

Principles of Database Systems (CS307)

Lecture 10: Normalization: A Deeper Look

Yuxin Ma

Department of Computer Science and Engineering
Southern University of Science and Technology

- Most contents are from slides made by Stéphane Faroult and the authors of Database System Concepts (7th Edition).
- Their original slides have been modified to adapt to the schedule of CS307 at SUSTech.

Prerequisites

Relation Schema and Instance

- A_1, A_2, \dots, A_n are attributes
- $R = (A_1, A_2, \dots, A_n)$ is a **relation schema**
 - Example on the right side:
instructor = (*ID*, *name*, *dept_name*, *salary*)
- $r(R)$ denotes a relation instance r defined over schema R
 - Or to say, the entire table on the right side
- An element t of relation r is called a **tuple**
 - ... and is represented by a row in a table

The relation schema ("R")

A_1	A_2	A_3	A_4
<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000

$r(R)$

A tuple

Relation Schema and Instance

- An analogy to programming languages:
 - Relation - Variables
 - Relation schema – Variable types
 - Relation instance – Value(s) stored in the variable

Database Schema

- Database schema is the **logical structure** of the database
 - It contains a set of relation schemas
 - ... and a set of integrity constraints
- Database instance is a **snapshot** of the data in the database at a given instant in time

Keys

- Let $K \subseteq R$ K 必须得可以区分行与行之间的区别
 - K is a **superkey** of R if values for K are sufficient to identify a unique tuple of each possible relation $r(R)$
 - E.g., $\{ID\}$ and $\{ID, name\}$ are both **superkeys** of instructor
 - If K is a superkey, any superset K' of K where $K' \subseteq R$ is a superkey as well
 - Superkey K is a **candidate key** if K is minimal, i.e., no proper subset of K is a superkey
 - E.g., $\{ID\}$ is a candidate key for *instructor*
- One of the candidate keys is selected to be the **primary key**
 - We mark the primary key with an underline:
instructor = (ID, name, dept_name, salary)

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000

instructor

Normalization

Features of Good Relational Designs

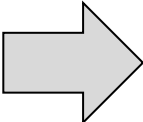
- Suppose we combine *instructor* and *department* into *in_dep*, which represents the natural join on the relations *instructor* and *department*
 - There is repetition of information
 - Need to use nulls (if we add a new department with no instructors)

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000

instructor

<i>dept_name</i>	<i>building</i>	<i>budget</i>
Physics	Watson	70000
Finance	Painter	120000
History	Painter	50000
Comp. Sci.	Taylor	100000
Elec. Eng.	Taylor	85000
Biology	Watson	90000
Comp. Sci.	Taylor	100000
History	Painter	50000
Comp. Sci.	Taylor	100000
Music	Packard	80000
Physics	Watson	70000
Finance	Painter	120000

department



<i>ID</i>	<i>name</i>	<i>salary</i>	<i>dept_name</i>	<i>building</i>	<i>budget</i>
22222	Einstein	95000	Physics	Watson	70000
12121	Wu	90000	Finance	Painter	120000
32343	El Said	60000	History	Painter	50000
45565	Katz	75000	Comp. Sci.	Taylor	100000
98345	Kim	80000	Elec. Eng.	Taylor	85000
76766	Crick	72000	Biology	Watson	90000
10101	Srinivasan	65000	Comp. Sci.	Taylor	100000
58583	Califieri	62000	History	Painter	50000
83821	Brandt	92000	Comp. Sci.	Taylor	100000
15151	Mozart	40000	Music	Packard	80000
33456	Gold	87000	Physics	Watson	70000
76543	Singh	80000	Finance	Painter	120000

