# Lecture 12: Database Design (Theory)

CS348 Spring 2025: Introduction to Database Management

> Instructor: **Xiao Hu** Sections: 001, 002, 003

#### Announcements

- Appealing of Assignment 1
  - Check remark request guidelines on Piazza
  - Reach out to corresponding TA, IA (Guy), ISC (Sylvie)
  - Check sample solutions on Learn
- Milestone 1 of Group Project
  - Due on June 19
- Switch-type cut-off for assessment
  - Due on June 19
- Assignment 2
  - Coverage: Lecture 4 Lecture 12
  - Check online office hours on Piazza
  - Due on June 24

# (Recap) Data Redundancy!

 Possible Design: one large table InstructorDep with one row for each instructor

<u>instructorID</u>	name	salary	depName	bldng	budget
111	Alice	5000	CS	DC	20000
222	Bob	4000	Physics	РНҮ	30000
333	Carl	5200	CS	DC	20000
444	Diana	5500	CS	DC	20000

• Decomposition: Break down a complex database schema into smaller, more manageable pieces while maintaining data integrity and relationships

## (Recap) Decompositions <sub>R</sub>



- R is decomposed into  $R_1$  and  $R_2$ 
  - Attribute  $(R_1) \cup \text{Attribute}(R_2) = \text{Attribute}(R)$
  - *R*<sub>1</sub> and *R*<sub>2</sub> are the projections of *R* onto Attribute (*R*<sub>1</sub>) and Attribute(*R*<sub>2</sub>)
- Any decomposition gives  $R \subseteq R_1 \bowtie R_2$ 
  - Lossless decomposition if  $R = R_1 \bowtie R_2$
  - Lossy decomposition if  $R \subset R_1 \bowtie R_2$

For any tuple  $(a, b, c, d) \in R$ ,  $(a, b, c) \in R_1$  and  $(c, d) \in R_2$ ; then  $(a, b, c) \bowtie$  $(c, d) \in R_1 \bowtie R_2$ 

# (Recap) A good decomposition

			structorID	name	salary	depName	bldng	budget
			111	Alice	5000	CS	DC	20000
			222	Bob	4000	Physics	PHY	30000
			333	Carl	5200	CS	DC	20000
			444	Diana	5500	CS	DC	20000
						·		
instructorID	name	salary	depName	2				
instructorID 111	name Alice	salary 5000	depName CS	2	[	depName	bldng	budget
instructorID 111 222	name Alice Bob	salary 5000 4000	depName CS Physics		4	depName CS	bldng DC	budget 20000
instructorID 111 222 333	name Alice Bob Carl	salary 5000 4000 5200	depName CS Physics CS		4	depName CS Physics	bldng DC PHY	budget 20000 30000
instructorID 111 222 333 444	name Alice Bob Carl Diana	salary 5000 4000 5200 5500	depName CS Physics CS CS	≥ ►	<b>⊲</b>	depName CS Physics	bldng DC PHY	budget 20000 30000

# Why can we recover all original tuples by joins?

instructorID	name	salary	depName	bldng	budget
111	Alice	5000	CS	DC	20000
222	Bob	4000	Physics	PHY	30000
333	Carl	5200	CS	DC	20000
444	Diana	5500	CS	DC	20000

## (Recap) A bad decomposition



instructorID	name	salary	depName	bldng	budget
111	Alice	5000	CS	DC	20000
111	Bob	5200	CS	PHY	30000
•••			•••		

Lossy? But I got more rows! Can't tell what's fact and what's not, so we lose information!

## **Outline for Today**

- Functional Dependencies
- Decomposition:
  - Lossless
  - Dependency-preserving
- Boyce-Codd Normal Form (BCNF)
- Third Normal Form (3NF)

#### Lossless decomposition

• We should be able to reconstruct the instance of the original table from the instances of the tables in the decomposition

A decomposition  $\{R_1, R_2\}$  of R is lossless if and only if the common attributes of  $R_1$  and  $R_2$  form a superkey for either schema, i.e.,  $R_1 \cap R_2 \to R_1$  or  $R_1 \cap R_2 \to R_2$ 

> (See Lecture 11) If X is a superkey of R, then  $X \rightarrow R$

### Is this a lossless decomposition?

• Consider  $R = \{sno, sname, city, pno, pname, price\}$  and a decomposition  $(R_1, R_2)$  below:

<u>sno</u>	sname	city	<u>pno</u>	pname	price
S1	Apple	K-W	P1	А	\$25
S1	Apple	K-W	P2	В	\$34
S2	BestBuy	London	Р3	А	\$5
S2	BestBuy	London	•••	•••	•••

Lossy since  $R_1 \cap R_2 =$ {pname} is not a superkey of either  $R_1$  or  $R_2$ 

## Which one is better?

 ${\mathcal F}$  includes:

 $project \rightarrow department$  $department \rightarrow division$  $project \rightarrow division$ 

 Consider Company = {project, department, division} under a set F of FDs, and two different decompositions:

 $R_1 = \{ project, department \}$  $R_2 = \{ department, division \}$ 

 $R_1 = \{ project, department \}$  $R_2 = \{ project, division \}$ 

- Both are lossless. (Why?)
- However, testing FDs is easier on one of them. (Which?)

