

COMP9120

Week 9: Schema Refinement and Normalisation

Semester 1, 2022

Professor Athman Bouguettaya
School of Computer Science



THE UNIVERSITY OF
SYDNEY



Acknowledgement of Country

I would like to acknowledge the Traditional Owners of Australia and recognise their continuing connection to land, water and culture. I am currently on the land of the Darug people and pay my respects to their Elders, past, present and emerging.

I further acknowledge the Traditional Owners of the country on which you are on and pay respects to their Elders, past, present and future.



COMMONWEALTH OF AUSTRALIA

Copyright Regulations 1969

WARNING

This material has been reproduced and communicated to you by or on behalf of the University of Sydney pursuant to Part VB of the Copyright Act 1968 (**the Act**).

The material in this communication may be subject to copyright under the Act. Any further copying or communication of this material by you may be the subject of copyright protection under the Act.

Do not remove this notice.

› Redundancy

- Update/Insertion/deletion anomalies

› Functional Dependencies and Normal Forms

- Functional dependencies
- Attribute closure, candidate keys
- 1NF, 2NF, 3NF, BCNF
- Multivalued dependencies and 4NF

› Schema Decomposition

- How to decompose a relation into BCNF relations
 - Lossless decomposition
 - Dependency-preservation

› Redundancy

- Update/Insertion/deletion anomalies

› Functional Dependencies and Normal Forms

- Functional dependencies
- Attribute closure, candidate keys
- 1NF, 2NF, 3NF, BCNF
- Multivalued dependencies and 4NF

› Schema Decomposition

- Dependency-preservation and Lossless decomposition into BCNF

Redundant Data Can Cause Anomalies

Student	UnitOfStudy	Room
Mary	COMP9120	R101
Joe	COMP9120	R101
Sam	COMP9120	R101
..

If each unit of study is in only one room, this table contains **redundant** information!

Student	UnitOfStudy	Room
Mary	COMP9120	R203
Joe	COMP9120	R101
Sam	COMP9120	R203
..

If we update the room number for all but one tuple, we get inconsistent data = an **update anomaly**

If everyone drops the unit of study, we lose what room this unit is in! = a **delete anomaly**

Similarly, we can't reserve a room for a unit of study without students = an **insert anomaly**



Redundant Data Can Cause Anomalies

Student	UnitOfStudy
Mary	COMP9120
Joe	COMP9120
Sam	COMP9120
..	..

UnitOfStudy	Room
COMP9120	R101
..	..

Is this better?

- Redundancy?
- Update anomaly?
- Delete anomaly?
- Insert anomaly?

- › Example: a bank database called *Lending*:

Table 1: *LENDING*

branchname	assets	city	loan	customer	amount
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>17</i>	<i>Jones</i>	<i>1000</i>
<i>Logan</i>	<i>5000000</i>	<i>SF</i>	<i>23</i>	<i>Smith</i>	<i>2000</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>15</i>	<i>Hayes</i>	<i>1500</i>
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>14</i>	<i>Jackson</i>	<i>1500</i>
<i>King George</i>	<i>10000000</i>	<i>Detroit</i>	<i>93</i>	<i>Curry</i>	<i>500</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>25</i>	<i>Glenn</i>	<i>2500</i>
<i>Andrew</i>	<i>15000000</i>	<i>NY</i>	<i>10</i>	<i>Brooks</i>	<i>2500</i>
<i>Logan</i>	<i>5000000</i>	<i>SF</i>	<i>30</i>	<i>Johnson</i>	<i>750</i>

- › Consider the following scenario: suppose we want to insert one more entry reflecting that Harris is requesting a loan (the id is 40) of 2100 dollars from the Mall St branch:

(Mall St, 40, Harris, 2100)

- › This means that we would have to insert the following tuple to *Lending*:

(Mall St, 9000000, LA, 40, Harris, 2100)

- › But then the branch information is repeated *many times* in the relation *Lending*!

› What are the problems?

