

# **Chapter 6: Database Normalization**

# Outline

- Key of a relational schema
- Normalization
- Normal Forms

# Keys

## □ Key/candidate key:

□ We say a set of **one or more attributes**  $\{A_1, A_2, \dots, A_n\}$  is a key/candidate key for a relation  $R$  if:

- ▶ Those attributes functionally determine all other attributes of the relation.
- ▶ **No proper subset** of  $\{A_1, A_2, \dots, A_n\}$  functionally determines all other attributes of  $R$ ; i.e., **a key must be minimal.**

□ **Superkey:** A set of attributes that contains a **key** is called a superkey.

# The algorithm for computing a candidate key of R.

- **Input:** a relation  $R = \{A_1, A_2, \dots, A_n\}$ , and  $F$ , a set of functional dependencies.
- **Output:**  $K$ , a candidate key of  $R$ .
- **Algorithm**

```
set  $K \leftarrow R$ ;  
for each attribute  $A \in K$   
{  
    compute  $\{K - A\}^+$  with respect to  $F$   
    if  $\{K - A\}^+$  contains all attributes of  $R$  then  
        set  $K = \{K - A\}$   
}  
return  $K$ 
```

# The algorithm for computing a candidate key of R.

## □ Example:

- $R = \{A, B, C, D, E, F, G\}$
- $F = \{B \rightarrow A, D \rightarrow C, D \rightarrow BE, DF \rightarrow G\}$
- $K = ?$

- B1:

$K = ABCDEFG.$

- B2:

- Lặp 1:  $(BCDEFG)_F^+ = BCDEFGA \Rightarrow K = BCDEFG.$
- Lặp 2:  $(CDEFG)_F^+ = CDEFGBA \Rightarrow K = CDEFG.$
- Lặp 3:  $(DEFG)_F^+ = DEFGCBA \Rightarrow K = DEFG.$
- Lặp 4:  $(EFG)_F^+ = EFG.$
- Lặp 5:  $(DFG)_F^+ = DFGCBEA \Rightarrow K = DFG.$
- Lặp 6:  $(DG)_F^+ = DGCBEA.$
- Lặp 7:  $(DF)_F^+ = DFCBEAG \Rightarrow K = DF.$

- B3:

Khóa là  $K = DF.$

- B1:  
 $K = ABCDEFG$ .
- B2:
  - Lặp 1:  $(BCDEFG)_F^+ = BCDEFGA \Rightarrow K = BCDEFG$ .
  - Lặp 2:  $(CDEFG)_F^+ = CDEFGBA \Rightarrow K = CDEFG$ .
  - Lặp 3:  $(DEFG)_F^+ = DEFGCBA \Rightarrow K = DEFG$ .
  - Lặp 4:  $(EFG)_F^+ = EFG$ .
  - Lặp 5:  $(DFG)_F^+ = DFGCBEA \Rightarrow K = DFG$ .
  - Lặp 6:  $(DG)_F^+ = DGCBEA$ .
  - Lặp 7:  $(DF)_F^+ = DFCBEAG \Rightarrow K = DF$ .
- B3:  
 Khóa là  $K = DF$ .

□  $R(A, B, C, D, E),$

□  $F = \{AB \rightarrow C, B \rightarrow D, CD \rightarrow E\}$

# Normalization

## □ What is normalization?

- ▶ Basically, it's the process of efficiently organizing data in a database.

## □ Normalization theory is based on the observation that relations with certain properties are *more effective in inserting, updating and deleting data than other sets of relations containing the same data*

## □ Normalization is a multi-step process beginning with an “unnormalized” relation

- This is the process which allows you to winnow out redundant data within your database.
- This involves restructuring the tables to successively meeting higher forms of Normalization.



# Normalization

- Normalization theory is based on the concepts of **normal forms**.
- A relational table is said to be in a **particular normal form** if it satisfied a **certain set of constraints**.
- A properly normalized database should have the following characteristics
  - ▶ Scalar values in each fields
  - ▶ Absence of redundancy.
  - ▶ Minimal use of null values.
  - ▶ Minimal loss of information.

