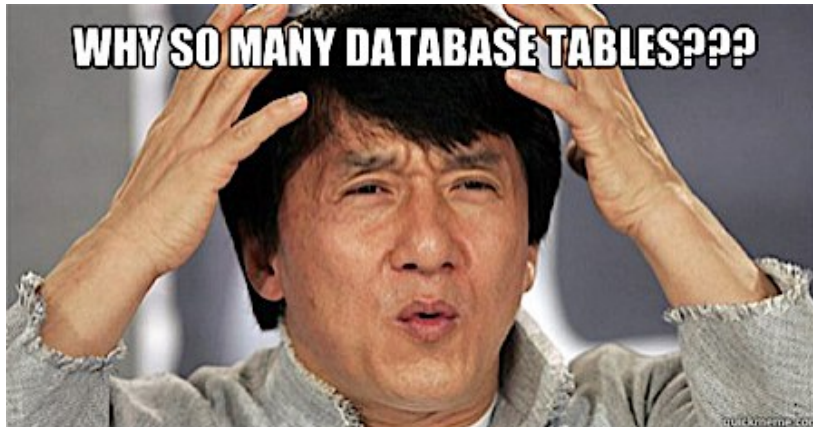
But most DBMSs do not implement assertions  
Because it is hard to support them efficiently  
Instead, they provide triggers

# Introduction to Data Management

## CSE 344

### Design Theory and BCNF

# 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


## Anomalies:

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

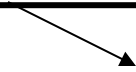
# Relation Decomposition

**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**  $A_1, \dots, A_m \rightarrow B_1, \dots, B_n$  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		$A_1$	...	$A_m$		$B_1$	...	$B_n$		
t										
t'										

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

# Example

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

<b>EmpID</b>	<b>Name</b>	<b>Phone</b>	<b>Position</b>
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

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

Position → Phone

# Example

<b>EmpID</b>	<b>Name</b>	<b>Phone</b>	<b>Position</b>
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 ?

# Buzzwords

- 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, we are **stating a constraint on R**

# Why bother with FDs?

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

## Anomalies:

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

# 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

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** is the set of attributes  $B$ , notated  $\{A_1, \dots, A_n\}^+$ ,  
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  $\{A, F\}^+$   $X = \{A, F, \quad \}$

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

# 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:

A, B	→	C
A, D	→	B
B	→	D

# 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
  - A 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 reducing to the minimal  $X$ 's to get the key

# Example

Product(name, price, category, color)

name, category → price  
category → color

What is 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 (name, category) is a 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 distinct 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 distinct 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?

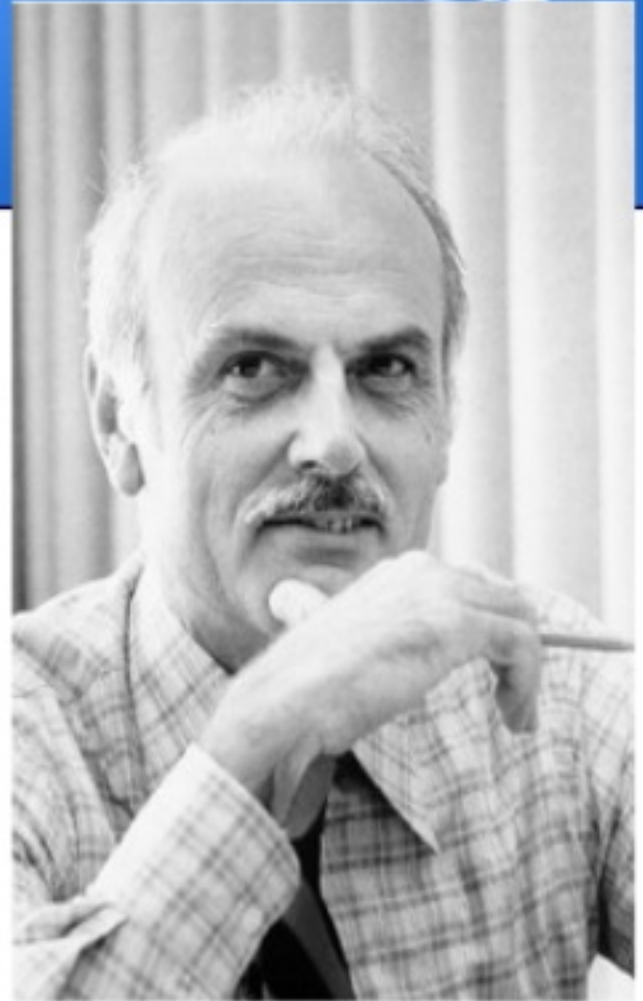
## Boyce-Codd Normal Form

# Boyce-Codd Normal Form

Dr. Raymond F. Boyce

Edgar Frank “Ted” Codd

"A Relational Model of Data for  
Large Shared Data Banks"



# Boyce-Codd Normal Form

There are no  
“bad” FDs:

**Definition.** A relation  $R$  is in BCNF if:

Whenever  $X \rightarrow B$  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)

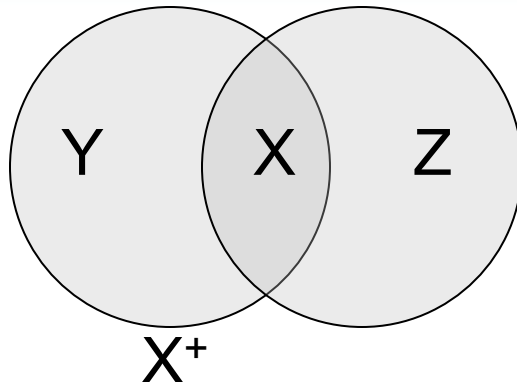
find  $X$  s.t.:  $X \neq X^+$  and  $X^+ \neq [\text{all attributes}]$

**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( $R_2$ );

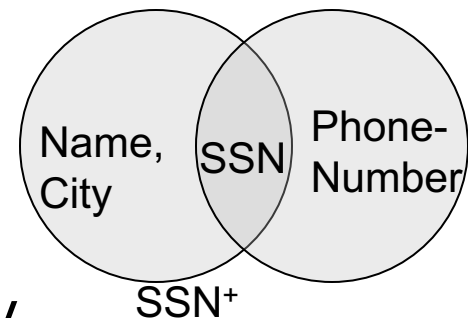




# Example

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

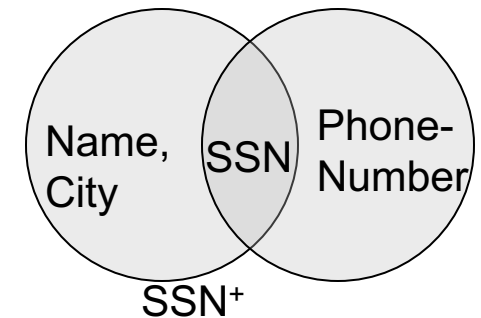
In other words:

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

# Example BCNF Decomposition

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

SSN  $\rightarrow$  Name, City



<b><u>SSN</u></b>	<b><u>PhoneNumber</u></b>
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)

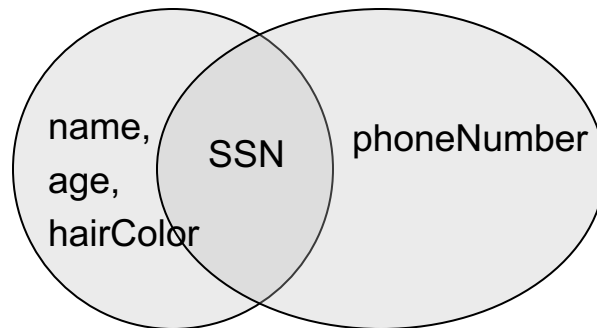
SSN  $\rightarrow$  name, age

age  $\rightarrow$  hairColor

Iteration 1: **Person**: SSN<sup>+</sup> = SSN, name, age, hairColor

Decompose into: **P**(SSN, name, 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

What are  
the keys ?

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)

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<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)

R(A,B,C,D)

## Example: BCNF

A	→	B
B	→	C

R(A,B,C,D)

R(A,B,C,D)