in_dep

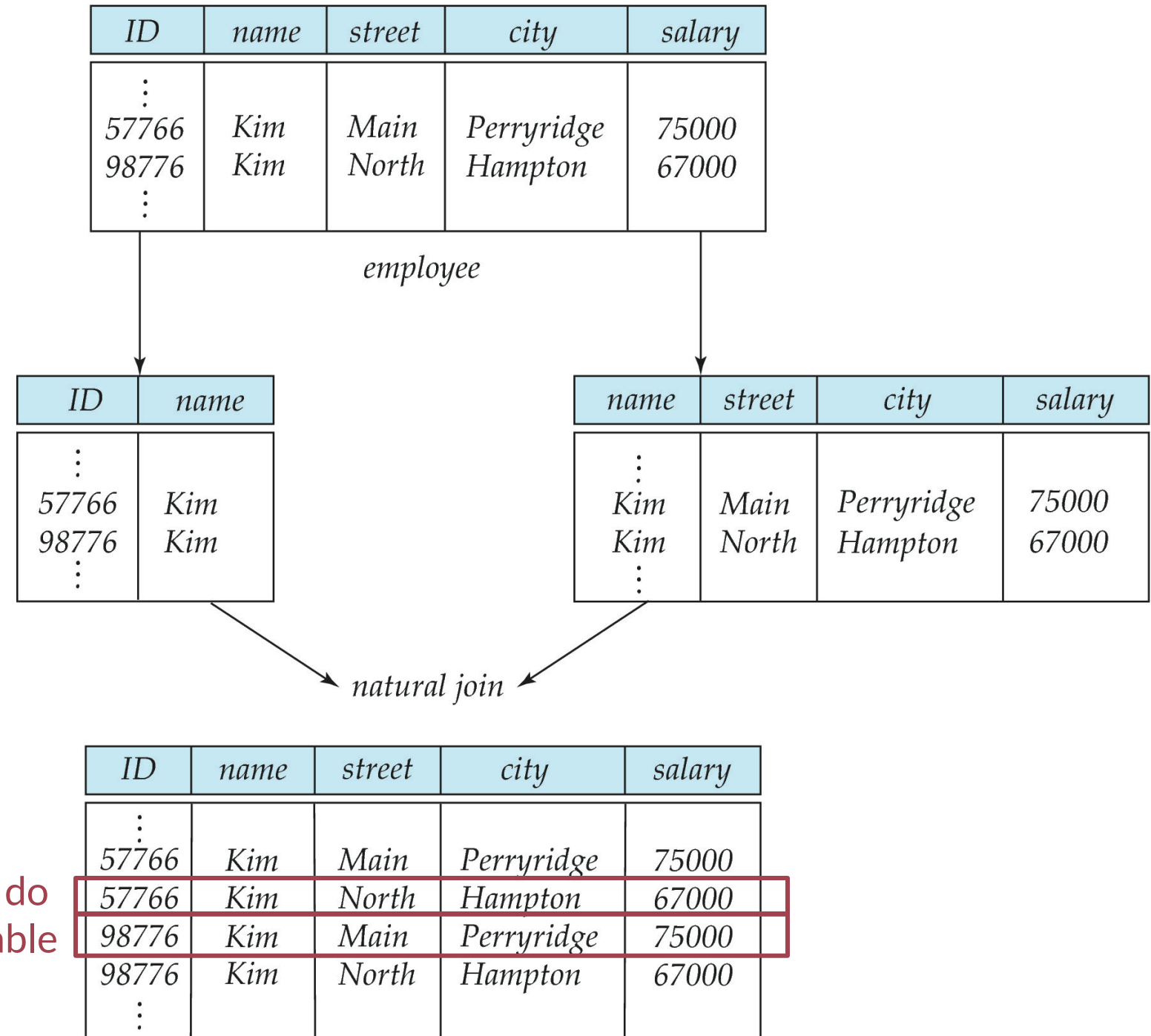
Decomposition

- Avoid the repetition-of-information problem
 - Decompose *in_dep* into two schemas: *instructor* and *department*
- However, not all decompositions are good
 - E.g., decompose *employee*(ID, name, street, city, salary) into:
 - employee1(ID, name)
 - employee2(name, street, city, salary)

The problem arises when we have two employees with the same name

A Lossy Decomposition

- (Continue) we **cannot** reconstruct the original employee relation with the join operation
 - We call it a **lossy decomposition**



