CS460: Intro to Database Systems

# Class 7: *Decomposition & Schema Normalization*

Instructor: Manos Athanassoulis

https://bu-disc.github.io/CS460/

# Review: Database Design

<span style="color:red">Requirements Analysis</span>

user needs; what must database do?

<span style="color:red">Conceptual Design</span>

high level description (often done w/ ER model)

<span style="color:red">Logical Design</span>

translate ER into DBMS data model

<span style="color:red">**Schema Refinement**</span>

**consistency, normalization**

<span style="color:red">Physical Design</span>

indexes, disk layout

# Why schema refinement

## what is a bad schema?
*a schema with redundancy!*

## why?
redundant storage & insert/update/delete anomalies

## how to fix it?
**normalize** the schema by decomposing
normal forms: BCNF, 3NF, …

# Motivating Example

| SSN | Name | Salary | Telephone |
|---|---|---|---|
| 987-00-8761 | John | 65K | 857-555-1234 |
| 987-00-8761 | John | 65K | 857-555-8800 |
| 123-00-9876 | Anna | 80K | 617-555-9876 |
| 787-00-4321 | John | 25K | 617-555-3761 |

SSN $\rightarrow$ Name, Salary

# Motivating Example

| SSN | Name | Salary | Telephone |
|---|---|---|---|
| 987-00-8761 | John | 65K | 857-555-1234 |
| 987-00-8761 | John | 65K | 857-555-8800 |
| 123-00-9876 | Anna | 80K | 617-555-9876 |
| 787-00-4321 | John | 25K | 617-555-3761 |

SSN → Name, Salary

| SSN | Name | Salary |
|---|---|---|
| 987-00-8761 | John | 65K |
| 123-00-9876 | Anna | 80K |
| 787-00-4321 | John | 25K |

| SSN | Telephone |
|---|---|
| 987-00-8761 | 857-555-1234 |
| 987-00-8761 | 857-555-8800 |
| 123-00-9876 | 617-555-9876 |
| 787-00-4321 | 617-555-3761 |

5

# Motivating Example 2

| name | category | color | price | department |
|---|---|---|---|---|
| iPhone | smartphone | black | 600 | phones |
| Lenovo Yoga | laptop | grey | 800 | computers |
| unifi | networking | white | 150 | computers |
| unifi | cables | white | 10 | stationary |
| OnePlus | smartphone | silver | 450 | phones |

name, category → price, color                     category → department

6

# Motivating Example 2

| name | category | color | price | department |
|------|----------|-------|-------|------------|
| iPhone | smartphone | black | 600 | phones |
| Lenovo Yoga | laptop | grey | 800 | computers |
| unifi | networking | white | 150 | computers |
| unifi | cables | white | 10 | stationary |
| OnePlus | smartphone | silver | 450 | phones |

name, category → price, color

category → department

| name | category | color | price |
|------|----------|-------|-------|
| iPhone | smartphone | black | 600 |
| Lenovo Yoga | laptop | grey | 800 |
| unifi | networking | white | 150 |
| unifi | cables | white | 10 |
| OnePlus | smartphone | silver | 450 |

| category | department |
|----------|------------|
| laptop | computers |
| networking | computers |
| cables | stationary |
| smartphone | phones |

7

# Reminder: Reasoning for FDs

Assume a relation R with attributes A, B, C

(1) reflexivity             e.g., AB → B

(2) augmentation          e.g., if A → B then AC → BC

(3) transitivity           e.g., if A → B and B → C then A → C

(4) union                 e.g., if A → B and A → C then A → BC

(5) decomposition          e.g., if A → BC then A → B  and  A → C

**FD closure of F, $F^+$**: is the set of all FDs that are implied by F
**attr. closure of X**: the set of all attributes that are determined by X
**minimal cover**: subset $S$ of $F^+$ such that $S^+ = F^+$

*"chopping the relation into pieces using FDs"*

# DECOMPOSITION

# Decomposition

**Formally**

we decompose $R(A_1, ..., A_n)$ by creating:

$R_1(B_1, ..., B_m)$

$R_2(C_1, ..., C_k)$

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

the instance of $R_1$ is the <u>projection</u> of $R$ onto $B_1, ..., B_m$

the instance of $R_2$ is the <u>projection</u> of $R$ onto $C_1, ..., C_k$

# "Good" Decomposition

(1) minimize redundancy

(2) avoid information loss (*lossless-join*)