# Normalization

- *The goal of normalization is to create a set of relational tables that are **free of redundant** data and that can be **consistently** and **correctly modified or updated**.*

# Definitions

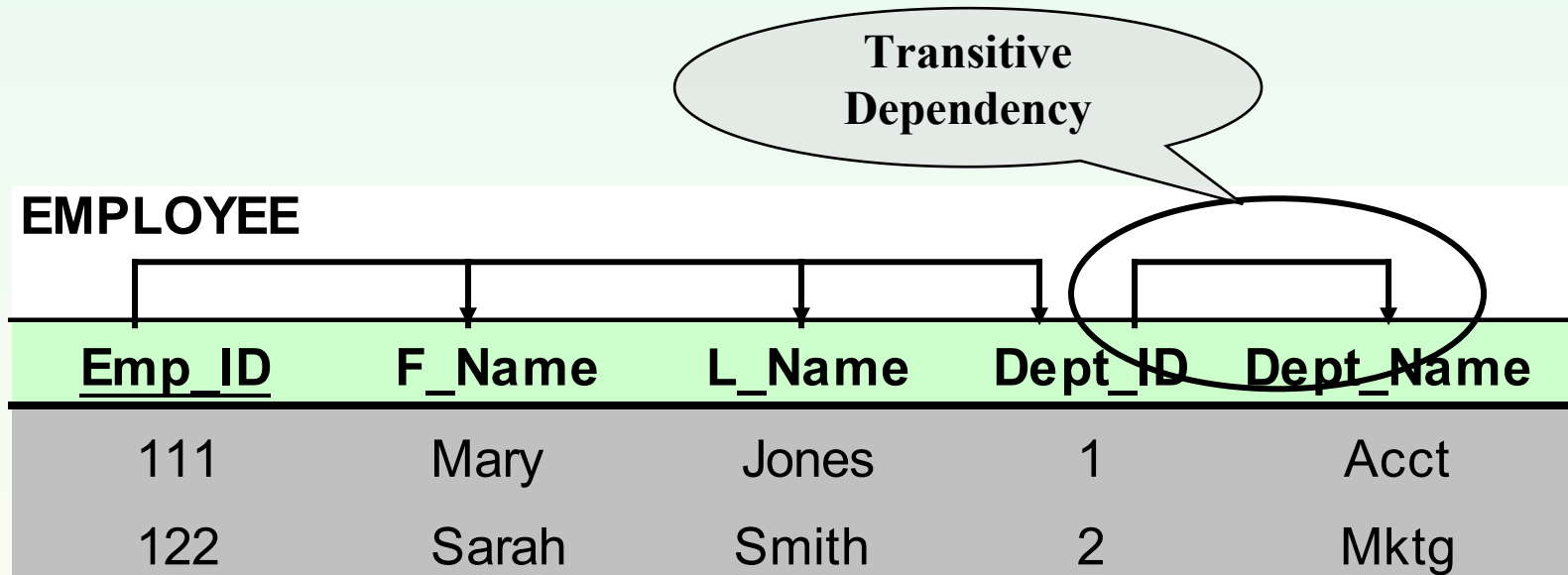
- **Partial FD** – *when a non-key attribute is determined by a part, but not the whole, of a **COMPOSITE** primary key.*
- **Example:** Key= {Cust\_ID, Order\_ID}
  - ▶ Cust\_ID → Name is **partial FD**

**CUSTOMER**

<u>Cust_ID</u>	Name	<u>Order_ID</u>
101	AT&T	1234
101	AT&T	156
125	Cisco	1250

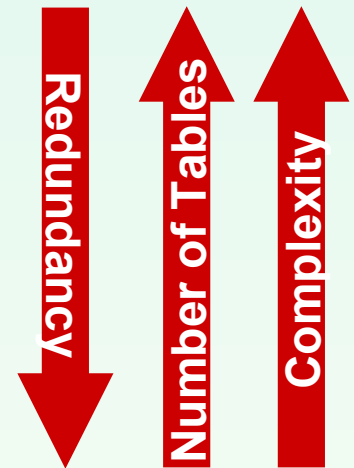
# Definitions

- **Transitive FD** – *when a non-key attribute determines another non-key attribute.*
- Example: Key= {Emp\_ID}
  - ▶ Dept\_ID → Dept\_Name is **transitive FD**



# Normal Forms

- Levels of normalization based on the amount of redundancy in the database.
- Various levels of normalization are:
  - First Normal Form (1NF)
  - Second Normal Form (2NF)
  - Third Normal Form (3NF)
  - Boyce-Codd Normal Form (BCNF)