Lossless Decomposition

- Let R be a relation schema and let R_1 and R_2 form a decomposition of R
 - That is, $R = R_1 \cup R_2$
 - The decomposition is a **lossless decomposition** if there is no loss of information by replacing R with the two relation schemas $R = R_1 \cup R_2$
- Formally, $\Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) = r$
 - ... and a decomposition is lossy if $r \subset \Pi_{R_1}(r) \bowtie \Pi_{R_2}(r)$
proper subset
- Or to say, the two SQL queries on the right side generate identical results:

```
select * -- 1
from (select R1 from r)
    natural join
    (select R2 from r);

select * from R; -- 2
```

Normalization Theory

- Decide whether a particular relation R is in “good” form
- In the case that a relation R is not in “good” form, decompose it into a set of relations $\{R_1, R_2, \dots, R_n\}$ such that
 - Each relation is in good form
 - The decomposition is a lossless decomposition
- Our theory is based on:
 - Functional dependencies
 - * Multivalued dependencies (self study)

Functional Dependencies

- There are usually a variety of constraints (rules) on the data in the real world
- For example, some of the constraints that are expected to hold in a university database are:
 - **Students** and **instructors** are uniquely identified by their ID
 - Each **student** and **instructor** has only one name
 - Each **instructor** and **student** is (primarily) associated with only one department
 - Each **department** has only one value for its budget, and only one associated building

Functional Dependencies

- An instance of a relation that satisfies all such real-world constraints is called a legal instance of the relation
 - A legal instance of a database is one where all the relation instances are legal instances
- Constraints on the set of legal relations
 - Require that the value for a certain set of attributes **determines uniquely** the value for another set of attributes
- A **functional dependency** is a generalization of the notion of a **key**

Definition of Functional Dependencies

- Let R be a relation schema, and $\alpha \subseteq R$ and $\beta \subseteq R$,
the **functional dependency**

$$\alpha \rightarrow \beta$$

holds on R if and only if for any legal relations $r(R)$, whenever any two tuples t_1 and t_2 of r agree on the attributes α , they also agree on the attributes β . That is,

$$t_1[\alpha] = t_2[\alpha] \Rightarrow t_1[\beta] = t_2[\beta]$$

- Example: Consider $r(A, B)$ with the following instance of r ,
 - On this instance, $A \rightarrow B$ does **NOT** hold, but $B \rightarrow A$ does hold

A	B
1	4
1	5
3	7

Closure of a Set of Functional Dependencies

- Given a set F set of functional dependencies, there are certain other functional dependencies that are logically implied by F :
 - If $A \rightarrow B$ and $B \rightarrow C$, then we can infer that $A \rightarrow C$
- The set of all functional dependencies logically implied by F is the **closure** of F
 - We denote the closure of F by F^+
 - F^+ is a superset of F

Keys and Functional Dependencies

- K is a superkey for relation schema R if and only if $K \rightarrow R$
 - K is a candidate key for R if and only if
 - $K \rightarrow R$, and
 - for **no** $\alpha \subset K$, $\alpha \rightarrow R$
proper subset
 - Functional dependencies allow us to express constraints that cannot be expressed using superkeys
 - Consider the schema: $inst_dept(\underline{ID}, name, salary, \underline{dept_name}, building, budget)$
 - We expect these functional dependencies to hold:
$$dept_name \rightarrow building, ID \rightarrow building$$
- ... but would not expect the following to hold:
- $$dept_name \rightarrow salary$$

Use of Functional Dependencies

- We use functional dependencies to
 - To test relations to see if they are legal under a given set of functional dependencies
 - If a relation r is legal under a set F of functional dependencies, we say that r satisfies F
 - To specify constraints on the set of legal relations
 - We say that F holds on R if all legal relations on R satisfy the set of functional dependencies F

Use of Functional Dependencies

- Example: List some functional dependencies that the table satisfies

A	B	C	D
a1	b1	c1	d1
a1	b2	c1	d2
a2	b2	c2	d2
a2	b3	c2	d3
a3	b3	c2	d4

Use of Functional Dependencies

- Example: List some functional dependencies that the table satisfies
 - $A \rightarrow C$

A	B	C	D
a1	b1	c1	d1
a1	b2	c1	d2
a2	b2	c2	d2
a2	b3	c2	d3
a3	b3	c2	d4

Use of Functional Dependencies

- Note: A specific instance of a relation schema may satisfy a functional dependency even if the functional dependency does not hold on all legal instances.
- Example: we see that $room_number \rightarrow capacity$ is satisfied.
 - However, in real world, two classrooms in different buildings can have the same room number but with different room capacity
 - We prefer $\{building, room_number\} \rightarrow capacity$