# Testing FDs

 $\mathcal{F}$  includes: project  $\rightarrow$  department department  $\rightarrow$  division Project  $\rightarrow$  division

 Consider Company = {project, department, division} under a set F of FDs, and two different decompositions:

#### decomposition 1:

 $R_1 = \{ project, department \}$  $R_2 = \{ department, division \}$ 

#### test in $R_1$

project  $\rightarrow$  department department  $\rightarrow$  division project  $\rightarrow$  division

> No need to test since it is implied by the other two FDs

test in  $R_2$ 

decomposition 2:

 $R_1 = \{ project, department \}$  $R_2 = \{ project, division \}$ 

#### test in $R_1$

 $\begin{array}{c} \text{project} \rightarrow \text{department} \\ \text{department} \rightarrow \text{division} \\ \text{project} \rightarrow \text{division} \end{array}$ 

test in  $R_2$ 

11

#### Dependency-preserving decomposition

• We should be able to (explicitly and implicitly) test all dependencies in each base table of the decomposition

Given a schema R and a set  $\mathcal{F}$  of FDs, a decomposition of R is dependency preserving if there is an equivalent set  $\mathcal{F}'$  of FDs to  $\mathcal{F}$ , none FD in  $\mathcal{F}'$  is cross-table in the decomposition.

• Given a set of FDs  $\mathcal{F}$ , we say  $\mathcal{F}'$  is equivalent to  $\mathcal{F}$  if their closures are the same, i.e.,  $\mathcal{F}^+ = \mathcal{F}'^+$ .

# **Outline for Today**

- Functional Dependencies
- Decomposition:
  - Lossless
  - Dependency-preserving
- Boyce-Codd Normal Form (BCNF)
- Third Normal Form (3NF)

## Boyce-Codd Normal Form (BCNF)

- A relation R is in BCNF under  $\mathcal{F}$  if each FD  $X \to Y \in \mathcal{F}^+$  with  $XY \subseteq R$  satisfies:
  - either  $X \to Y$  is trivial, i.e.,  $Y \subseteq X$
  - or X is a super key of R, i.e.,  $X \rightarrow R$

• Is *R* = {sno, sname, city, pno, pname, price} under *F* in BCNF?

 $\mathcal{F}$  includes: sno  $\rightarrow$  sname, city pno  $\rightarrow$  pname sno, pno  $\rightarrow$  price

(See Lecture 11) How to determine if K is a superkey of R?

## Compute BCNF decomposition

Repeat the following until all relations are in BCNF

- Step 1: Find a BCNF violation
  - A relation *R* in the decomposition
  - A non-trivial FD  $X \to Y$  in  $\mathcal{F}^+$  with  $XY \subseteq R$ , where X is not a super key of R
- Step 2: Decompose *R* into *R*<sub>1</sub> and *R*<sub>2</sub>
  - $R_1$  has attributes  $X \cup Y$ ;
  - $R_2$  has attributes  $X \cup Z$ , where Z contains all attributes of R that are in neither X nor Y

#### Example of BCNF decomposition

(See Lecture 11) How to determine if K is a superkey of R?  $\mathcal{F}$  includes: sno  $\rightarrow$  sname, city pno  $\rightarrow$  pname sno, pno  $\rightarrow$  price



#### BCNF helps remove redundancy

<u>sno</u>	sname	city	<u>pno</u>	pname	price
<b>S</b> 1	Apple	K-W	P1	А	\$25
S1	Apple	K-W	P2	В	\$34
S1	Apple	K-W	Р3	А	\$20
S2	BBuy	London	•••		•••
BCNF violation: sno $\rightarrow$ sname, city					

 $\mathcal{F}$  includes: sno  $\rightarrow$  sname, city pno  $\rightarrow$  pname sno, pno  $\rightarrow$  price

<u>sno</u>	<u>pno</u>	pname	price
S1	P1	А	\$25
S1	P2	В	\$34
S1	Р3	А	\$20
S2	•••	•••	•••

<u>sno</u>	sname	city
S1	Apple	K-W
S2	BBuy	London

#### Another example

 $\mathcal{F}$  includes: uid  $\rightarrow$  uname, twitterid twitterid  $\rightarrow$  uid uid, gid  $\rightarrow$  fromDate

UserJoinsGroup (uid, uname, twitterid, gid, fromDate) BCNF violation: uid  $\rightarrow$  uname, twitterid

User (uid, uname, twitterid)

uid  $\rightarrow$  uname, twitterid twitterid  $\rightarrow$  uid

In BCNF (why?)