A → B  
B → 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

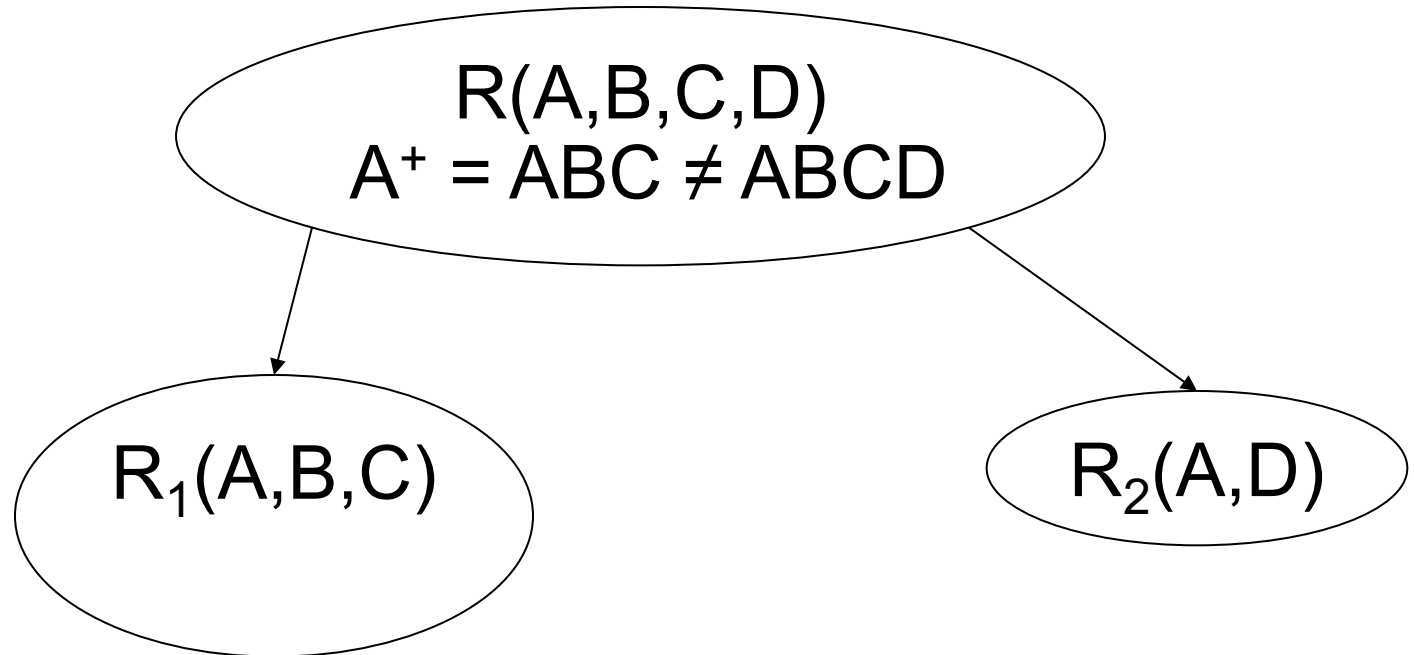
## Example: BCNF

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

R(A,B,C,D)

A → B  
B → C

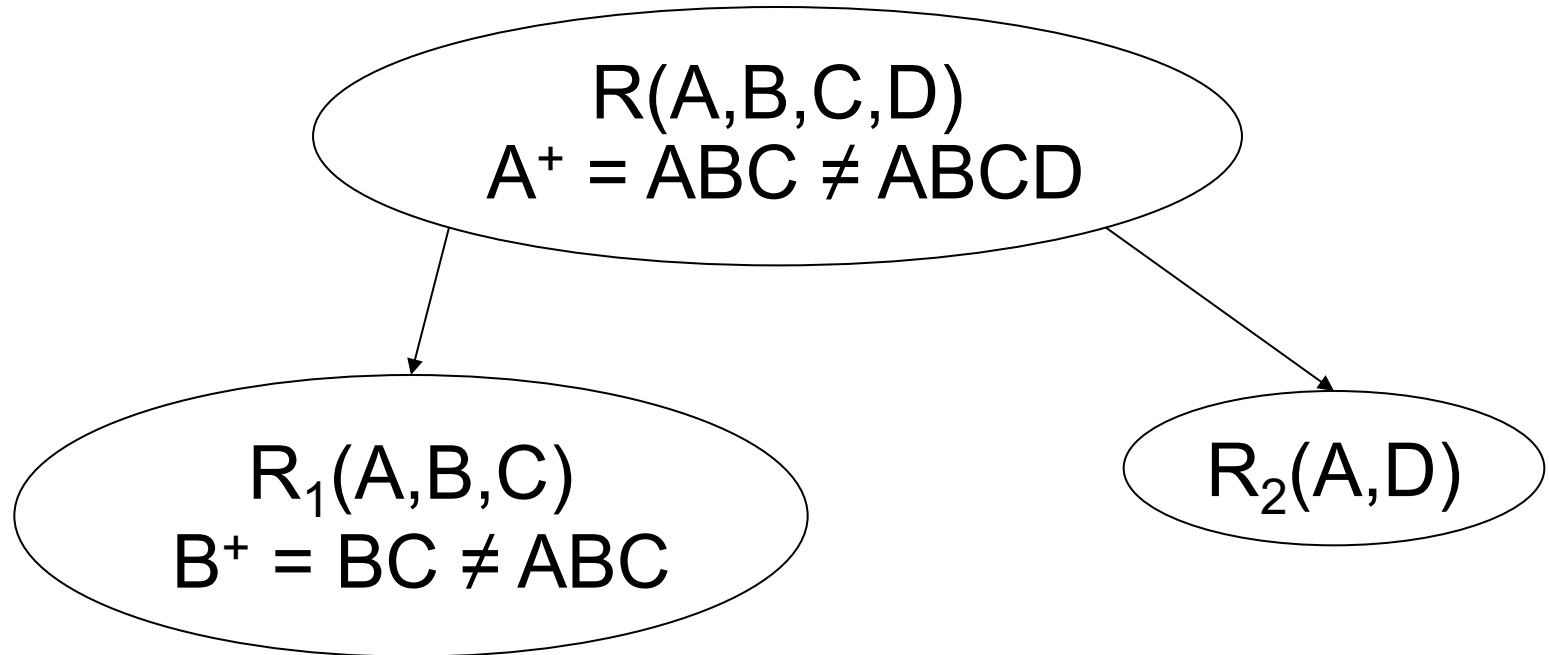
## Example: BCNF



R(A,B,C,D)