**Most databases should be 3NF or BCNF in order to avoid the database anomalies.**

**Each higher level is a subset of the lower level**

**FDs and keys are used to define normal forms for relations**

# Normal Forms

- Unnormalized – There are multivalued attributes or repeating groups
- 1 NF – No multivalued attributes or repeating groups.
- 2 NF – 1 NF plus no partial dependencies
- 3 NF – 2 NF plus no transitive dependencies
- BCNF – 3 NF plus all determinants are superkeys

# First Normal Form (1NF)

- A **table** is considered to be in 1NF if **all the fields contain only atomic or scalar values** (as opposed to list of values or there are no sets of values within a column).
  - ▶ No repeating groups
  - ▶ A column or set of columns is called a Candidate Key when its values can uniquely identify the row in the relation.
- Note: **All *relations* are in 1<sup>st</sup> Normal Form**

# First Normal Form (1NF)

□ **Example:** Table with multivalued attributes, not in 1<sup>st</sup> normal form

<u>Order_ID</u>	Order_ Date	Customer_ ID	Customer_ Name	Customer_ Address	<u>Product_ID</u>	Product_ Description	Product_ Finish	Unit_ Price	Ordered_ Quantity
1006	10/24/2004	2	Value Furniture	Plano, TX	7	Dining Table	Natural Ash	800.00	2
					5	Writer's Desk	Cherry	325.00	2
					4	Entertainment Center	Natural Maple	650.00	1
1007	10/25/2004	6	Furniture Gallery	Boulder, CO	11	4-Dr Dresser	Oak	500.00	4
					4	Entertainment Center	Natural Maple	650.00	3



# First Normal Form (1NF)

- **Example:** Table with no multivalued attributes and unique rows, in 1st normal form

INVOICE relation (1NF) (Pine Valley Furniture Company)

<u>Order_ID</u>	<u>Order_</u> Date	<u>Customer_</u> ID	<u>Customer_</u> Name	<u>Customer_</u> Address	<u>Product_ID</u>	<u>Product_</u> Description	<u>Product_</u> Finish	<u>Unit_</u> Price	<u>Ordered_</u> Quantity
1006	10/24/2004	2	Value Furniture	Plano, TX	7	Dining Table	Natural Ash	800.00	2
1006	10/24/2004	2	Value Furniture	Plano, TX	5	Writer's Desk	Cherry	325.00	2
1006	10/24/2004	2	Value Furniture	Plano, TX	4	Entertainment Center	Natural Maple	650.00	1
1007	10/25/2004	6	Furniture Gallery	Boulder, CO	11	4-Dr Dresser	Oak	500.00	4
1007	10/25/2004	6	Furniture Gallery	Boulder, CO	4	Entertainment Center	Natural Maple	650.00	3

Note: this is relation, but not a well-structured one

# Second Normal Form (2NF)

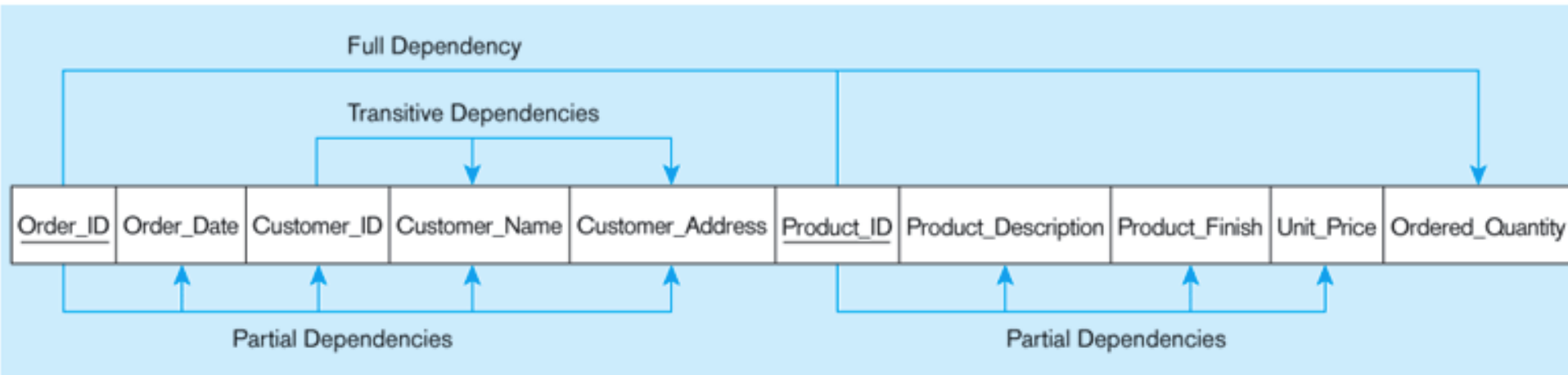
- Where the 1NF deals with atomicity of data, the Second Normal Form (or 2NF) deals with relationships between composite key columns and non-key columns:
  - ▶ Meet all the requirements of the first normal form.
  - ▶ Any non-key columns must depend on the entire primary key. In the case of a composite primary key, this means that a non-key column cannot depend on only part of the composite key.
  - ▶ Create relationships between these new tables and their predecessors through the use of foreign keys.
  - ▶ A relation R is in 2nf if every non-primary attribute A in R is fully Functionally dependent on the primary key.