Member (uid, gid, fromDate)

uid, gid  $\rightarrow$  fromDate

In BCNF (why?)



#### Is BCNF lossless?

#### YES!

Given a non-trivial FD  $X \rightarrow Y$  in  $\mathcal{F}^+$  with  $XY \subseteq R$ , where X is not a super key of R, we need to prove:

Let Z =Attribute (R) - X - Y

- Anything we project always comes back in the join:  $R \subseteq (\pi_{XY}R) \bowtie (\pi_{XZ}R)$ 
  - Sure; and it doesn't depend on the FD
- Anything that comes back in the join must be in the original relation:

 $R \supseteq (\pi_{XY}R) \bowtie (\pi_{XZ}R)$ 

• Proof uses the fact that  $X \to Y$ 

## Is BCNF dependency-preserving?

- NO!
- Consider a simple example R under  $\mathcal{F}$  :



# **Outline for Today**

- Functional Dependencies
- Decomposition:
  - Lossless
  - Dependency-preserving
- Boyce-Codd Normal Form (BCNF)
  - Lossless decomposition
  - Not necessarily dependency-preserving decomposition
- Third Normal Form (3NF)
  - Lossless decomposition
  - Dependency-preserving decomposition

## Third normal form (3NF)

- A relation R is in 3NF under  $\mathcal{F}$  if each FD  $X \rightarrow Y \in \mathcal{F}^+$  with  $XY \subseteq R$  satisfies:
  - either  $X \to Y$  is trivial, i.e.,  $Y \subseteq X$ ,
  - or X is a super key of R, i.e.,  $X \rightarrow R$  or,
  - or each attribute in Y X is in a key of R

BCNF only allows the first two cases, so 3NF is looser than BCNF

- Is  $R = \{A, B, C\}$  under  $\mathcal{F}$  in 3NF?
  - A,  $B \rightarrow C$  is satisfied since AB is a super key
  - $C \rightarrow B$  is satisfied since B is part of key  $\{A, B\}$

 $\mathcal{F}$  includes:  $A, B \rightarrow C$  $C \rightarrow B$ 

#### Compute 3NF decomposition

- Step 1: Finding the minimal cover of the FD set  ${\mathcal F}$ 



- Given a set of FDs  $\mathcal{F}$ , we say  $\mathcal{F}'$  is equivalent to  $\mathcal{F}$  if their closures are the same, i.e.,  $\mathcal{F}^+ = \mathcal{F}'^+$ .
- The smallest equivalent set of FDs
- Step 2: Decompose based on the minimal cover

#### Minimal cover of $\mathcal F$

A set of FDs  $\mathcal F$  is minimal if

- every RHS of a FD in  ${\mathcal F}$  is a single attribute

 Consider R = {sno, sname, city, pno, pname, price, ptype} under a set F of FDs



## Minimal cover of $\mathcal F$

A set of FDs  $\mathcal F$  is minimal if

- every RHS of a FD in  ${\mathcal F}$  is a single attribute
- No FD  $X \to A$  with  $Z \subseteq X$  such that  $(\mathcal{F} - \{X \to A\} + \{Z \to A\})^+ = \mathcal{F}^+$

How to check efficiently? if A is in the attribute closure of Z under  $\mathcal{F}$ 

No redundant attributes in *X* 

• Consider  $R = \{$ sno, sname, city, pno, pname, price, ptype $\}$  under a set  $\mathcal{F}$  of FDs

 $\mathcal{F}$  includes:  $sno \rightarrow sname, city$   $pno \rightarrow pname$   $sno, pno \rightarrow price$   $sno, pname \rightarrow price$  $pno, pname \rightarrow ptype$ 

Fail condition 2: removing pname since pno  $\rightarrow$  pname and pno  $\rightarrow$  ptype together imply it