A → B  
B → C

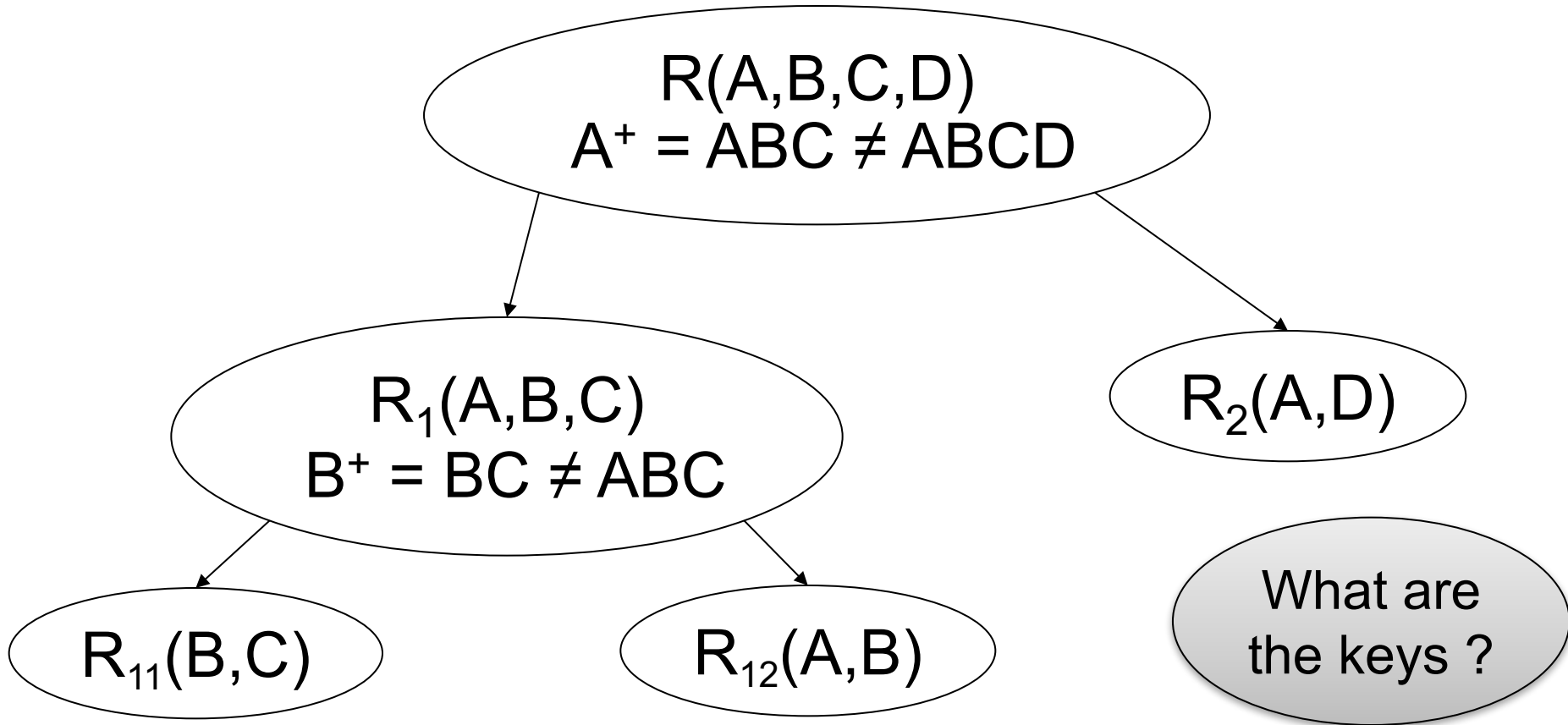
## Example: BCNF



R(A,B,C,D)

