

# DECOMPOSITION & SCHEMA NORMALIZATION

---

*CS 564 - Fall 2021*

---

*ACKs: Dan Suciu, Jignesh Patel, AnHai Doan*

---

# WHAT IS THIS LECTURE ABOUT?

---

- Bad schemas lead to redundancy
- To “correct” bad schemas: **decompose** relations
  - lossless-join
  - dependency preserving
- BCNF : a desired **normal forms**

---

# DB DESIGN THEORY

---

- Helps us identify the “bad” schemas and improve them
  1. express constraints on the data: **functional dependencies (FDs)**
  2. use the FDs to decompose the relations
- The process, called **normalization**, obtains a schema in a “normal form” that guarantees certain properties
  - examples of normal forms: **BCNF, 3NF, ...**

---

# SCHEMA DECOMPOSITION

---

# WHAT IS A DECOMPOSITION?

We **decompose** a relation  $\mathbf{R}(A_1, \dots, A_n)$  by creating

- $\mathbf{R}_1(B_1, \dots, B_m)$
- $\mathbf{R}_2(C_1, \dots, C_k)$

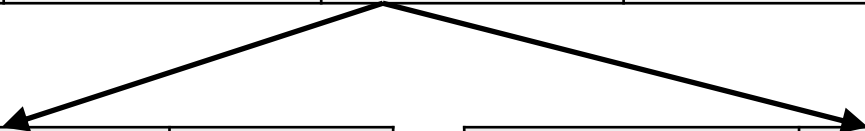
where  $\{B_1, \dots, B_m\} \cup \{C_1, \dots, C_k\} = \{A_1, \dots, A_n\}$

- The instance of  $\mathbf{R}_1$  is the projection of  $\mathbf{R}$  onto  $B_1, \dots, B_m$
- The instance of  $\mathbf{R}_2$  is the projection of  $\mathbf{R}$  onto  $C_1, \dots, C_k$

In general we can decompose a relation into multiple relations.

# EXAMPLE: DECOMPOSITION

SSN	name	age	phoneNumber
934729837	Paris	24	608-374-8422
934729837	Paris	24	603-534-8399
123123645	John	30	608-321-1163
384475687	Arun	20	206-473-8221



SSN	name	age
934729837	Paris	24
123123645	John	30
384475687	Arun	20

SSN	phoneNumber
934729837	608-374-8422
934729837	603-534-8399
123123645	608-321-1163
384475687	206-473-8221

---

# DECOMPOSITION DESIDERATA

---

What should a **good** decomposition achieve?

1. minimize redundancy
2. avoid information loss (**lossless-join**)
3. preserve the FDs (**dependency preserving**)
4. ensure good query performance

# EXAMPLE: INFORMATION LOSS

name	age	phoneNumber
Paris	24	608-374-8422
John	24	608-321-1163
Arun	20	206-473-8221

Decompose into:

$R_1(\text{name, age})$

$R_2(\text{age, phoneNumber})$

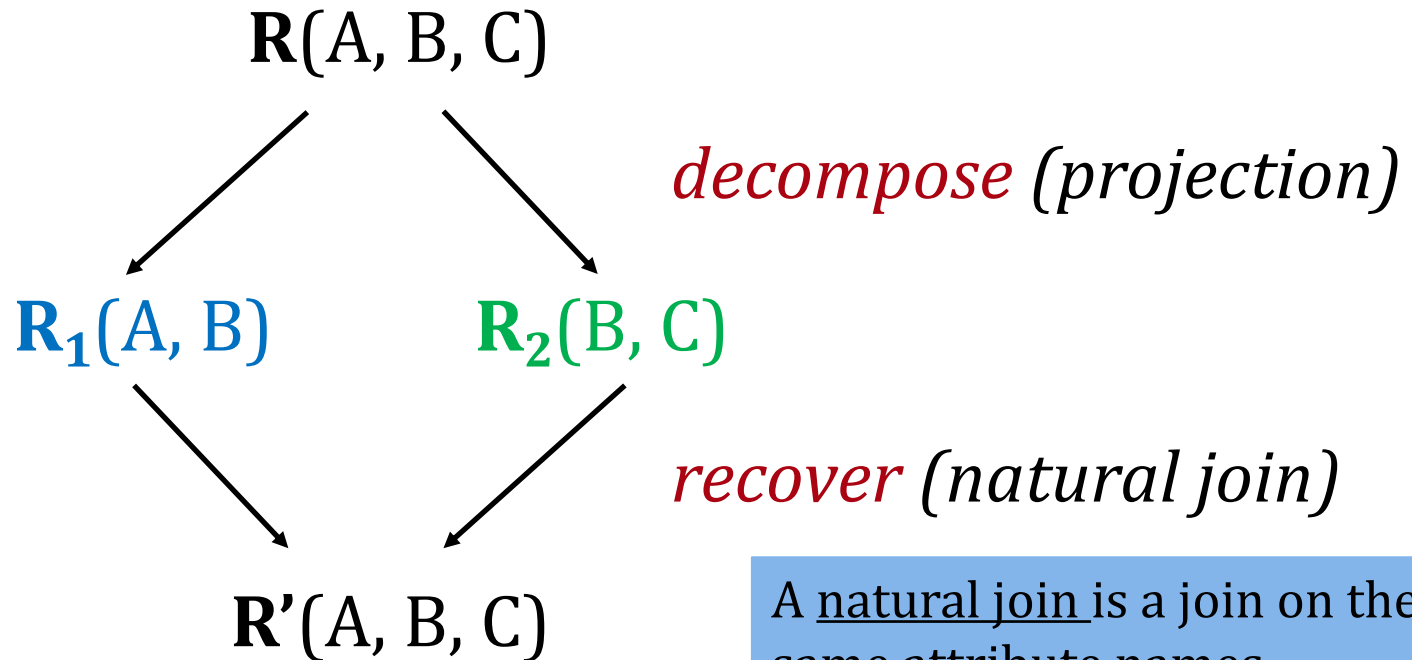
name	age
Paris	24
John	24
Arun	20

age	phoneNumber
24	608-374-8422
24	608-321-1163
20	206-473-8221

We can't figure out which phoneNumber corresponds to which person!



# LOSSLESS-JOIN DECOMPOSITION



A natural join is a join on the same attribute names

A schema decomposition is **lossless-join** if for any initial instance  $R$ ,  $R = R'$

# THE CHASE ALGORITHM

The **chase** algorithm is a classic database technique that can be used to check for lossless-join decomposition

## Running example

- relation  $\mathbf{R}(A, B, C, D, E)$
- FDs:  $A \rightarrow B, C$      $D \rightarrow E$

**Question:** is the following decomposition lossless-join?

$\mathbf{R}_1(A, D)$        $\mathbf{R}_2(A, B, C)$        $\mathbf{R}_3(D, E)$

# CHASE: INITIALIZATION

- We create a table with the attributes of the original relation
- We add one row for each relation we split to

	A	B	C	D	E
$R_1(A, D)$	a	$b_1$	$c_1$	d	$e_1$
$R_2(A, B, C)$	a	b	c	$d_2$	$e_2$
$R_3(D, E)$	$a_3$	$b_3$	$c_3$	d	e

Attribute not in the relation  
gets a subscript

Attribute in the relation - no subscript

# CHASE: MAIN ALGORITHM

At every iteration, we check whether an FD is violated, and if so, we “force” it to hold

- If one has a subscript and the other not, we remove the subscript
- If both have a subscript, we make one subscript equal to the other

	A	B	C	D	E
$R_1(A, D)$	a	$b_1$	$c_1$	d	$e_1 \rightarrow e$
$R_2(A, B, C)$	a	b	c	$d_2$	$e_2$
$R_3(D, E)$	$a_3$	$b_3$	$c_3$	d	e

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

The FD  $D \rightarrow E$  is violated, so we need to drop the subscript from the first row

# CHASE: MAIN ALGORITHM

$A \rightarrow B, C$

$D \rightarrow E$

A	B	C	D	E
a	$b_1$	$c_1$	d	$e_1$
a	b	c	$d_2$	$e_2$
$a_3$	$b_3$	$c_3$	d	e

$D \rightarrow E$

A	B	C	D	E
a	$b_1$	$c_1$	d	e
a	b	c	$d_2$	$e_2$
$a_3$	$b_3$	$c_3$	d	e

$A \rightarrow B, C$

A	B	C	D	E
a	b	c	d	e
a	b	c	$d_2$	$e_2$
$a_3$	$b_3$	$c_3$	d	e

At the end of the chase:

- If there is a row without subscripts, we can say that the decomposition is lossless-join
- otherwise, it is not