## Second Normal Form (2NF)

- 1NF PLUS *every non-key attribute is fully functionally dependent on the ENTIRE primary key*
  - ▶ Every non-key attribute must be defined by the entire key, not by only part of the key
  - ▶ No partial functional dependencies

# Second Normal Form (2NF)

Functional dependency diagram for INVOICE



**Order\_ID → Order\_Date, Customer\_ID, Customer\_Name, Customer\_Address**

**Customer\_ID → Customer\_Name, Customer\_Address**

**Product\_ID → Product\_Description, Product\_Finish, Unit\_Price**

**Order\_ID, Product\_ID → Order\_Quantity**

**Therefore, NOT in 2<sup>nd</sup> Normal Form**

# Second Normal Form (2NF)

Removing partial dependencies

ORDER\_LINE (3NF)

<u>Order_ID</u>	<u>Product_ID</u>	Ordered_Quantity
-----------------	-------------------	------------------

PRODUCT (3NF)

<u>Product_ID</u>	Product_Description	Product_Finish	Unit_Price
-------------------	---------------------	----------------	------------

CUSTOMER\_ORDER (2NF)

<u>Order_ID</u>	Order_Date	Customer_ID	Customer_Name	Customer_Address
-----------------	------------	-------------	---------------	------------------

Transitive Dependencies



Partial Dependencies are removed, but there are still transitive dependencies

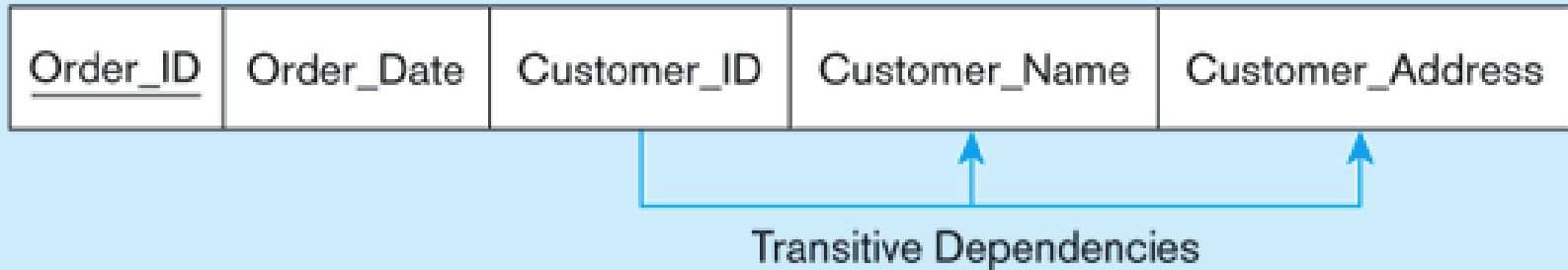
# Third Normal Form (3NF)

# Third Normal Form (3NF)

- **2NF PLUS *no transitive dependencies*** (functional dependencies on non-primary-key attributes)
- Note: this is called transitive, because the primary key is a determinant for another attribute, which in turn is a determinant for a third
- **Solution:** non-key determinant with transitive dependencies go into a new table; non-key determinant becomes primary key in the new table and stays as foreign key in the old table

# Third Normal Form (3NF)

## Customer\_Order (2NF)



Removing transitive dependencies

### ORDER (3NF)

<u>Order_ID</u>	Order_Date	<u>Customer_ID</u>
-----------------	------------	--------------------