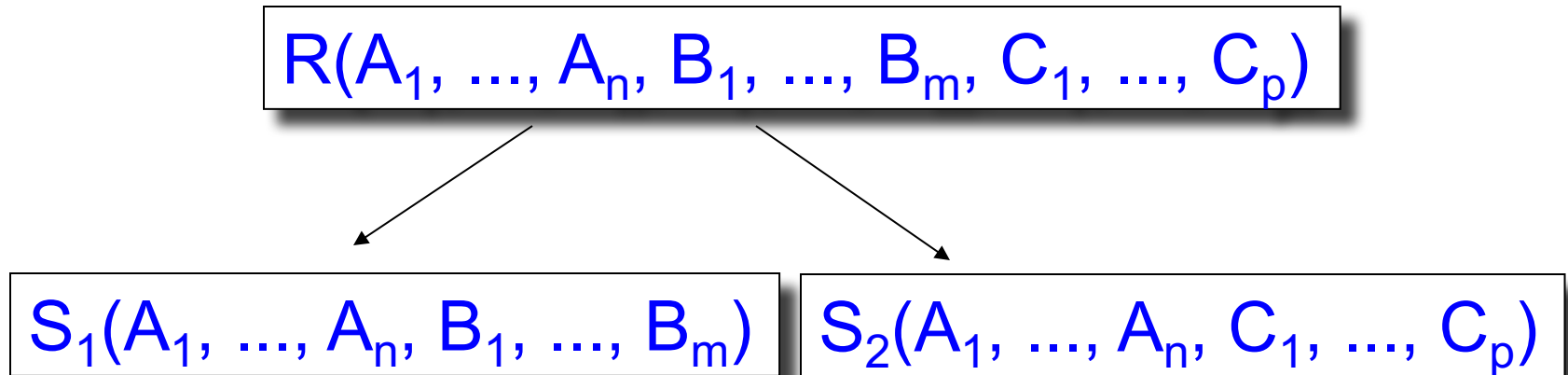
A → B  
B → C

## Example: BCNF



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

# 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
<del>Gizmo</del>	<del>19.99</del>

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

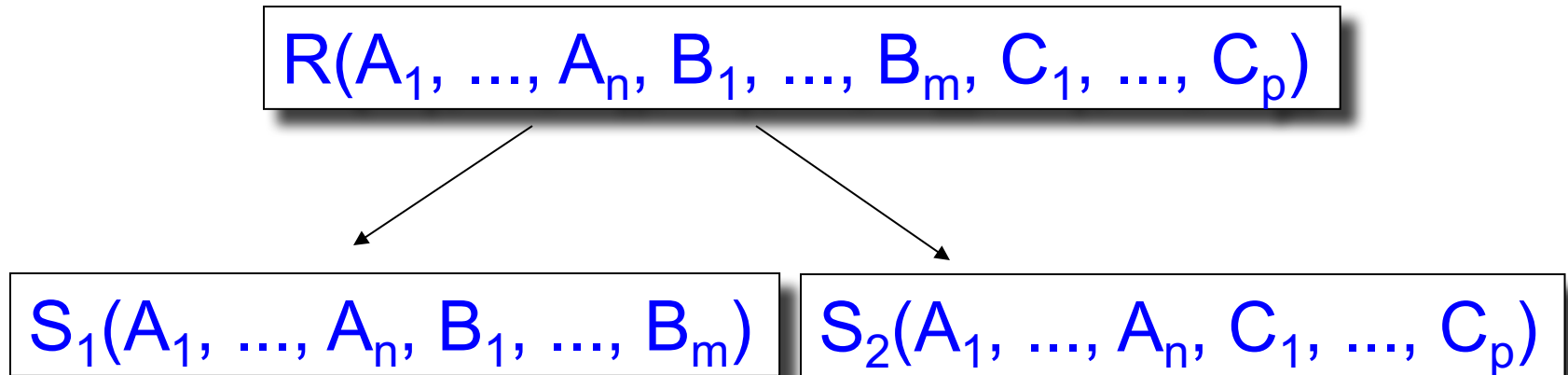
Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera

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 109

# Testing for Lossless Join

If we decompose  $R$  into  $\Pi_{S_1}(R)$ ,  $\Pi_{S_2}(R)$ ,  $\Pi_{S_3}(R)$ , ...  
Is it true that  $S_1 \bowtie S_2 \bowtie S_3 \bowtie \dots = R$  ?

That is true if we can show that:

$R \subseteq S_1 \bowtie S_2 \bowtie S_3 \bowtie \dots$  always holds (why?)

$R \supseteq S_1 \bowtie S_2 \bowtie S_3 \bowtie \dots$  **need to check**

Example from textbook Ch. 3.4.2

# The Chase Test for Lossless Join

$$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$$

R satisfies:  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $CD \rightarrow A$

$$S1 = \Pi_{AD}(R), S2 = \Pi_{AC}(R), S3 = \Pi_{BCD}(R),$$

$$\text{hence } R \subseteq S1 \bowtie S2 \bowtie S3$$

$$\text{Need to check: } R \supseteq S1 \bowtie S2 \bowtie S3$$

# The Chase Test for Lossless Join

$$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$$

R satisfies:  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $CD \rightarrow A$

$$S1 = \Pi_{AD}(R), S2 = \Pi_{AC}(R), S3 = \Pi_{BCD}(R),$$

$$\text{hence } R \subseteq S1 \bowtie S2 \bowtie S3$$

$$\text{Need to check: } R \supseteq S1 \bowtie S2 \bowtie S3$$

Suppose  $(a,b,c,d) \in S1 \bowtie S2 \bowtie S3$  Is it also in R?

# The Chase Test for Lossless Join

$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$   
R satisfies:  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $CD \rightarrow A$

$S1 = \Pi_{AD}(R)$ ,  $S2 = \Pi_{AC}(R)$ ,  $S3 = \Pi_{BCD}(R)$ ,

hence  $R \subseteq S1 \bowtie S2 \bowtie S3$

Need to check:  $R \supseteq S1 \bowtie S2 \bowtie S3$

Suppose  $(a,b,c,d) \in S1 \bowtie S2 \bowtie S3$  Is it also in R?

