

Schema Decomposition and Normal Forms

Database Design

- Some schemas are better than others
- Design by decomposition
 - Start with "mega" relations containing all information
 - Then decompose them into smaller, better relations containing the same information
 - Do the decomposition automatically
- Automatic decomposition
 1. Designer specifies mega-relations plus *properties of the data*
 2. System decomposes based on properties
 3. Final set of relations satisfies *normal form* -- guarantees no anomalies, no lost information

Database Anomalies

- Redundancy
- Update anomaly
- Insertion anomaly
- Deletion anomaly

Roll	Name	CID	CNAME	ClassRoom	Slot	Book	Dept
23CS10003	Aryan	CS30202	DBMS	Nalanda 123	Tuesday	Korth	CSE
23CS10004	Ashutosh	CS30202	DBMS	Nalanda 123	Tuesday	Korth	CSE
23CS10063	Sanjana	CS60045	AI	Nalanda 224	Wednesday	Russel	CSE
23CS10091	Sibasis	CS60070	SLT	Nalanda 411	Wednesday	Ben David	History

Schema Decomposition

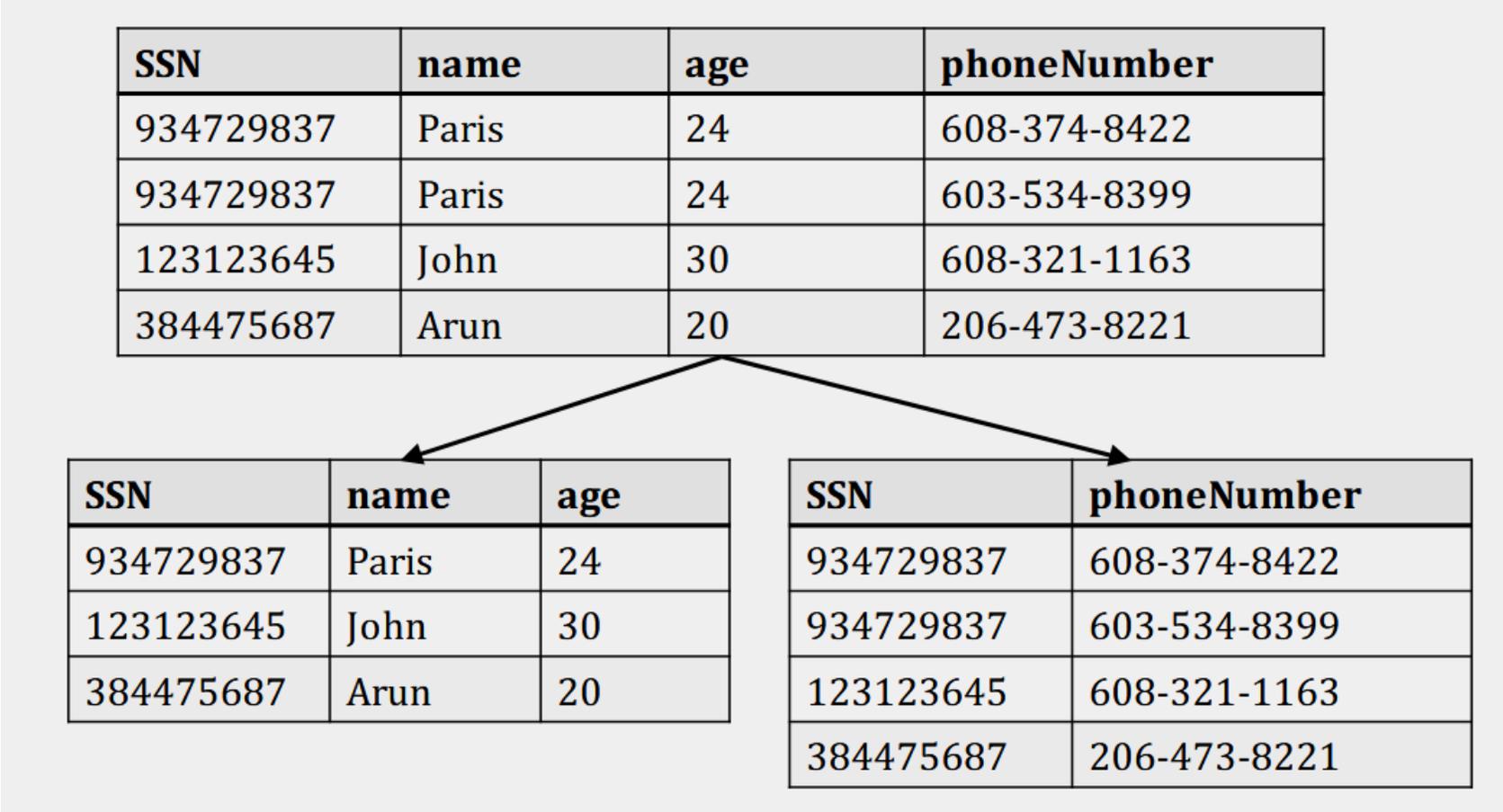
We **decompose** a relation $\mathbf{R}(A_1, \dots, A_n)$ by creating