- Information is *needlessly repeated*, causing:
 - Waste of space
 - *branchname* and *city* are repeated for each new *loan* in the relation *Lending*, thus wasting disk space.
- Potential inconsistencies due to updates, deletions and insertions
 - If *branchname* changes, there is a need to update more than one tuple. Possible inconsistencies.
 - Assume *Lending* is the *only* relation in the database so far: if a new branch is inserted and no loans have been awarded yet, how are we to record this information?
 - May be add null values but are usually hard to handle because of their ambiguous semantics.
 - If all loans have been paid out, information about the branch is also deleted!

- › Bad design usually means
 - Redundancy (waste of disk space)
 - Update/insertion/deletion anomalies
- › Good design usually means
 - Minimal redundancy
 - No update/insertion/deletion anomalies

Develop a theory to understand ***whether*** a relation needs to be ***decomposed*** and if yes, ***how*** to decompose it

› Redundancy

- Update/Insertion/deletion anomalies

› Functional Dependencies and Normal Forms

- **Functional dependencies**
- Attribute closure, candidate keys
- 1NF, 2NF, 3NF, BCNF
- Multivalued dependencies and 4NF

› Schema Decomposition

- Dependency-preservation and Lossless decomposition into BCNF

- › Solution to bad design: *Normalize* the design to avoid these pitfalls

Normalization is the process that defines what is an *acceptable as good relational design*.

Note that normalization is:

- *Geared towards update*

Before we talk about the normalization process, we need to define some needed tools :

- **Functional Dependencies:** What are they?

They are a tool:

- to capture a *semantic relationship* between attributes
- used to *eliminate bad* design problems

In database design, functional dependencies are used to *break up relations* to eliminate update/insertion/deletion anomalies and some redundancy.

- › **Dependency:** value of attribute X **determines** the value of attribute Y
 - Every UoS is taught in a single room
 - COMP9120 is always held in Room R101
- › **Functional Dependency (FD):**
 - We write $X \rightarrow Y$, and say “X (functionally) determines Y”
 - *Definition of a functional dependency:* Assuming that X and Y are two sets of attributes, the relationship between X and Y values is modelled using a *function*. This essentially means that *a value of X cannot be mapped to more than one value of Y*. Only (n-1) (and thus (1-1)) relationships may exist between X and Y.
 - The functional dependency is denoted, $X \rightarrow Y$

$$\text{UoS} \rightarrow \text{Room}$$
- › FDs are used to identify *sources of redundancy* within schemas.



- › Using the example of the bank database

Table 1: *LENDING*

branchname	assets	city	loan	customer	amount
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>17</i>	<i>Jones</i>	<i>1000</i>
<i>Logan</i>	<i>5000000</i>	<i>SF</i>	<i>23</i>	<i>Smith</i>	<i>2000</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>15</i>	<i>Hayes</i>	<i>1500</i>
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>14</i>	<i>Jackson</i>	<i>1500</i>
<i>King George</i>	<i>10000000</i>	<i>Detroit</i>	<i>93</i>	<i>Curry</i>	<i>500</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>25</i>	<i>Glenn</i>	<i>2500</i>
<i>Andrew</i>	<i>15000000</i>	<i>NY</i>	<i>10</i>	<i>Brooks</i>	<i>2500</i>
<i>Logan</i>	<i>5000000</i>	<i>SF</i>	<i>30</i>	<i>Johnson</i>	<i>750</i>

- › How do we determine functional dependencies? By looking either at the:
 - *Meaning* of the attributesor
 - *Data*
- › In most cases, it is the *former* (i.e., meaning of attributes) that we use. In the *latter* case (looking at the data), the process is called *knowledge mining*.
- › Example of using *data*:
 - In relation *Lending* table, there is a functional dependency:
$$\text{branchname} \rightarrow \text{city}$$



Functional Dependencies

Table 1: *LENDING*

branchname	assets	city	loan	customer	amount
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>17</i>	<i>Jones</i>	<i>1000</i>
<i>Logan</i>	<i>5000000</i>	<i>SF</i>	<i>23</i>	<i>Smith</i>	<i>2000</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>15</i>	<i>Hayes</i>	<i>1500</i>
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>14</i>	<i>Jackson</i>	<i>1500</i>
<i>King George</i>	<i>10000000</i>	<i>Detroit</i>	<i>93</i>	<i>Curry</i>	<i>500</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>25</i>	<i>Glenn</i>	<i>2500</i>