R must contain the following tuples:

A	B	C	D
a	b1	c1	d

Why ?

$(a,d) \in S1 = \Pi_{AD}(R)$

# The Chase Test for Lossless Join

$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$   
R satisfies:  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $CD \rightarrow A$

$S1 = \Pi_{AD}(R)$ ,  $S2 = \Pi_{AC}(R)$ ,  $S3 = \Pi_{BCD}(R)$ ,

hence  $R \subseteq S1 \bowtie S2 \bowtie S3$

Need to check:  $R \supseteq S1 \bowtie S2 \bowtie S3$

Suppose  $(a,b,c,d) \in S1 \bowtie S2 \bowtie S3$  Is it also in R?

R must contain the following tuples:

A	B	C	D
a	b1	c1	d
a	b2	c	d2

Why ?

$(a,d) \in S1 = \Pi_{AD}(R)$

$(a,c) \in S2 = \Pi_{AC}(R)$

# The Chase Test for Lossless Join

$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$   
R satisfies:  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $CD \rightarrow A$

$S1 = \Pi_{AD}(R)$ ,  $S2 = \Pi_{AC}(R)$ ,  $S3 = \Pi_{BCD}(R)$ ,

hence  $R \subseteq S1 \bowtie S2 \bowtie S3$

Need to check:  $R \supseteq S1 \bowtie S2 \bowtie S3$

Suppose  $(a,b,c,d) \in S1 \bowtie S2 \bowtie S3$  Is it also in R?

R must contain the following tuples:

A	B	C	D
a	b1	c1	d
a	b2	c	d2
a3	b	c	d

Why ?

$(a,d) \in S1 = \Pi_{AD}(R)$

$(a,c) \in S2 = \Pi_{AC}(R)$

$(b,c,d) \in S3 = \Pi_{BCD}(R)$

# The Chase Test for Lossless Join

$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$   
 R satisfies:  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $CD \rightarrow A$

$S1 = \Pi_{AD}(R)$ ,  $S2 = \Pi_{AC}(R)$ ,  $S3 = \Pi_{BCD}(R)$ ,

hence  $R \subseteq S1 \bowtie S2 \bowtie S3$

Need to check:  $R \supseteq S1 \bowtie S2 \bowtie S3$

Suppose  $(a,b,c,d) \in S1 \bowtie S2 \bowtie S3$  Is it also in R?

R must contain the following tuples:

A	B	C	D
a	b1	c1	d
a	b2	c	d2
a3	b	c	d

Why ?

$(a,d) \in S1 = \Pi_{AD}(R)$

$(a,c) \in S2 = \Pi_{AC}(R)$

$(b,c,d) \in S3 = \Pi_{BCD}(R)$

“Chase” them (apply FDs):

$A \rightarrow B$

A	B	C	D
a	b1	c1	d
a	b1	c	d2
a3	b	c	d





# The Chase Test for Lossless Join

$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$   
 R satisfies:  $A \rightarrow B, B \rightarrow C, CD \rightarrow A$

$S1 = \Pi_{AD}(R), S2 = \Pi_{AC}(R), S3 = \Pi_{BCD}(R)$ ,  
 hence  $R \subseteq S1 \bowtie S2 \bowtie S3$

Need to check:  $R \supseteq S1 \bowtie S2 \bowtie S3$

Suppose  $(a,b,c,d) \in S1 \bowtie S2 \bowtie S3$  Is it also in R?

R must contain the following tuples:

A	B	C	D
a	b1	c1	d
a	b2	c	d2
a3	b	c	d

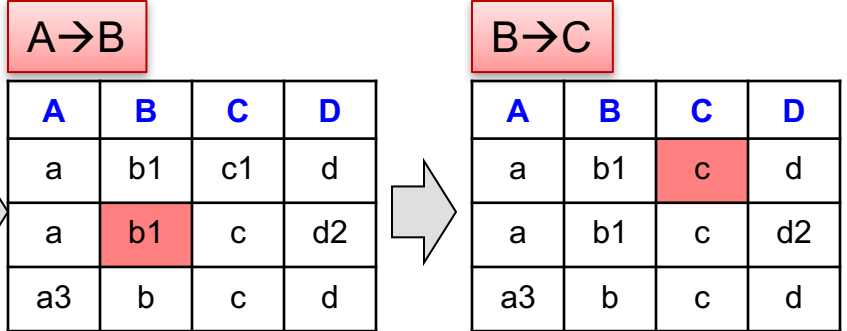
Why ?