- $\mathbf{R}_1(B_1, \dots, B_m)$
- $\mathbf{R}_2(C_1, \dots, C_l)$
- where $\{B_1, \dots, B_m\} \cup \{C_1, \dots, C_l\} = \{A_1, \dots, A_n\}$

- 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_l

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



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 Properties

- 1.No anomalies in decomposed relations
- 2.Can reconstruct all original information

Decomposition Goals

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

Information Loss

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

Decompose into:

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

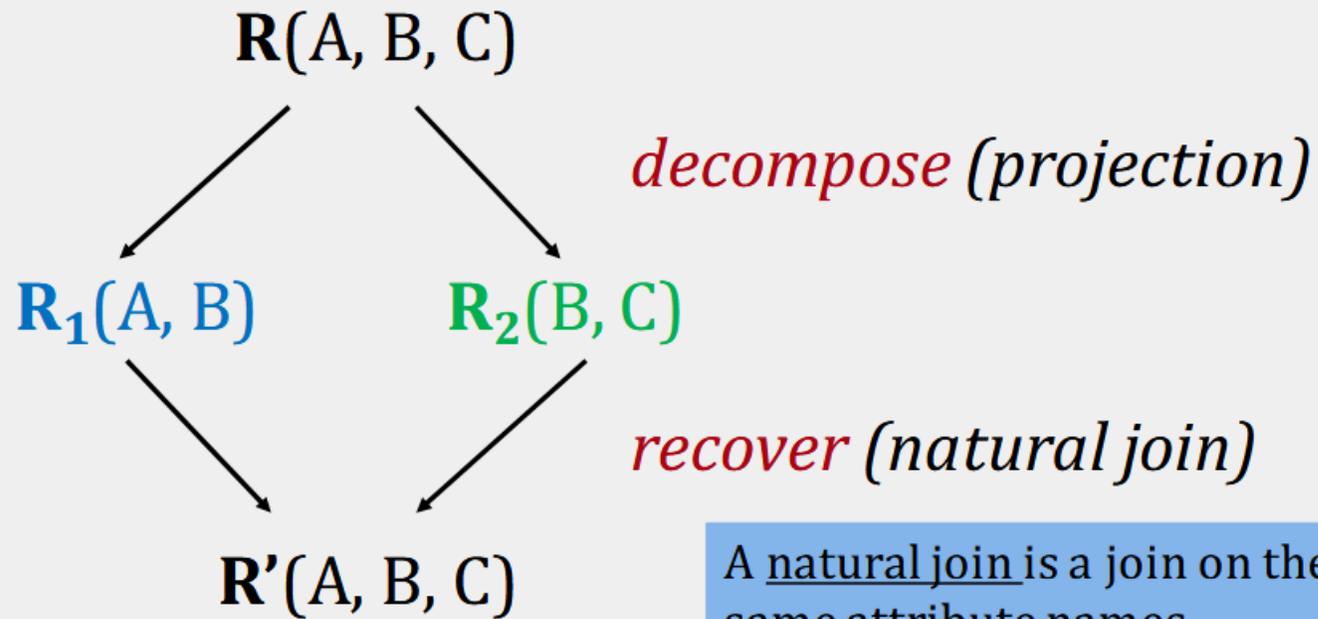
$R_2(\text{age, phoneNumber})$

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 \mathbf{R} , $\mathbf{R} = \mathbf{R}'$

Lossless-Join Criterion

Starting with:

- a relation $\mathbf{R}(\mathbf{A})$ + set F of FDs
- a decomposition of \mathbf{R} into $\mathbf{R}_1(\mathbf{A}_1)$ and $\mathbf{R}_2(\mathbf{A}_2)$

we say that a decomposition is **lossless-join if and only if** at least one of the following FDs is in F^+ (the closure of F) :

1. $\mathbf{A}_1 \cap \mathbf{A}_2 \rightarrow \mathbf{A}_1$
2. $\mathbf{A}_1 \cap \mathbf{A}_2 \rightarrow \mathbf{A}_2$

Example

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

Lossless-join

- decomposition into $\mathbf{R}_1(A, B, C)$ and $\mathbf{R}_2(A, D)$
- $\{A, B, C\} \cap \{A, D\} = \{A\}$
- For \mathbf{R}_1 we have indeed $A \rightarrow B, C$

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}

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 ₁	b
a ₂	b

R_2

A	C
a ₁	c
a ₂	c



recover

A	B	C
a ₁	b	c
a ₂	b	c

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

Database Design

- Design one single mega table to include all information
- Mega table may have certain anomalies which leads to inconsistencies in database
- Decompose table into smaller tables automatically
 - Smaller tables should have certain desired properties
 - Smaller tables should be in Normal forms

Normal Forms

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

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

more
restrictive



Boyce-Codd Normal Form (BCNF)

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

Equivalent definition: for every attribute set X

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

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

$SSN \rightarrow name, age$

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

Example

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

$SSN \rightarrow name, age$

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

Example

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

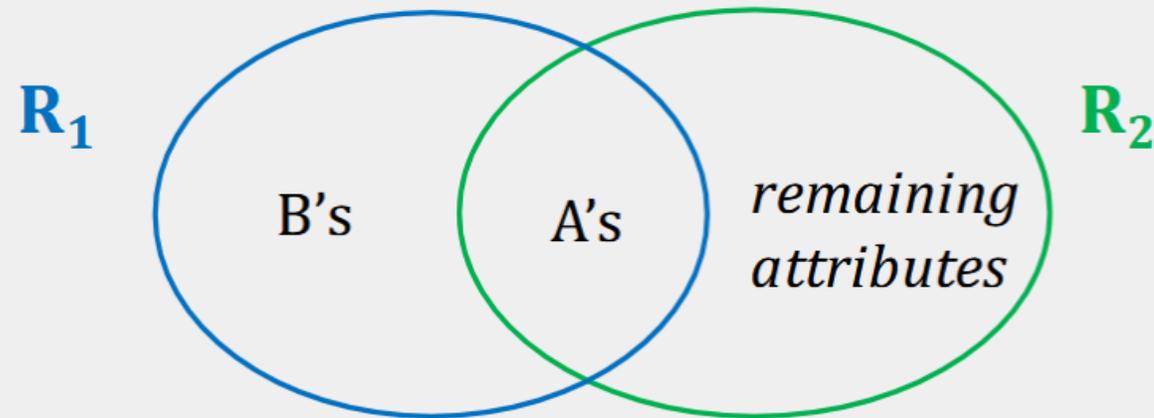
- **key** = $\{SSN, phoneNumber\}$
- The above relation is in BCNF!
- Is it possible that a binary relation is not in BCNF?

BCNF Decomposition

- Find an FD that violates the BCNF condition

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

- Decompose \mathbf{R} to \mathbf{R}_1 and \mathbf{R}_2 :