### CUSTOMER (3NF)

<u>Customer_ID</u>	Customer_Name	Customer_Address
--------------------	---------------	------------------



# Boyce-Codd Normal Form (BCNF)

- BCNF eliminates all redundancy that can be discovered by functional dependencies
- Most 3NF relations are also BCNF relations.
- A 3NF relation is NOT in BCNF if:
  - ▶ Candidate keys in the relation are composite keys (they are not single attributes)
  - ▶ There is more than one candidate key in the relation, and
  - ▶ The keys are not disjoint, that is, some attributes in the keys are common

# Boyce-Codd Normal Form (BCNF)

## □ Definition

- A relation schema  $R$  is in BCNF with respect to a set  $F$  if:
  - ▶ For all functional dependencies of  $F$  of the form  $\alpha \rightarrow \beta$ , where  $\alpha \subseteq R$  and  $\beta \subseteq R$ 
    - $\alpha \rightarrow \beta$  is a trivial functional dependency ( $\beta \subseteq \alpha$ )
    - $\alpha$  is a superkey for schema  $R$
- *A database design is in BCNF if each member of the set of relational schemas that constitute the design is in BCNF*

# Boyce-Codd Normal Form (BCNF)

## □ Rule for schema not in BCNF

□ Let ***R*** be a schema **not in BCNF**, then there is at least one nontrivial functional dependency  $\alpha \rightarrow \beta$  such that  $\alpha$  is not a superkey

## □ Example of not BCNF:

- ▶ *bor\_loan = (customer\_id, loan\_number, amount)*
- ▶ *FD: loan\_number  $\rightarrow$  amount*
- ▶ *But **loan\_number** is not a superkey*

# BCNF Decomposition

- The definition of BCNF can be used to directly test if a relationship is in BCNF
- If a relation is not in BCNF it can be decomposed to create relations that are in BCNF
- Example:
  - *borrower = (customer id, loan number)*
    - Is BCNF because no nontrivial functional dependency hold onto it
  - *loan = (loan number, amount)*
    - Has **one nontrivial functional dependency** that holds, *loan\_number → amount*,  
But **loan\_number** is a **superkey** so loan is in BCNF

# 3NF vs BCNF

- BCNF requires that all nontrivial dependencies be of the form  $\alpha \rightarrow \beta$ , where  $\alpha$  is a **superkey**
- **3NF** relaxes this constraint a little bit by **allowing nontrivial functional dependencies**

# BCNF Decomposition

- If ***R*** is not in BCNF, we can decompose *R* into a collection of BCNF schemas:

$$R_1, R_2, \dots, R_n$$

- **Example:**

***lending*** = (*branch\_name*, *branch\_city*, *assets*, *customer\_name*, *loan\_number*, *amount*)

**FDs:**

*branch\_name* → *assets branch\_city*

*loan\_number* → *amount branch\_name*

**Candidate key** : {*loan\_number*, *customer\_name*}

***branch\_name*** is not superkey so ***lending*** is not BCNF

# BCNF Decomposition

□ So we replace *lending* by:

***branch*** = (*branch\_name*, *branch\_city*, *assets*)

***loan\_info*** = (*branch\_name*, *customer\_name*,  
*loan\_number*, *amount*)

◆ The only nontrivial functional dependencies that hold on ***branch*** include ***branch\_name*** on the left side of the arrow.

=> Since ***branch\_name*** is a key for ***branch***, the relation ***branch*** is in BCNF

# BCNF Decomposition

## □ For *loan\_info*

- The functional dependency

*loan\_number* → *amount branch\_name*

holds on *loan\_info*

- But *loan\_number* is not a key for *loan\_info*

→ So we replace *loan\_info* by

*loanb* = (*loan\_number*, *branch\_name*, *amount*)

*borrower* = (*customer\_name*, *loan\_number*)