Table 1: *LENDING*

branchname	assets	city	loan	customer	amount
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>17</i>	<i>Jones</i>	<i>1000</i>
<i>Logan</i>	<i>5000000</i>	<i>SF</i>	<i>23</i>	<i>Smith</i>	<i>2000</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>15</i>	<i>Hayes</i>	<i>1500</i>
<i>Mall St</i>	<i>9000000</i>	<i>LA</i>	<i>14</i>	<i>Jackson</i>	<i>1500</i>
<i>King George</i>	<i>10000000</i>	<i>Detroit</i>	<i>93</i>	<i>Curry</i>	<i>500</i>
<i>Queen</i>	<i>500000</i>	<i>DC</i>	<i>25</i>	<i>Glenn</i>	<i>2500</i>
<i>Andrew</i>	<i>15000000</i>	<i>NY</i>	<i>10</i>	<i>Brooks</i>	<i>2500</i>
<i>Logan</i>	<i>5000000</i>	<i>SF</i>	<i>30</i>	<i>Johnson</i>	<i>750</i>

- › We might be able to tell that a certain FD does *not* hold by looking at an instance of a relation,
- › However, we can **never** deduce that an FD *does* hold by looking at one or more instances of the relation, because an FD is a statement about *all* possible legal instances of the relation.

› **Armstrong's Axioms** (A, B, C are sets of attributes):

1. **Reflexivity:** If $B \subseteq A$, then $A \rightarrow B$

- Example: $cpoints, uos_name \rightarrow uos_name$

2. **Augmentation:** If $A \rightarrow B$, then $AC \rightarrow BC$ for any C

- Example: $cpoints \rightarrow wload$ **implies** $cpoints, uos_name \rightarrow wload, uos_name$

3. **Transitivity:** If $A \rightarrow B$ and $B \rightarrow C$, then $A \rightarrow C$

- Example: $uos_code \rightarrow cpoints, cpoints \rightarrow wload$ **implies** $uos_code \rightarrow wload$

› **Trivial FDs:** When all RHS attributes appear on LHS: *all reflexive* FDs are trivial.

› Example

Products

Name	Color	Category	Dept	Price
Gizmo	Green	Gadget	Toys	49
Widget	Black	Gadget	Toys	59
Gizmo	Green	Whatsit	Garden	99

Provided FDs:

1. Name \rightarrow Color
2. Category \rightarrow Department
3. Color, Category \rightarrow Price

Given the provided FDs, we can see that *Name, Category \rightarrow Price* must also hold on **any instance** (use augmentation and transitivity)...

Name \rightarrow Color
Color, Category \rightarrow Price

- › Let F be a set of functional dependencies:
- › We say that F *logically implies* the FD $X \rightarrow Y$, (denoted $F \models X \rightarrow Y$) if for all functional dependencies in F are satisfied by a relation R , this implies that the FD $X \rightarrow Y$ is also satisfied by the relation R .

Example:

- › If the set $F = \{ X \rightarrow Y \text{ and } Y \rightarrow Z \}$ then $F \models X \rightarrow Z$ (using the transitivity rule)
- › **Basis:** A set of FDs is called a *basis* for a relation R if we can infer from this set all FDs for the relation R .
- › **Minimality:** If there is no subset of a *basis* from which we can derive the set of all FDs for a relation R , then the *basis* is *minimal*.

- › *Closure* of a set F is the set F^+ , the set of *all* FDs that *can be deduced* from the set F using the inference rules (Armstrong's axioms): reflexivity, augmentation, and transitivity inference rules.
- › All FDs logically implied by an initial set of given FDs F , is known as the **closure of F , labelled F^+**

- › The closure of a set F is the set F^+ defined by:

$$F^+ = \{X \rightarrow Y \mid F \models X \rightarrow Y\}$$

- › In general, the calculation of the closure of F is, in the worst case, exponential in the number of attributes in F .