# MORE EXAMPLES

- relation  $\mathbf{R}(A, B, C, D)$
- FD  $A \rightarrow B, C$

	A	B	C	D
$\mathbf{R}_1$				
$\mathbf{R}_2$				

Lossless-join

- decomposition into  $\mathbf{R}_1(A, B, C)$  and  $\mathbf{R}_2(A, D)$

**Not** lossless-join

- decomposition into  $\mathbf{R}_1(A, B, C)$  and  $\mathbf{R}_2(D)$

# DEPENDENCY PRESERVING

Given  $\mathbf{R}$  and a set of FDs  $F$ , we decompose  $\mathbf{R}$  into  $\mathbf{R}_1$  and  $\mathbf{R}_2$ . Suppose:

- $\mathbf{R}_1$  has a set of FDs  $F_1$
- $\mathbf{R}_2$  has a set of FDs  $F_2$
- $F_1$  and  $F_2$  are computed from  $F$

A decomposition is **dependency preserving** if by enforcing  $F_1$  over  $\mathbf{R}_1$  and  $F_2$  over  $\mathbf{R}_2$ , we can enforce  $F$  over  $\mathbf{R}$

# A NOTE ON FDS OF SPLIT RELATIONS

Given  $\mathbf{R}$  and a set of FDs  $F$ , we decompose  $\mathbf{R}$  into  $\mathbf{R}_1$  and  $\mathbf{R}_2$ . How do we find the FDs  $F_1$  that hold for  $\mathbf{R}_1$ ?

- It is not enough to only keep the FDs from  $F$  with attributes in  $\mathbf{R}_1$
- Instead, we need to find the non-trivial FDs in the fd closure of  $F$  with attributes in  $\mathbf{R}_1$

**Example:**  $\mathbf{R}(A, B, C)$  with FDs:  $A \rightarrow B \quad B \rightarrow C$

- For  $\mathbf{R}_1(A, C) \quad F_1 = A \rightarrow C$



# GOOD EXAMPLE

---

**Person**(SSN, name, age, canDrink)

- $SSN \rightarrow name, age$
- $age \rightarrow canDrink$

decomposes into

- $R_1(SSN, name, age)$ 
  - $SSN \rightarrow name, age$
- $R_2(age, canDrink)$ 
  - $age \rightarrow canDrink$

# BAD EXAMPLE

$R(A, B, C)$

- $A \rightarrow B$
- $B, C \rightarrow A$

Decomposes into:

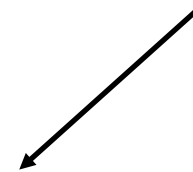
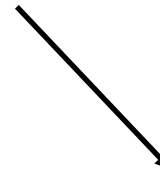
- $R_1(A, B)$ 
  - $A \rightarrow B$
- $R_2(A, C)$ 
  - no FDs here!!

$R_1$

A	B
a <sub>1</sub>	b
a <sub>2</sub>	b

$R_2$

A	C
a <sub>1</sub>	c
a <sub>2</sub>	c



*recover*

A	B	C
a <sub>1</sub>	b	c
a <sub>2</sub>	b	c

The recovered table  
violates  $B, C \rightarrow A$

# NORMAL FORMS

---

A **normal form** represents a “good” schema design:

- 1NF (flat tables/atomic values)
- 2NF
- 3NF
- **BCNF**
- 4NF
- ...

more  
restrictive



---

# BCNF DECOMPOSITION

---

# BOYCE-CODD NORMAL FORM (BCNF)

A relation **R** is in **BCNF** if whenever  $X \rightarrow B$  is a non-trivial FD, then  $X$  is a **superkey** in **R**

**Equivalent definition:** for every attribute set  $X$

- either  $X^+ = X$
- or  $X^+ = \text{all attributes}$

# BCNF EXAMPLE 1

SSN	name	age	phoneNumber
934729837	Paris	24	608-374-8422
934729837	Paris	24	603-534-8399
123123645	John	30	608-321-1163
384475687	Arun	20	206-473-8221

$SSN \rightarrow name, age$

- **key** =  $\{SSN, phoneNumber\}$
- $SSN \rightarrow name, age$  is a “bad” FD
- The above relation is **not** in BCNF!

# BCNF EXAMPLE 2

SSN	name	age
934729837	Paris	24
123123645	John	30
384475687	Arun	20

$SSN \rightarrow name, age$

- **key** = { $SSN$ }
- The above relation is in BCNF!