(3) preserve FDs (*dependency preserving*)

(4) ensure good query performance

# Information Loss

| SSN | Name | Salary | Telephone |
|---|---|---|---|
| 987-00-8761 | John | 65K | 857-555-1234 |
| 987-00-8761 | John | 65K | 857-555-8800 |
| 123-00-9876 | Anna | 80K | 617-555-9876 |
| 787-00-4321 | John | 25K | 617-555-3761 |

Decompose into:
$R_1$(SSN, Name, Salary)
$R_2$(Name, Telephone)

| SSN | Name | Salary |
|---|---|---|
| 987-00-8761 | John | 65K |
| 123-00-9876 | Anna | 80K |
| 787-00-4321 | John | 25K |

| Name | Telephone |
|---|---|
| John | 857-555-1234 |
| John | 857-555-8800 |
| Anna | 617-555-9876 |
| John | 617-555-3761 |

can we reconstruct R?

# Lossless Decomposition

**R**(A,B,C)

*decompose*

**R₁**(A,B)          **R₂**(B,C)

*recover (join on B)*

**R'**(A,B,C)

the decomposition is *lossless-join* if
for any initial instance **R, R = R'**

# Lossless Criterion

given:

- a relation **R**(A)
- a set *F* of FDs
- a decomposition of **R** into $R_1(A_1)$ and $R_2(A_2)$

the decomposition is *lossless-join* **if and only if**

at least one of the following FDs is in $F^+$ (closure of F):

(1) $A_1 \cap A_2 \rightarrow A_1$

(2) $A_1 \cap A_2 \rightarrow A_2$

15

# Example

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

$A \rightarrow BC \quad A \rightarrow AB \quad A \rightarrow AC$

$A \rightarrow ABC$

Relation **R**(A, B, C, D)

FD $A \rightarrow BC$

*what is the $F^+$?*

**lossy**

decomposition into **R₁**(A, B, C) and **R₂**(D)

$A_1 \cap A_2$ empty set

**lossless-join**

decomposition into **R₁**(A, B, C) and **R₂**(A, D)

$A_1 \cap A_2 = A$  and $A_1 = A,B,C$

$A \rightarrow ABC$ *is in $F^+$*

16

# Dependency Preserving

given **R** and a set of FDs *F*, we decompose
**R** into **R$_1$** and **R$_2$**. Suppose:

**R$_1$** has a set of FDs *F$_1$*
**R$_2$** has a set of FDs *F$_2$*
*F$_1$ and F$_2$ are computed from F*

it is <u>dependency preserving</u> if by enforcing
*F$_1$* over **R$_1$** and *F$_2$* over **R$_2$** , we can enforce *F* over **R**