<i>building</i>	<i>room_number</i>	<i>capacity</i>
Packard	101	500
Painter	514	10
Taylor	3128	70
Watson	100	30
Watson	120	50

Trivial Functional Dependencies

- A functional dependency is **trivial** if it is satisfied by all relations
- Example:
 - $ID, name \rightarrow ID$
 - $name \rightarrow name$
- In general,
 - $\alpha \rightarrow \beta$ is trivial if $\beta \subseteq \alpha$

Lossless Decomposition

- We can use functional dependencies to show when certain decomposition are **lossless**

- For the case of $R = (R_1, R_2)$, we require that for all possible relations r on schema R

$$r = \Pi_{R_1}(r) \bowtie \Pi_{R_2}(r)$$

- A decomposition of R into R_1 and R_2 is a **lossless decomposition** if at least one of the following dependencies is in F^+ :

- $R_1 \cap R_2 \rightarrow R_1$
- $R_1 \cap R_2 \rightarrow R_2$

In other words, if $R_1 \cap R_2$ forms a **superkey** for either R_1 or R_2 , the decomposition of R is a lossless decomposition

Lossless Decomposition

- Example:
 - *in_dep* (ID, name, salary, dept_name, building, budget)
 - ... and the decomposed schemas, *instructor* and *department*:
 - *instructor*(ID, name, dept_name, salary)
 - *department*(dept_name, building, budget)

$instructor \cap department = dept_name$
 $dept_name \rightarrow dept_name, building, budget$

(... which means the decomposition is lossless)

Lossless Decomposition

- Note: the above functional dependencies are a sufficient condition for lossless join decomposition
 - The dependencies are a necessary condition only if all constraints are functional dependencies

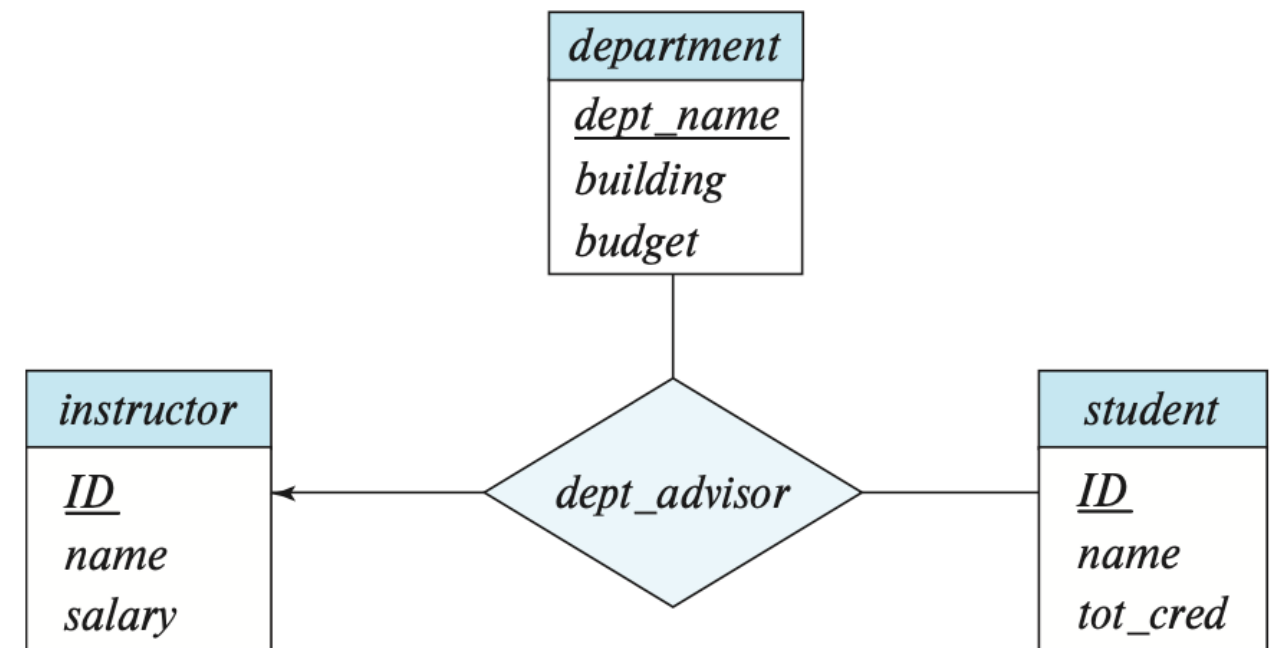
*(There are other types of constraints, e.g., **multivalued dependencies**, that can ensure that a decomposition is lossless even if no functional dependencies are present)*