→ *loanb* and *borrower* are in BCNF



## Advantage / Disadvantage of 3NF

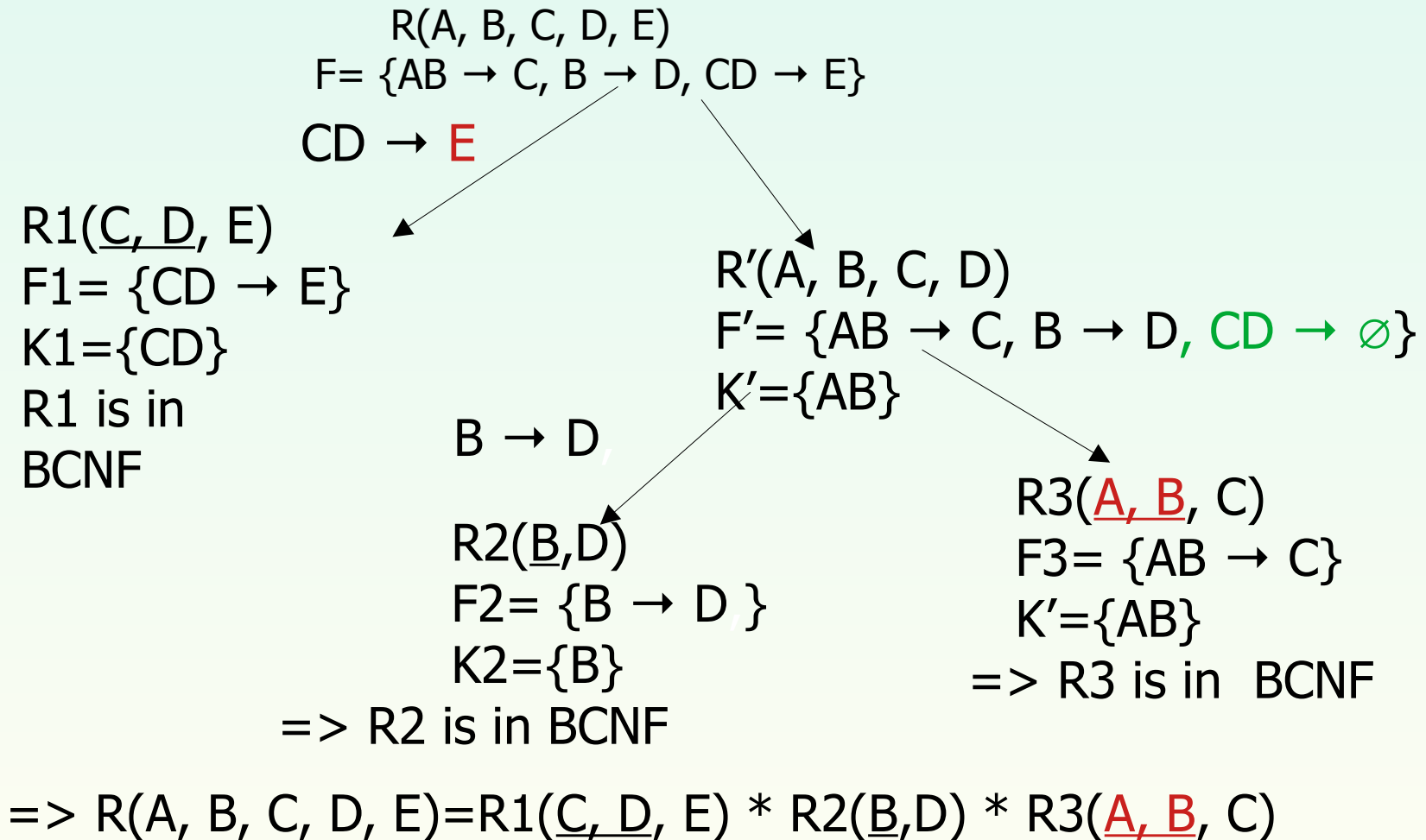
- **Advantage of 3NF:** it is always possible to obtain a 3NF design without sacrificing losslessness or dependency preservation
- **Disadvantage of 3NF:** we may have to use null values to represent some of the possible meaningful relationships among data items, and there is the problem of repetition of information

# Conclusion

- Goals of database design with functional dependencies are
  - 1) BCNF
  - 2) Losslessness
  - 3) Dependency preservation
- Not possible to get all 3, we have to choose between BCNF or dependency preservation