# (Good) Example

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

*SSN* $\rightarrow$ *name, age*

*age* $\rightarrow$ *canDrink*

*what is a **dependency preserving** decomposition?*

**R$_1$**(SSN, name, age)        and      **R$_2$**(age, canDrink)

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

**Is it also lossless-join?**

**Yes!** $A_1 \cap A_2$ = age  and $A_2$ = age, canDrink

age $\rightarrow$ age, canDrink *is in F$^+$*

18

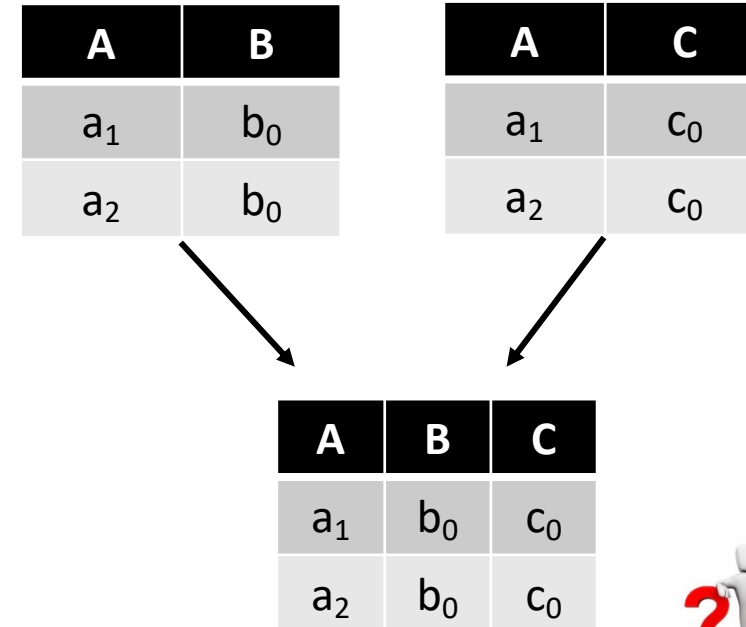# (Bad) Example

**R** (A, B, C)

$A \rightarrow B$

$BC \rightarrow A$

*not dependency preserving*

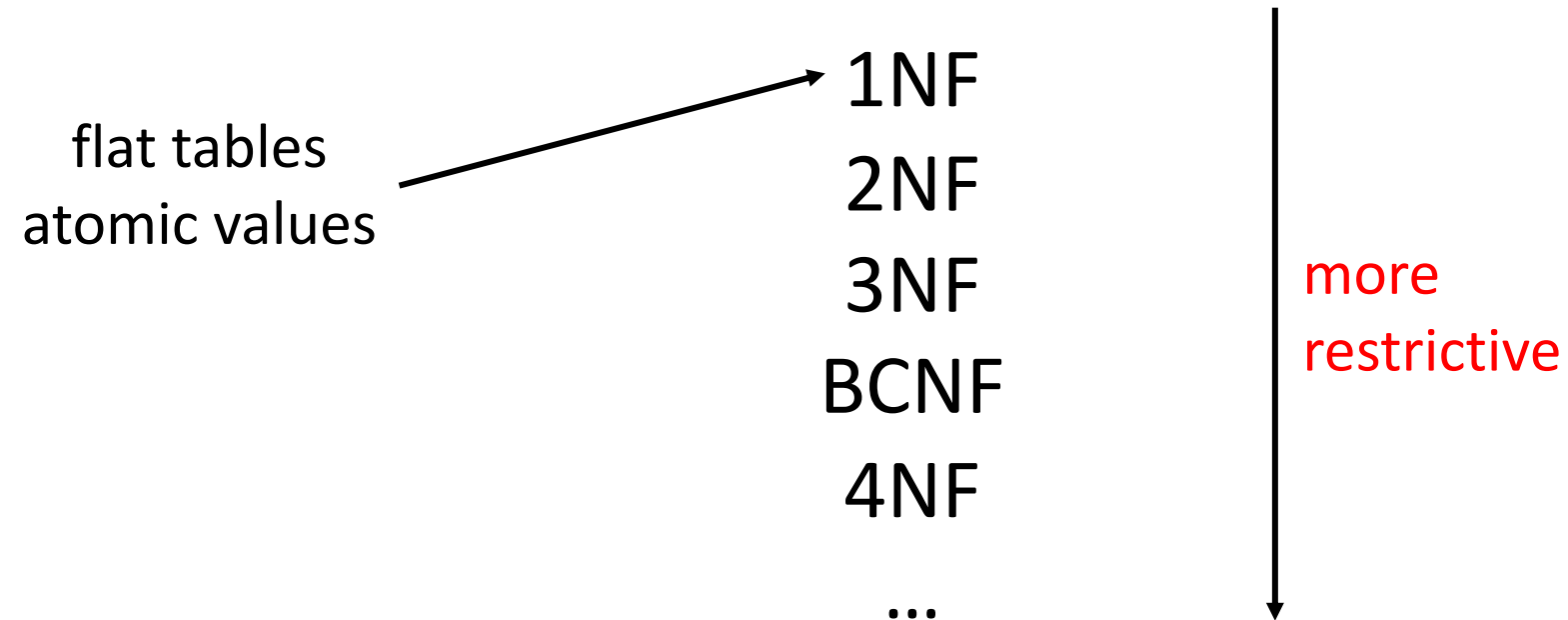**R$_1$**(A, B)      and      **R$_2$**(A, C)