# BCNF EXAMPLE 3

SSN	phoneNumber
934729837	608-374-8422
934729837	603-534-8399
123123645	608-321-1163
384475687	206-473-8221

- **key** =  $\{SSN, phoneNumber\}$
- The above relation is in BCNF!
- Is it possible that a binary relation is not in BCNF?

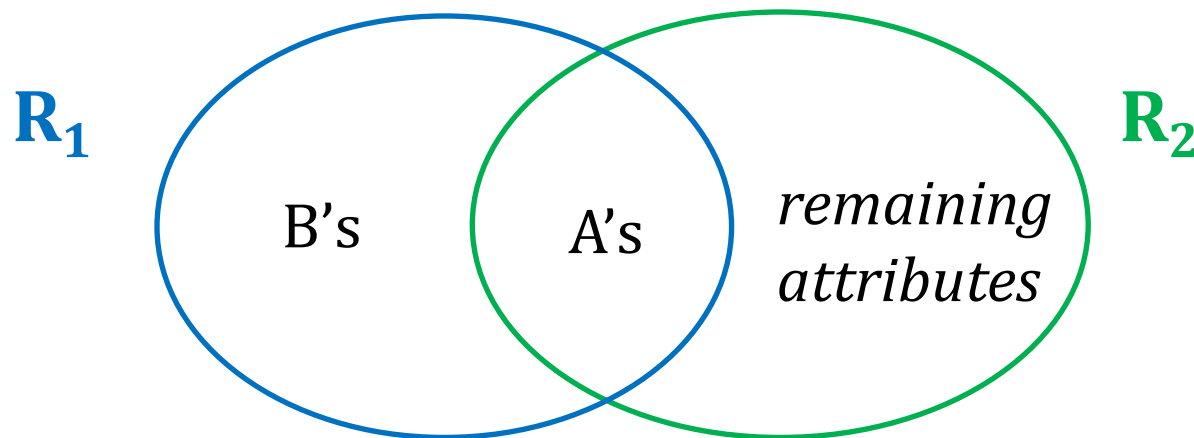


# BCNF DECOMPOSITION

- Find an FD that violates the BCNF condition

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

- Decompose **R** to **R<sub>1</sub>** and **R<sub>2</sub>**:

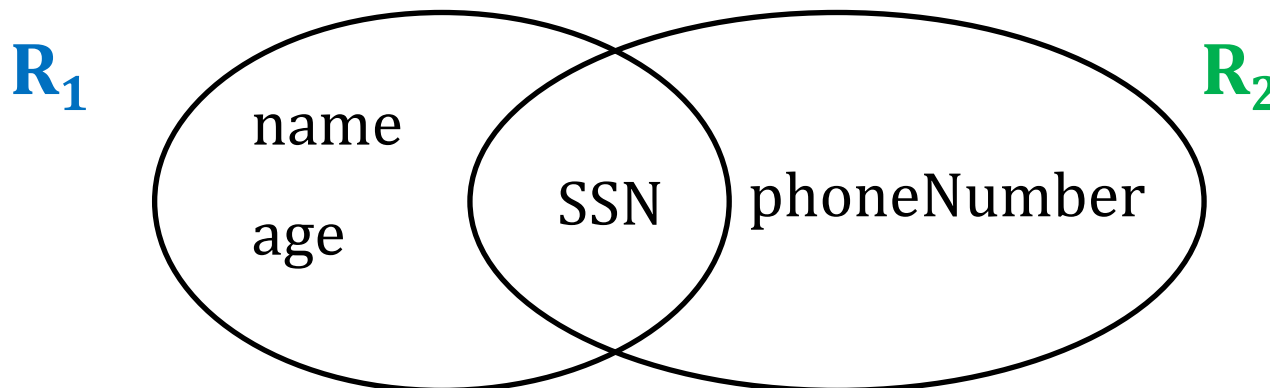


- Continue until no BCNF violations are left

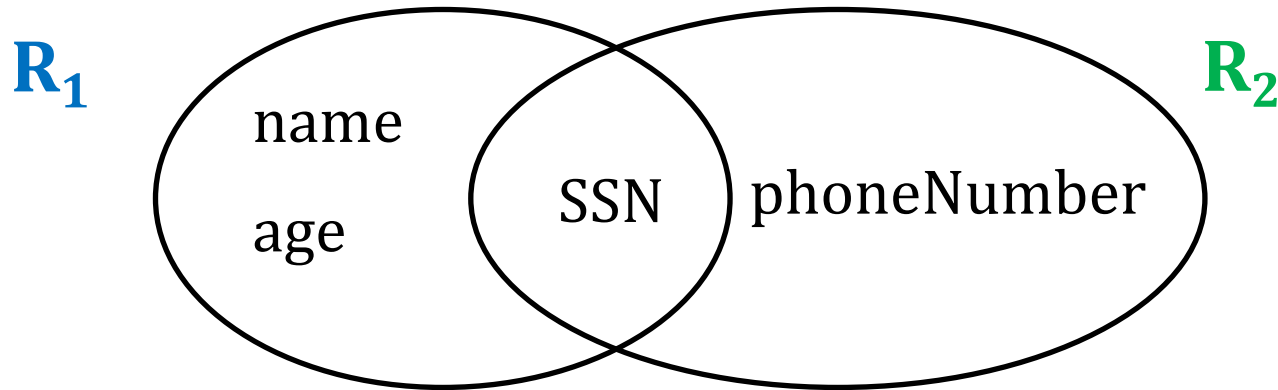
# EXAMPLE

SSN	name	age	phoneNumber
934729837	Paris	24	608-374-8422
934729837	Paris	24	603-534-8399
123123645	John	30	608-321-1163
384475687	Arun	20	206-473-8221

- The FD  $SSN \rightarrow name, age$  violates BCNF
- Split into two relations  $R_1$ ,  $R_2$  as follows:



# EXAMPLE CONT'D



$SSN \rightarrow name, age$

SSN	name	age
934729837	Paris	24
123123645	John	30
384475687	Arun	20

SSN	phoneNumber
934729837	608-374-8422
934729837	603-534-8399
123123645	608-321-1163
384475687	206-473-8221

---

# BCNF DECOMPOSITION PROPERTIES

---

The BCNF decomposition:

- removes certain types of redundancy
- is **lossless-join**
- is **not always** dependency preserving

# BCNF IS LOSSLESS-JOIN

---

Example:

$R(A, B, C)$  with  $A \rightarrow B$  decomposes into:

$R_1(A, B)$  and  $R_2(A, C)$

- The BCNF decomposition always satisfies the lossless-join criterion!

# BCNF IS NOT DEPENDENCY PRESERVING

$R(A, B, C)$

- $A \longrightarrow B$
- $B, C \longrightarrow A$

There may not exist any BCNF decomposition that is FD preserving!

The BCNF decomposition is:

- $R_1(A, B)$  with FD  $A \longrightarrow B$
- $R_2(A, C)$  with no FDs

# BCNF EXAMPLE (1)

---

**Books** (author, gender, booktitle, genre, price)

- $author \rightarrow gender$
- $booktitle \rightarrow genre, price$

What is the candidate key?

- $(author, booktitle)$  is the only one!

Is is in BCNF?

- **No**, because the left hand side of both (not trivial) FDs is not a superkey!

# BCNF EXAMPLE (2)

**Books** (author, gender, booktitle, genre, price)

- $author \rightarrow gender$
- $booktitle \rightarrow genre, price$

Splitting **Books** using the FD  $author \rightarrow gender$ :

- **Author** (author, gender)

FD:  $author \rightarrow gender$  **in BCNF!**

- **Books2** (author, booktitle, genre, price)

FD:  $booktitle \rightarrow genre, price$  **not in BCNF!**



# BCNF EXAMPLE (3)

**Books** (author, gender, booktitle, genre, price)

- $author \rightarrow gender$
- $booktitle \rightarrow genre, price$

Splitting **Books** using the FD  $author \rightarrow gender$ :

- **Author** (author, gender)  
FD:  $author \rightarrow gender$  **in BCNF!**
- Splitting **Books2** (author, booktitle, genre, price):
  - **BookInfo** (booktitle, genre, price)  
FD:  $booktitle \rightarrow genre, price$  **in BCNF!**
  - **BookAuthor** (author, booktitle) **in BCNF!**

---

# IS NORMALIZATION ALWAYS GOOD?

---

- **Example:** suppose A and B are always used together, but normalization says they should be in different tables
  - decomposition might produce unacceptable performance loss
- **Example:** data warehouses
  - huge historical DBs, rarely updated after creation
  - joins expensive or impractical