Dependency Preservation

- Testing functional dependency constraints each time the database is updated can be costly
 - It is useful to design the database in a way that constraints can be tested efficiently.
- If a functional dependency in the original relation R does not exist in any of the decomposed relations, we say it is not **dependency-preserving**
 - In the dependency preservation, at least one decomposed table must satisfy every dependency

Dependency Preservation

- Consider a new E-R design for relationships between students, instructors, and departments
 - An instructor can only be associated with one department
 - A student can have multiple advisors but not more than one from a given department
 - Think about double-major students



Dependency Preservation

- Consider a schema
 - *dept_advisor*(*s_ID*, *i_ID*, *dept_name*)
 - ... with function dependencies: (1) $i_ID \rightarrow dept_name$ (2) $s_ID, dept_name \rightarrow i_ID$

In the above design, we are forced to repeat the department name once for each time an instructor participates in a *dept_advisor* relationship.

Dependency Preservation

- Consider a schema
 - $dept_advisor(\underline{s_ID}, \underline{i_ID}, \underline{dept_name})$
 - ... with function dependencies: (1) $i_ID \rightarrow dept_name$ (2) $s_ID, dept_name \rightarrow i_ID$

In the above design, we are forced to repeat the department name once for each time an instructor participates in a $dept_advisor$ relationship.

- To fix this problem, we need to decompose $dept_advisor$
 - However, any decomposition will not include all the attributes in
 $s_ID, dept_name \rightarrow i_ID$
 - Thus, the decomposition will **NOT** be **dependency-preserving**

Dependency Preservation

- Problem when not meeting dependency preservation
 - Every time the database wants to check the integrity of the functional dependency
$$s_ID, dept_name \rightarrow i_ID,$$
the decomposed tables must be joined
 - ... where the computational cost could be very high with join operations

Normal Forms: Revisited

- Boyce-Codd Normal Form (BCNF)
- 3NF
- Higher-order normal forms

Boyce-Codd Normal Form

- A relation schema R is in **BCNF** with respect to a set F of functional dependencies if for all functional dependencies in F^+ of the form

$$\alpha \rightarrow \beta$$

where $\alpha \subseteq R$ and $\beta \subseteq R$, at least one of the following holds:

- $\alpha \rightarrow \beta$ is **trivial** (i.e., $\beta \subseteq \alpha$)
 - α is a **superkey** for R
-
- * A database design is in BCNF if each member of the set of relation schemas that constitutes the design is in BCNF

Boyce-Codd Normal Form

- Example schema that is **not** in BCNF:
 - *in_dep* (ID, name, salary, dept_name, building, budget)

Because,

$dept_name \rightarrow building, budget$

- holds in *in_dep*, however, *dept_name* is not a superkey
 - ... where {ID, dept_name} is
- When decompose *in_dept* into *instructor* and *department*
 - *instructor* is in BCNF
 - *department* is in BCNF

Decomposing a Schema into BCNF

- Let R be a schema R that is not in BCNF
- Let $\alpha \rightarrow \beta$ be the functional dependency that causes a violation of BCNF
 - We decompose R into:
 - $(\alpha \cup \beta)$
 - $(R - (\beta - \alpha))$
- Example: $in_dep(\underline{ID}, name, salary, \underline{dept_name}, building, budget)$
 - $\alpha = dept_name, \beta = building, budget$
 - Thus, in_dep is replaced by:
 - $(\alpha \cup \beta) = (dept_name, building, budget)$
 - $(R - (\beta - \alpha)) = (ID, name, dept_name, salary)$