$A \rightarrow B$                    *no FDs!*

| A | B |
|---|---|
| $a_1$ | $b_0$ |
| $a_2$ | $b_0$ |

| A | C |
|---|---|
| $a_1$ | $c_0$ |
| $a_2$ | $c_0$ |

| A | B | C |
|---|---|---|
| $a_1$ | $b_0$ | $c_0$ |
| $a_2$ | $b_0$ | $c_0$ |

the table violates
$BC \rightarrow A$

# Normal Forms

How "good" is a schema design?
follows normal forms

flat tables
atomic values →

1NF
2NF
3NF
BCNF
4NF
...

more
restrictive

# Normal Forms

How "good" is a schema design?
follows normal forms

flat tables
atomic values

1NF
2NF
**3NF**
**BCNF**
4NF
...

<span style="color:red">more
restrictive</span>

# Boyce-Codd Normal Form (BCNF)

given a relation **R**($A_1,...,A_n$),
a set of FDs *F*, and X⊆{$A_1,...,A_n$}
**R** is in BCNF if *∀ X→A* one of the two holds:

- A ∈ X (that is, it is a trivial FD)
- X is a <u>superkey</u>

<u>in other words</u>: *∀ **non-trivial FD** X→A, **X is a superkey in R***

# BCNF - Example

| SSN | Name | Salary | Telephone |
|-----|------|--------|-----------|
| 987-00-8761 | John | 65K | 857-555-1234 |
| 987-00-8761 | John | 65K | 857-555-8800 |
| 123-00-9876 | Anna | 80K | 617-555-9876 |
| 787-00-4321 | John | 25K | 617-555-3761 |

*SSN → Name, Salary*
*key: {SSN, Telephone}*
    *FD is not trivial!*
    *so, is SSN a superkey?*
    *no! it is **not** in **BCNF***

# BCNF - Example 2

| SSN | Name | Salary |
|---|---|---|
| 987-00-8761 | John | 65K |
| 123-00-9876 | Anna | 80K |
| 787-00-4321 | John | 25K |

*SSN  → Name, Salary*
*key: {SSN}*

*FD is not trivial!*
*so, is SSN a superkey?*
*yes! it is in **BCNF***

# BCNF - Example 3

| SSN | Telephone |
|---|---|
| 987-00-8761 | 857-555-1234 |
| 987-00-8761 | 857-555-8800 |
| 123-00-9876 | 617-555-9876 |
| 787-00-4321 | 617-555-3761 |

*key: {SSN, Telephone}    the relation is in **BCNF***

*why?*

*no FDs*

*Is it possible a binary relation*
*to <u>not</u> be in **BCNF**?*

25

# Binary Relations always BCNF

**R** (A,B)
excluding all trivial FDs, there are three cases:
(1) **R** has no FD
(2) **R** has one FD, either $A \rightarrow B$ or $B \rightarrow A$, or,
(3) **R** has two FDs, $A \rightarrow B$ and $B \rightarrow A$

(1) trivially in BCNF
(2) in either LHS is the key (hence, superkey)
(3) both, A and B candidate keys

# BCNF Decomposition Algorithm

(1) find a FD that violates BCNF:

$$A_1, ..., A_n \rightarrow B_1, ..., B_m$$

(2) decompose **R** to **R$_1$** and **R$_2$**

**R$_1$**$(A_1, ..., A_n, B_1, ..., B_m)$
**R$_2$**$(A_1, ..., A_n,$ all other attributes of **R**$)$

(3) repeat until no BCNF violations are left
(in new tables as well)

27

# Our favorite example!

| SSN | Name | Salary | Telephone |
|---|---|---|---|
| 987-00-8761 | John | 65K | 857-555-1234 |
| 987-00-8761 | John | 65K | 857-555-8800 |
| 123-00-9876 | Anna | 80K | 617-555-9876 |
| 787-00-4321 | John | 25K | 617-555-3761 |

*SSN → Name, Salary* violates BCNF

$A_1$ = SSN, $B_1$ = *Name,* $B_2$ = *Salary*

Split in two relations:

**R₁**(SSN, Name, Salary)

**R₂**(SSN, Telephone)

# Our favorite example!

| SSN | Name | Salary | Telephone |
|---|---|---|---|
| 987-00-8761 | John | 65K | 857-555-1234 |
| 987-00-8761 | John | 65K | 857-555-8800 |
| 123-00-9876 | Anna | 80K | 617-555-9876 |
| 787-00-4321 | John | 25K | 617-555-3761 |

| SSN | Name | Salary |
|---|---|---|
| 987-00-8761 | John | 65K |
| 123-00-9876 | Anna | 80K |
| 787-00-4321 | John | 25K |

| SSN | Telephone |
|---|---|
| 987-00-8761 | 857-555-1234 |
| 987-00-8761 | 857-555-8800 |
| 123-00-9876 | 617-555-9876 |
| 787-00-4321 | 617-555-3761 |

29

# BCNF Decomposition Properties

removes [certain types of] redundancy

is **lossless-join**

is <u>not always</u> <span style="color:red">dependency preserving</span>

# BCNF – Lossless Join

## Example

**R** (A, B, C) and FD: $A \rightarrow B$

superkey(s) of the relation?

$\{A, C\}^+$, $\{A, B, C\}^+ = \{A, B, C\}$

$A \rightarrow B$ *violates* BCNF (A is not a superkey)

so, the BCNF decomposition is :

**R$_1$**(A, B) and **R$_2$**(A, C)

we can reconstruct it!

31

# BCNF – **not** dependency preserving

**Example**

**R** (A, B, C), FDs: $A \rightarrow B$ and $BC \rightarrow A$

superkey(s) of the relation?

$\{A, C\}^+$, $\{B, C\}^+$, $\{A, B, C\}^+ = \{A, B, C\}$

$BC \rightarrow A$ *is ok, but* $A \rightarrow B$ *violates BCNF*

so, the BCNF decomposition is :

**R$_1$** (A, B) and **R$_2$** (A, C)

$A \rightarrow B$ *is preserved in* **R$_1$**

$BC \rightarrow A$ *is not preserved!*

# BCNF Decomposition Examples

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

candidate key(s)?
{author, booktitle} is the only one

Is it in BCNF?
**No**, because LHS of both FD are not a superkey!

# BCNF Decomposition Examples

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

Splitting using: *author → gender*
**AuthorInfo** (author, gender)
*FD author → gender (in BCNF!)*

**Book2** (author, booktitle, genre, price)
*FD booktitle → genre, price*
is booktitle a superkey?

No! {booktitle, author} is.
So not in BCNF!

# BCNF Decomposition Examples

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

**AuthorInfo** (author, gender)
Further splitting with *booktitle ⟶ genre, price*
~~**Book2** (author, booktitle, genre, price)~~

**BookInfo** (booktitle, genre, price)
  *FD booktitle ⟶ genre, price*  in BCNF!
  is booktitle a superkey?  Yes!
**BookAuthor** (booktitle, author)  binary is in BCNF!

# what if not dependency preserving?

in some cases BCNF decomposition is not dependency preserving

how to address this?

relax the normalization requirements

# Third Normal Form (3NF)

given a relation **R** $(A_1,...,A_n)$,
a set of FDs *F*, and $X \subseteq \{A_1,...,A_n\}$
**R** is in 3NF if $\forall\ X \rightarrow A$ one of the three holds:

- $A \in X$ (that is, it is a trivial FD)
- X is a <u>superkey</u>
- A is <u>part of some candidate key</u> for **R**

is a relation in 3NF also in BCNF?
   *No, but a relation in BCNF is always in 3NF!*

37

# Third Normal Form (3NF)

**Example**
**R** (A, B, C), FDs $C \to A$ and $AB \to C$
  is in 3NF but not in BCNF. Why?

superkeys?
{A, B}, {B, C}, and {A, B, C}

candidate keys?
{A, B} and {B, C}

**Compromise:** aim for BCNF but settle for 3NF
lossless-join & dependency preserving possible

# 3NF Algorithm

(1) apply BCNF until all relations are in 3NF

(2) compute a <u>minimal cover</u> *F'* of *F*

(3) for each non-preserved FD X → A in *F'*
add a new relation **R** (X, A)

# 3NF algorithm example

Assume **R** (A, B, C, D)

$A \rightarrow D$

$AB \rightarrow C$

$AD \rightarrow C$

$B \rightarrow C$

$D \rightarrow AB$

superkeys?