## Minimal cover of $\mathcal F$

#### A set of FDs $\mathcal{F}$ is minimal if

- every RHS of a FD in  $\mathcal{F}$  is a single attribute
- No FD  $X \to A$  with  $Z \subseteq X$  such that  $(\mathcal{F} - \{X \to A\}) + \{Z \to A\})^+ = \mathcal{F}^+$
- No FD  $X \to A$  such that  $(\mathcal{F} \{X \to A\})^+ = \mathcal{F}^+$



Fail condition 3: pno  $\rightarrow$  pname and sno, pname  $\rightarrow$  price can imply it

 $\mathcal{F}$  includes:

sno  $\rightarrow$  sname, city

How to check efficiently?

if A is in the attribute closure

of X under  $\mathcal{F} - \{X \to A\}$ 

No redundant

FD in  $\mathcal{F}$ 

- pno  $\rightarrow$  pname
- sno, pno  $\rightarrow$  price
  - sno, pname  $\rightarrow$  price
  - pno, pname  $\rightarrow$  ptype

#### Compute minimal cover of $\mathcal F$

Repeat the following steps until  ${\mathcal F}$  does not change

- Step 1: Replace  $X \rightarrow YZ$  with  $X \rightarrow Y$  and  $X \rightarrow Z$
- Step 2: Remove *B* from the LHS of  $X \rightarrow A$  if *A* is in the attribute closure of  $X \{B\}$  under  $\mathcal{F}$
- Step 3: Remove  $X \to A$  if A is in the attribute closure of X under  $\mathcal{F} \{X \to A\}$

Why equivalent set of FDs?

#### Example of minimal cover

• Consider  $R = \{sno, sname, city, pno, pname, price, ptype\}$  under a set  $\mathcal{F}$  of FDs

 $\mathcal{F}$  includes:  $sno \rightarrow sname, city$   $pno \rightarrow pname$   $sno, pno \rightarrow price$   $sno, pname \rightarrow price$  $pno, pname \rightarrow ptype$ 

 $\mathcal{F}^*$  includes:  $sno \rightarrow sname$   $sno \rightarrow city$   $pno \rightarrow pname$   $sno, pname \rightarrow price$  $pno \rightarrow ptype$  Remove it as it is implied by pno  $\rightarrow$  pname and sno, pname  $\rightarrow$  price

Split it into sno  $\rightarrow$  sname and sno  $\rightarrow$  city

Update it by pno  $\rightarrow$  ptype as it is implied by pno  $\rightarrow$  pname and pno  $\rightarrow$  ptype

#### Compute 3NF decomposition

Given a relation R with a set  $\mathcal{F}$  of FDs:

Step 1: Find a minimal cover  $\mathcal{F}^*$  for  $\mathcal{F}$ 

Step 2: For every  $X \to Y$  in  $\mathcal{F}^*$ , add a relation {X, Y} to the decomposition

Step 3: If no relation contains a key for R, add a relation containing an arbitrary key for R to the decomposition

## Example for 3NF decomposition

 Consider R = {sno, sname, city, pno, pname, price, ptype} under a set F of FDs

#### ${\mathcal F}$ includes:

sno  $\rightarrow$  sname, city pno  $\rightarrow$  pname sno, pno  $\rightarrow$  price sno, pname  $\rightarrow$  price pno, pname  $\rightarrow$  ptype

#### $\mathcal{F}^*$ includes:

 $sno \rightarrow sname$   $sno \rightarrow city$   $pno \rightarrow pname$   $sno, pname \rightarrow price$  $pno \rightarrow ptype$ 

$R_1$ (sno, sname) $R_2$ (sno, city)	(Optional) possible optimization: combine R <sub>1</sub> and R <sub>2</sub>	
$R_3$ (pno, pname) $R_4$ (sno, pname, p $R_5$ (pno, ptype)	orice)	
R <sub>6</sub> (sno, pno)	Add th	e key (sno,pno)
		a Lastura 11)

(See Lecture 11) How to determine if K is or is not a key of R?

## Summary (this week)

- Functional dependencies: provide clues towards the elimination of (some) redundancies in a schema
  - Armstrong's axioms
  - Closure of FDs
  - Closure of attributes
- Decomposition (this lecture):
  - BCNF: lossless (but not necessarily dependencypreserving)
  - 3NF: lossless and dependency-preserving (but with more redundancy)

#### Where are we now?

- Overview (Lecture 1)
- Relational data model and Relational Algebra (Lectures 2 - 3)
- SQL (Lectures 4 8)
- Database Design (Lectures 9 12)
- Database Internals
  - Storage & Indexing (Lectures 13 15)
  - Query Processing & Optimization (Lectures 16 18)
  - Transactions (Lectures 19 20)

#### What's next?

• Midterm Review