BCNF and Dependency Preservation

- It is not always possible to **achieve** both BCNF and dependency preservation
- Consider the schema (that we have visited before)
 - $dept_advisor(\underline{s_ID}, \underline{i_ID}, \underline{dept_name})$
 - ... with function dependencies: (1) $i_ID \rightarrow dept_name$ (2) $s_ID, dept_name \rightarrow i_ID$
 - $dept_advisor$ is not in BCNF since for $i_ID \rightarrow dept_name$, i_ID is not a superkey
 - (where $\{s_ID, i_ID, dept_name\}$ is)
- To fix this problem, we need to decompose $dept_advisor$
 - However, any decomposition will **not** include all the attributes in
 $s_ID, dept_name \rightarrow i_ID$
 - Thus, the decomposition will **NOT** be **dependency-preserving**

Third Normal Form (3NF)

- A relation schema R is in **third normal form (3NF)** if for all

$$\alpha \rightarrow \beta \text{ in } F^+$$

at least one of the following holds:

- $\alpha \rightarrow \beta$ is trivial (i.e., $\beta \subseteq \alpha$)
 - α is a superkey for R
 - Each attribute A in $\beta - \alpha$ is contained in a candidate key for R
- Notes
 - Each attribute A may be in a different candidate key
 - If a relation is in BCNF, it is in 3NF (... since in BCNF, one of the first two conditions above must hold)
 - The third condition above is a minimal relaxation of BCNF to ensure dependency preservation

3NF Example

- Consider the schema (that we have visited before)
 - *dept_advisor*(*s_ID*, *i_ID*, *department_name*)
 - ... with function dependencies: (1) $i_ID \rightarrow dept_name$ (2) $s_ID, dept_name \rightarrow i_ID$
 - We have two candidate keys: $\{s_ID, dept_name\}$ and $\{s_ID, i_ID\}$
- *dept_advisor* is not in BCNF, but it can be in 3NF
 - $\{s_ID, dept_name\}$ is a superkey
 - $i_ID \rightarrow dept_name$ and *i_ID* is **NOT** a superkey (which violates BCNF), but:
 - α is *i_ID*, β is *dept_name*
 - $\{dept_name\} - \{i_ID\} = \{dept_name\}$
 - *dept_name* is contained in a candidate key ($\rightarrow \{s_ID, dept_name\}$)

Redundancy in 3NF

- Consider the schema R below, which is in 3NF
 - $R = (J, K, L)$, $F = \{JK \rightarrow L, L \rightarrow K\}$, and an instance table:
- Problems in this table:
 - Repetition of information
 - Row 1-3: L and K
 - Need to use nulls
 - Row 4: Represent the relationship l_2, k_2 with no corresponding value for J

J	L	K
j_1	l_1	k_1
j_2	l_1	k_1
j_3	l_1	k_1
null	l_2	k_2

Comparison of BCNF and 3NF

- Advantages to 3NF over BCNF
 - It is always possible to **obtain a 3NF design** without sacrificing *losslessness* or *dependency preservation*
- Disadvantages to 3NF
 - We may have to use **nulls** to represent some of the possible meaningful relationships among data items
 - There is a problem of potential repetition of information

Goals of Normalization

- Let R be a relation scheme with a set F of functional dependencies
 - Decide whether a relation scheme R is in “good” form.
 - In the case that a relation scheme R is not in “good” form, need to decompose it into a set of relation scheme $\{R_1, R_2, \dots, R_n\}$ such that:
 - Each relation scheme is in good form
 - The decomposition is a lossless decomposition
 - Preferably, the decomposition should be dependency preserving

How good is BCNF?

- There are database schemas in BCNF that do not seem to be sufficiently normalized
 - Consider a relation *inst_info*(*ID*, *child_name*, *phone*)
 - ... where an instructor may have more than one phone and can have multiple children
 - Actually, we would better use two relations: (*ID*, *child_name*) and (*ID*, *phone*)

An instance of *inst_info*:

(99999, David, 512-555-1234)
(99999, David, 512-555-4321)
(99999, William, 512-555-1234)
(99999, William, 512-555-4321)

How good is BCNF?