{A} {D} {A, B} {A, D}, …

not {B}

**Step 1:** find a BCNF decomposition

**R$_1$** (B, C)

**R$_2$** (A, B, D)

# 3NF algorithm example

Assume **R** (A, B, C, D)

$A \to D$

$AB \to C$     *can be expressed via: $AB \to AC$ which gives $AB \to A$ and $AB \to C$*

$AD \to C$     *can be expressed via: $D \to AB$, which gives $D \to A$ and $D \to B$ & $B \to C$*

$B \to C$

$D \to AB$     *which is simplified to $D \to A$ and $D \to B$*

**Step 2:** find a minimal cover

$A \to D$

$B \to C$

$D \to A$

$D \to B$

# 3NF algorithm example

Assume **R** (A, B, C, D)

$A \rightarrow D$

$AB \rightarrow C$

$AD \rightarrow C$

$B \rightarrow C$

$D \rightarrow AB$

**Step 3:** add a new relation for not preserved FDs

$A \rightarrow D$                            $R_1$ (B, C)

$B \rightarrow C$                            $R_2$ (A, B, D)

$D \rightarrow A$            all FD are preserved!

$D \rightarrow B$           both are in BCNF!

42

# Is Normalization Always Good?

**Example 1:** suppose A and B are always used together, but normalization puts them in different tables (e.g., hours_worked and hourly_rate)

decomposition might produce <u>unacceptable performance</u> loss

**Example 2:** data warehouses
huge historical DBs, rarely updated after creation
joins expensive or impractical
[we want "flat" tables, a.k.a, denormalized]

# Example

R (C, S, J, D, P, Q, V)

C → SJDPQV

JP → C

~~SD → P~~

J → S

**Step 1:**
**R$_1$** (S, D, P)
**R$_2$** (C, S, J, D, Q, V)

superkeys?

{C}, {J, P}, {D, J}, …

not {S, D}

# Example

R (C, S, J, D, P, Q, V)

C → SJDPQV

JP → C

SD → P

$\boxed{J \rightarrow S}$

**Step 1b:**

**R₁** (S, D, P)

~~**R₂** (C, S, J, D, Q, V)~~

**R₂'** (J, S)

**R₃** (C, J, D, Q, V)

superkeys of **R₂** (C, S, J, D, Q, V) ?

{C}, … not {J}

# Example

R (C, S, J, D, P, Q, V)

C $\rightarrow$ SJDPQV

JP $\rightarrow$ C

SD $\rightarrow$ P

J $\rightarrow$ S

**Step 2:** Minimal Cover

C $\rightarrow$ J,  C $\rightarrow$ D,  C $\rightarrow$ Q,  C $\rightarrow$ V

JP $\rightarrow$ C

SD $\rightarrow$ P

J $\rightarrow$ S

**R$_1$** (S, D, P)
**R$_{2'}$** (J, S)
**R$_3$** (C, J, D, Q, V)
**R$_4$** (J, P, C)

are they all preserved?

**No!**
**Step 3:** need to add **R$_4$** (J, P, C)

# Example

R (C, S, J, D, P, Q, V)

C $\rightarrow$ SJDPQV

JP $\rightarrow$ C

SD $\rightarrow$ P

J $\rightarrow$ S

**Step 2:** Minimal Cover

    C $\rightarrow$ J,  C $\rightarrow$ D,  C $\rightarrow$ Q,  C $\rightarrow$ V

    JP $\rightarrow$ C

    SD $\rightarrow$ P

    J $\rightarrow$ S

**R$_1$** (S, D, P)
**R$_{2'}$** (J, S)
**R$_3$** (C, J, D, Q, V)
**R$_4$** (J, P, C)

are they all preserved?

**No!**
**Step 3:** need to add **R$_4$** (J, P, C)

***did we just introduce redundancy?***

# Lesson!

theory of normalization is a guide

cannot always give a "perfect" solution

redundancy

alternatives

query performance

# Summary

## fix bad schemas (redundancy) by decomposition

lossless-join

dependency preserving

## Desired normal forms

**BCNF:** only superkey FDs

**3NF:** superkey FDs + dependencies to prime attributes in RHS

## **Next:** SQL