$(a,d) \in S1 = \Pi_{AD}(R)$

$(a,c) \in S2 = \Pi_{AC}(R)$

$(b,c,d) \in S3 = \Pi_{BCD}(R)$

“Chase” them (apply FDs):



# The Chase Test for Lossless Join

$R(A,B,C,D) = S1(A,D) \bowtie S2(A,C) \bowtie S3(B,C,D)$   
 $R$  satisfies:  $A \rightarrow B, B \rightarrow C, CD \rightarrow A$

$S1 = \Pi_{AD}(R), S2 = \Pi_{AC}(R), S3 = \Pi_{BCD}(R),$

hence  $R \subseteq S1 \bowtie S2 \bowtie S3$

Need to check:  $R \supseteq S1 \bowtie S2 \bowtie S3$

Suppose  $(a,b,c,d) \in S1 \bowtie S2 \bowtie S3$  Is it also in  $R$ ?

$R$  must contain the following tuples:

A	B	C	D
a	b1	c1	d
a	b2	c	d2
a3	b	c	d

Why ?

$(a,d) \in S1 = \Pi_{AD}(R)$

$(a,c) \in S2 = \Pi_{AC}(R)$

$(b,c,d) \in S3 = \Pi_{BCD}(R)$

“Chase” them (apply FDs):

$A \rightarrow B$

A	B	C	D
a	b1	c1	d
a	b1	c	d2
a3	b	c	d

$B \rightarrow C$

A	B	C	D
a	b1	c	d
a	b1	c	d2
a3	b	c	d

$CD \rightarrow A$

A	B	C	D
a	b1	c	d
a	b1	c	d2
a	b	c	d

Hence  $R$   
contains  $(a,b,c,d)$

# Schema Refinements = Normal Forms

- 1st Normal Form = all tables are flat
- 2nd Normal Form = obsolete
- Boyce Codd Normal Form = no bad FDs
- 3rd Normal Form = see book
  - BCNF is lossless but can cause loss of ability to check some FDs (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

# Getting Practical

How to implement normalization in SQL

# Motivation

- We learned about how to normalize tables to avoid anomalies
- How can we implement normalization in SQL if we can't modify existing tables?
  - This might be due to legacy applications that rely on previous schemas to run

# Views

- A **view** in SQL =
  - A table computed from other tables, s.t., whenever the base tables are updated, the view is updated too
- More generally:
  - A **view** is derived data that keeps track of changes in the original data
- Compare:
  - A **function** computes a value from other values, but does not keep track of changes to the inputs

Purchase(customer, product, store)

Product(pname, price)

StorePrice(store, price)

# A Simple View

Create a view that returns for each store  
the prices of products purchased at that store

```
CREATE VIEW StorePrice AS
SELECT DISTINCT x.store, y.price
FROM Purchase x, Product y
WHERE x.product = y.pname
```

This is like a new table  
StorePrice(store, price)

Purchase(customer, product, store)

Product(pname, price)

StorePrice(store, price)

# We Use a View Like Any Table

- A "high end" store is a store that sell some products over 1000.
- For each customer, return all the high end stores that they visit.

```
SELECT DISTINCT u.customer, u.store
FROM Purchase u, StorePrice v
WHERE u.store = v.store
      AND v.price > 1000
```



# Types of Views

- Virtual views
  - Computed only on-demand – slow at runtime
  - Always up to date
- Materialized views
  - Pre-computed offline – fast at runtime
  - May have stale data (must recompute or update)
  - Indexes *are* materialized views
- A key component of physical tuning of databases is the selection of materialized views and indexes

# Vertical Partitioning

Resumes

<u>SSN</u>	Name	Address	Resume	Picture
234234	Mary	Huston	Clob1...	Blob1...
345345	Sue	Seattle	Clob2...	Blob2...
345343	Joan	Seattle	Clob3...	Blob3...
432432	Ann	Portland	Clob4...	Blob4...

T1

<u>SSN</u>	Name	Address
234234	Mary	Huston
345345	Sue	Seattle
...		

T2

<u>SSN</u>	Resume
234234	Clob1...
345345	Clob2...

T3

<u>SSN</u>	Picture
234234	Blob1...
345345	Blob2...

T2.SSN is a key *and* a foreign key to T1.SSN. Same for T3.SSN

T1(ssn,name,address)

Resumes(ssn,name,address,resume,picture)

T2(ssn,resume)

T3(ssn,picture)

# Vertical Partitioning

```
CREATE VIEW Resumes AS
  SELECT T1.ssn, T1.name, T1.address,
         T2.resume, T3.picture
  FROM   T1,T2,T3
  WHERE  T1.ssn=T2.ssn AND T1.ssn=T3.ssn
```

T1(ssn,name,address)

Resumes(ssn,name,address,resume,picture)

T2(ssn,resume)

T3(ssn,picture)

# Vertical Partitioning

```
CREATE VIEW Resumes AS
  SELECT T1.ssn, T1.name, T1.address,
         T2.resume, T3.picture
  FROM   T1,T2,T3
  WHERE  T1.ssn=T2.ssn AND T1.ssn=T3.ssn
```

```
SELECT address
FROM   Resumes
WHERE  name = 'Sue'
```