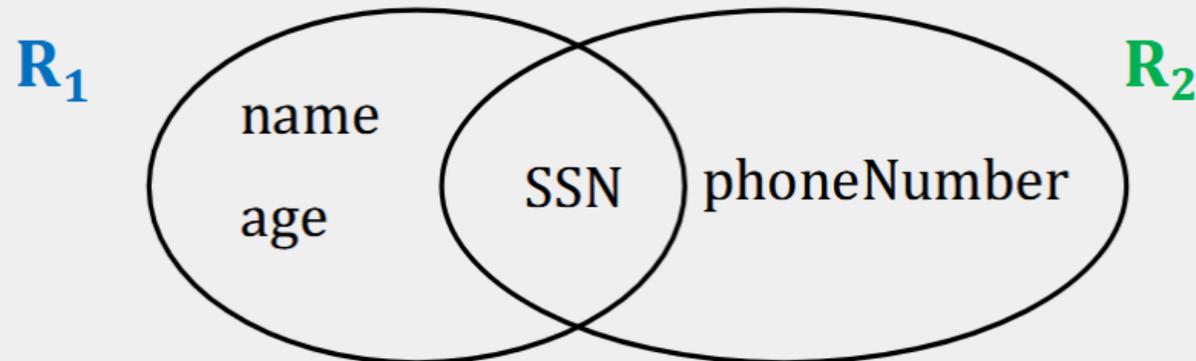


- Continue until no BCNF violations are left

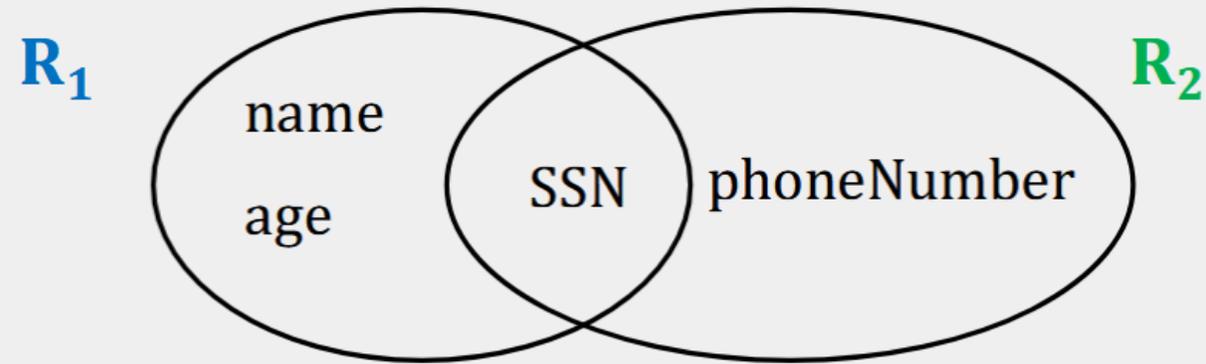
Example

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



$SSN \rightarrow name, age$

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

Properties of BCNF Decomposition

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 Not Dependency Preserving

R(A, B, C)

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

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

The BCNF decomposition is:

- **R₁(A, B)** with FD $A \rightarrow B$
- **R₂(A, C)** with no FDs

Example

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

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

Example

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!**

Example

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!

Third Normal Form (3NF)

A relation \mathbf{R} is in **3NF** if whenever $X \rightarrow A$, one of the following is true:

- $A \in X$ (trivial FD)
- X is a superkey
- A is part of some key of \mathbf{R} (prime attribute)

BCNF implies 3NF !!

3NF

- **Example:** $\mathbf{R}(A, B, C)$ with $A, B \rightarrow C$ and $C \rightarrow A$
 - is in 3NF. Why?
 - is not in BCNF. Why?
- Compromise used when BCNF not achievable: *aim for BCNF and settle for 3NF*
- Lossless-join and dependency preserving decomposition into a collection of 3NF relations is always possible!

3NF Decomposition Algorithm

1. Apply the algorithm for **BCNF decomposition** until all relations are in 3NF (we can stop earlier than BCNF)
2. Compute a **minimal basis** F' of F
3. For each non-preserved FD $X \rightarrow A$ in F' , add a new relation $R(X, A)$

Example

Start with relation **R** (A, B, C, D) with FDs:

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

Step 1: find a BCNF decomposition

- **R1** (B, C)
- **R2** (A, B, D)

Example

Start with relation \mathbf{R} (A, B, C, D) with FDs:

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

Step 2: compute a minimal basis of the original set of FDs:

- $A \rightarrow D$
- $B \rightarrow C$
- $D \rightarrow A$
- $D \rightarrow B$

Example

Start with relation **R** (A, B, C, D) with FDs:

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

Step 3: add a new relation for any FD in the basis that is not satisfied:

- all the dependencies in F' are satisfied!
- the resulting decomposition **R1, R2** is also BCNF!

Fourth Normal Form

- A relation schema R is in **4NF** with respect to a set D of functional and multivalued dependencies if for all multivalued dependencies in D^+ of the form $\alpha \twoheadrightarrow \beta$, where $\alpha \subseteq R$ and $\beta \subseteq R$, at least one of the following hold:
 - $\alpha \twoheadrightarrow \beta$ is trivial (i.e., $\beta \subseteq \alpha$ or $\alpha \cup \beta = R$)
 - α is a superkey for schema R
- If a relation is in 4NF it is 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