› Armstrong's Axioms are

- *Sound*
 - they generate *only* FDs in F^+ when applied to a set F of FDs
- *Complete*
 - repeated application of these rules will generate *all* FDs in the closure F^+

› A few additional rules (that follow from Armstrong's Axioms.):

- **Union**

If $A \rightarrow B$ and $A \rightarrow C$, then $A \rightarrow BC$

Proof:

$A \rightarrow B$

$A \rightarrow C$

$AA \rightarrow AB$ [Augmentation]

$AB \rightarrow BC$

$A \rightarrow AB$

$A \rightarrow BC$ [Transitivity]

- **Decomposition**

If $A \rightarrow BC$, then $A \rightarrow B$ and $A \rightarrow C$

› Redundancy

- Update/Insertion/deletion anomalies

› Functional Dependencies and Normal Forms

- Functional dependencies
- **Attribute closure, candidate keys**
- 1NF, 2NF, 3NF, BCNF
- Multivalued dependencies and 4NF

› Schema Decomposition

- Dependency-preservation and Lossless decomposition into BCNF

- › A **superkey** is a set of attributes that uniquely identify each tuple in a relation
 - If K is a superkey for a relation R , $K \rightarrow R$ i.e. all the attributes of R
- › A **candidate key** (or just key) is the minimal superkey
 - No subset of a candidate key is a candidate key itself
 - There may be many candidate keys for a relation, but *only* 1 primary key.
- › Given a relation R , with attributes **ABCDE** (each letter denotes an attribute) where:
 - **A** uniquely identifies each row in R
 - **BC** also uniquely identifies each row in R (but not B or C alone)
- › **A** is a superkey (and candidate key) for R
- › **BC** is a superkey (and candidate key) for R
- › **BCE** is a superkey (but not a candidate key)
 -as it is *not* minimal

› Attribute Closure of X (X^+)

- › Given a set F of FDs and a set X of attributes, X^+ is the set of all attributes that are determined by X under F .

$$X \rightarrow X^+$$

- › If we know the functional dependencies, then we can check whether an attribute (or a set of attributes) is a key for the relation
 1. Any set of attributes, whose **attribute closure** is the whole relation, is a **superkey**
 2. A superkey is a candidate key if it is **minimal** (i.e., none of its subset is a superkey)

Starting with the given set X of attributes, one repeatedly expands the set by adding the right side of an FD as soon as the left side is present:

Algorithm:

1. Initialise *result* with the given set of attributes: $X = \{A_1, \dots, A_n\}$ (*reflexivity rule*)
2. Repeatedly search for some FD $A_1 A_2 \dots A_m \rightarrow C$ such that all A_1, \dots, A_m are already in the set of attributes *result*, but C is not.
3. Add C to the set *result*. (*transitivity and decomposition rules*)
4. Repeat step 2 until no more attributes can be added to *result*
5. The set *result* is the correct value of X^+

Example: Test for Candidate keys

$R(A, B, C, G, H, I)$

$F = \{$

$A \rightarrow B$

$A \rightarrow C$

$CG \rightarrow H$

$CG \rightarrow I$

$B \rightarrow H$

$\}$

› Check $(AG)^+$

1. $result = AG$
2. $result = ABG (A \rightarrow B)$
3. $result = ABCG (A \rightarrow C)$
4. $result = ABCGH (CG \rightarrow H)$
5. $result = ABCGHI (CG \rightarrow I)$

› Is (AG) a candidate key?

1. Is (AG) a super key? YES! $(AG)^+ = R$
2. Is (AG) minimal? (Is any subset of (AG) a superkey?)
 - $(A)^+ = ABCH \neq R$
 - $(G)^+ = G \neq R$

› Yes, (AG) is a candidate key

What are the keys of our relation?