T1(ssn,name,address)

Resumes(ssn,name,address,resume,picture)

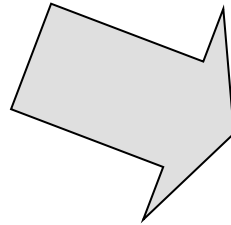
T2(ssn,resume)

T3(ssn,picture)

# Vertical Partitioning

```
CREATE VIEW Resumes AS
  SELECT T1.ssn, T1.name, T1.address,
         T2.resume, T3.picture
  FROM   T1,T2,T3
  WHERE  T1.ssn=T2.ssn AND T1.ssn=T3.ssn
```

```
SELECT address
FROM   Resumes
WHERE  name = 'Sue'
```



Original query:

```
SELECT T1.address
FROM T1, T2, T3
WHERE T1.name = 'Sue'
      AND T1.SSN=T2.SSN
      AND T1.SSN = T3.SSN
```

T1(ssn,name,address)

T2(ssn,resume)

T3(ssn,picture)

Resumes(ssn,name,address,resume,picture)

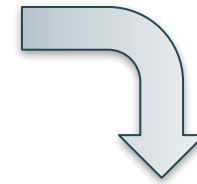
# Vertical Partitioning

```
CREATE VIEW Resumes AS
  SELECT T1.ssn, T1.name, T1.address,
         T2.resume, T3.picture
  FROM   T1,T2,T3
  WHERE  T1.ssn=T2.ssn AND T1.ssn=T3.ssn
```

```
SELECT address
FROM   Resumes
WHERE  name = 'Sue'
```

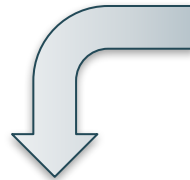
Final query:

```
SELECT T1.address
FROM   T1
WHERE  T1.name = 'Sue'
```



Modified query:

```
SELECT T1.address
FROM   T1, T2, T3
WHERE  T1.name = 'Sue'
AND T1.SSN=T2.SSN
AND T1.SSN = T3.SSN
```



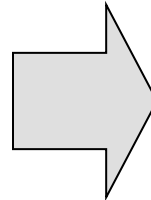
# Vertical Partitioning Applications

- **Advantages**
  - Speeds up queries that touch only a small fraction of columns
  - Single column can be compressed effectively, reducing disk I/O
- **Disadvantages**
  - Updates are expensive!
  - Need many joins to access many columns
  - Repeated key columns add overhead

# Horizontal Partitioning

## Customers

SSN	Name	City
234234	Mary	Houston
345345	Sue	Seattle
345343	Joan	Seattle
234234	Ann	Portland
--	Frank	Calgary
--	Jean	Montreal



## CustomersInHouston

SSN	Name	City
234234	Mary	Houston

## CustomersInSeattle

SSN	Name	City
345345	Sue	Seattle
345343	Joan	Seattle

.....



CustomersInHouston(ssn,name,city)

Customers(ssn,name,city)

CustomersInSeattle(ssn,name,city)

.....

# Horizontal Partitioning

```
CREATE VIEW Customers AS
  CustomersInHouston
  UNION ALL
  CustomersInSeattle
  UNION ALL
  ...
```

CustomersInHouston(ssn,name,city)

CustomersInSeattle(ssn,name,city)

Customers(ssn,name,city)

.....

# Horizontal Partitioning

```
SELECT name  
FROM Customers  
WHERE city = 'Seattle'
```

Which tables are inspected by the system ?

CustomersInHouston(ssn,name,city)

Customers(ssn,name,city)

CustomersInSeattle(ssn,name,city)

.....

# Horizontal Partitioning

Better: remove CustomerInHouston.city etc

```
CREATE VIEW Customers AS
  (SELECT SSN, name, 'Houston' as city
   FROM CustomersInHouston)
  UNION ALL
  (SELECT SSN, name, 'Seattle' as city
   FROM CustomersInSeattle)
  UNION ALL
  . . .
```

CustomersInHouston(ssn,name,city)

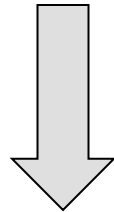
Customers(ssn,name,city)

CustomersInSeattle(ssn,name,city)

.....

# Horizontal Partitioning

```
SELECT name  
FROM Customers  
WHERE city = 'Seattle'
```



```
SELECT name  
FROM CustomersInSeattle
```

# Horizontal Partitioning Applications

- Performance optimization
  - Especially for data warehousing
  - E.g., one partition per month
  - E.g., archived applications and active applications
- Distributed and parallel databases
- Data integration

# Conclusion

- Poor schemas can lead to performance inefficiencies
- E/R diagrams are means to structurally visualize and design relational schemas
- Normalization is a principled way of converting schemas into a form that avoid such problems
- BCNF is one of the most widely used normalized form in practice