- inst_info is in BCNF
 - since $ID \rightarrow child_name$, $ID \rightarrow phone$, and ID is the superkey
- However,
 - Insertion anomalies
 - If we add a phone number 981-992-3443 to the instructor 99999, we need to add two tuples:

(99999, David, 981-992-3443)
(99999, William, 981-992-3443)
 - If we only add one of the two tuples above, it will imply that only David or William corresponds to 981-992-3443, which is not the functional dependency we need to keep

Higher Normal Forms

- It is better to decompose *inst_info* into *inst_child* and *inst_phone*:

<i>ID</i>	<i>child_name</i>
99999	David
99999	William

<i>ID</i>	<i>phone</i>
99999	512-555-1234
99999	512-555-4321

- This suggests a need for higher normal forms, such as Fourth Normal Form (4NF) that resolves such kind of **multivalued dependencies**

Self Study

- Database System Concepts , 7th Edition
 - Chapter 7.4 “Functional Dependency Theory”
 - Chapter 7.5 “Algorithms for Decomposition Using Functional Dependencies”
 - Chapter 7.6 “Decomposition Using Multivalued Dependencies”
 - Chapter 7.7 “More Normal Forms”

Wait, Where are 1NF and 2NF?

- 1NF is about the attribute domains but not decompositions
 - ... and hence not quite related to dependencies we have learned in this section

Wait, Where are 1NF and 2NF?

- 2NF: Partial dependency
 - A functional dependency $\alpha \rightarrow \beta$ is called a partial dependency if there is a proper subset γ of α such that $\gamma \rightarrow \beta$
 - We say that β is partially dependent on α
 - A relation schema R is in second normal form (2NF) if each attribute A in R meets one of the following criteria:
 - It appears in a candidate key
 - It is not partially dependent on a candidate key

Wait, Where are 1NF and 2NF?

- 2NF: Partial dependency
 - You can try to prove that a relation meeting 3NF also satisfies 2NF
 - Exercise 7.19 in “Database System Concepts, 7th Edition”
 - In practice, we usually choose to satisfy 3NF or BCNF

Summary

Overall Database Design Process

- We have assumed schema R is given
 - R could have been generated when converting E-R diagram to a set of tables
 - R could have been a single relation containing **all** attributes that are of interest (called **universal relation**)
 - Normalization breaks R into smaller relations
 - R could have been the result of some ad-hoc design of relations, which we then test/convert to normal form

E-R Model and Normalization

- When an E-R diagram is carefully designed, identifying all entities correctly, the tables generated from the E-R diagram should not need further normalization.
 - However, in a real (imperfect) design, there can be **functional dependencies** from **non-key attributes** of an entity to other attributes of the entity
 - Example: an *employee* entity with attributes *department_name* and *building*
 - ... but with functional dependency: *department_name* → *building*
 - Good design would have made department an entity

Denormalization for Performance

- We may want to use non-normalized schema for performance
- For example, displaying *prereqs* along with *course_id*, and *title* requires join of *course* with *prereq*
 - Alternative 1: Use **denormalized relation** containing attributes of *course* as well as *prereq* with all above attributes
 - faster lookup
 - extra space and extra execution time for updates
 - extra coding work for programmer and possibility of error in extra code
 - Alternative 2: use a materialized view defined a *course* ⋈ *prereq*
 - Benefits and drawbacks same as above, except no extra coding work for programmer and avoids possible errors

Other Design Issues

- Some aspects of database design are not caught by normalization
 - Examples of bad database design, to be avoided: Instead of *earnings* (*company_id*, *year*, *amount*), use
 - *earnings_2004*, *earnings_2005*, *earnings_2006*, etc., all on the schema (*company_id*, *earnings*).
 - Above are in BCNF, but make querying across years difficult and needs new table each year
 - *company_year* (*company_id*, *earnings_2004*, *earnings_2005*, *earnings_2006*)
 - Also in BCNF, but also makes querying across years difficult and requires new attribute each year
 - It is an example of a **crosstab**, where values for one attribute become column names
 - Such crosstabs are widely used in spreadsheets and data analysis tools