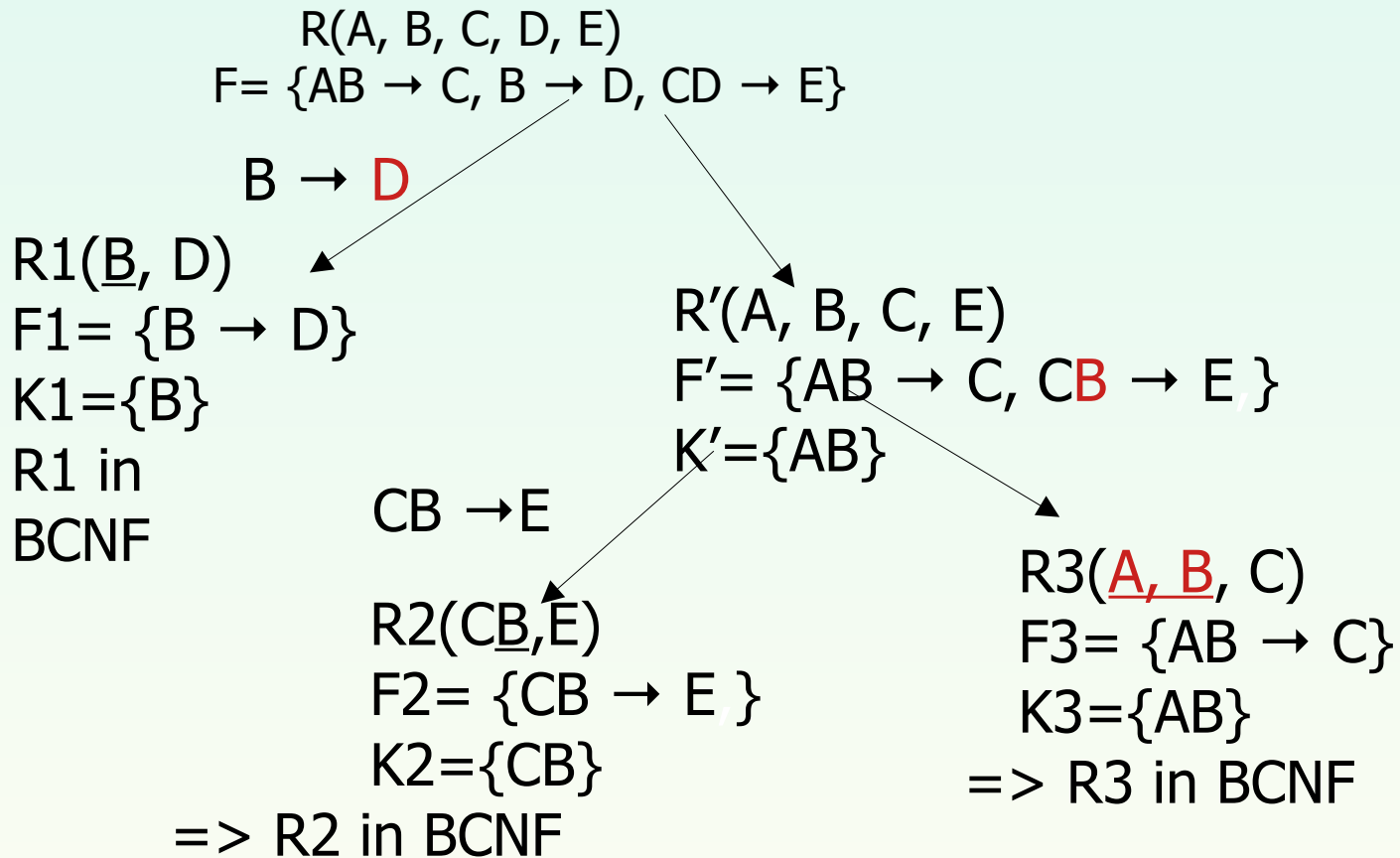
# Example

- $R(A, B, C, D, E), F = \{AB \rightarrow C, B \rightarrow D, CD \rightarrow E\}$
- $R$  is in BCNF ?
- $K = \{AB\} \Rightarrow R$  is in 1NF,  $B \rightarrow D$  is partial FD



# Bài tập

- $R(A, B, C, D, E), F = \{AB \rightarrow C, B \rightarrow D, CD \rightarrow E\}$
- R is in BCNF ?
- $K = \{AB\} \Rightarrow R$  is in 1NF,  $B \rightarrow D$  is partial FD



$$\Rightarrow R(A, B, C, D, E) = R1(\underline{B}, D) * R2(\underline{CB}, E) * R3(\underline{A}, \underline{B}, C)$$