PhD(sid, first, last, dept, advisor, award, description)

Given FDs

- a) sid \rightarrow first, last
- b) advisor \rightarrow dept
- c) description \rightarrow award
- d) sid,dept \rightarrow advisor,description

2 keys:

1. (sid, dept)
2. (sid, advisor)

Short 5 mn break:

please stand up, stretch, and move around



THE UNIVERSITY OF
SYDNEY

› Redundancy

- Update/Insertion/deletion anomalies

› Functional Dependencies and Normal Forms

- Functional dependencies
- Attribute closure, candidate keys
- **1NF, 2NF, 3NF, BCNF**
- Multivalued dependencies and 4NF

› Schema Decomposition

- Dependency-preservation and Lossless decomposition into BCNF

- › Database theory identifies several ***normal forms*** for relational schemas.
 - Normal Forms
 - Goal: reduce redundancy (hence potential anomalies)
 - In what follows, a key refers to a **candidate key**, unless otherwise stated. Without loss of generality, we assume that we have one single candidate key in all examples, unless specified otherwise.
 - Each normal form is characterized by a set of restrictions.
 - For a relation to be in a normal form, it must satisfy the restrictions associated with that form.
- › We focus on 1NF, 2NF, 3NF, BCNF, and 4NF
 - Restrictions based on *functional dependencies* and *multivalued dependencies*.
- › Normalisation is usually used to verify a design
 - Start a design with a set of relations and use normalisation to break them down into smaller relations

- › A relation R is in **first normal form (1NF)** if the domains of all attributes of R are *atomic*.
- *multivalued attributes* are an example of *non-atomic* domains

Student	UnitOfStudy
Mary	{COMP9120,COMP5318}
Joe	{COMP9120,COMP5313}

Violates 1NF

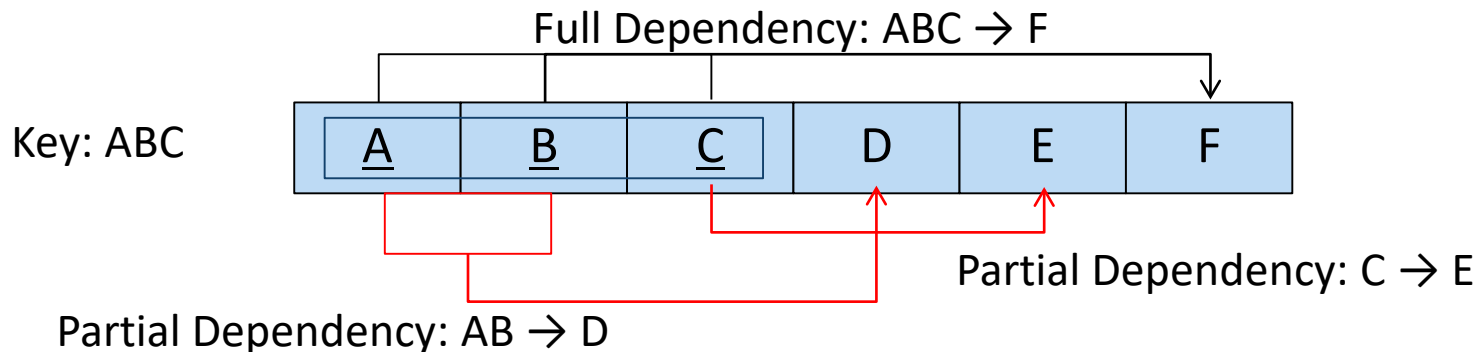
Student	UnitOfStudy
Mary	COMP9120
Mary	COMP5318
Joe	COMP9120
Joe	COMP5313

In 1NF

Second Normal Form (2NF)

> 1NF + No partial dependencies

- A partial dependency is a *non-trivial* FD $X \rightarrow Y$ for R where X is a strict (proper) subset of some key for R and Y is not part of a key. In other words, *There cannot be a functional dependency between a subset of a key to non-key attributes*. This is applicable only when the key consists of more than one attribute.



<u>Teacher_id</u>	<u>UnitOfStudy</u>	Teacher_position
Mary	COMP9120	Lecturer
Mary	COMP5313	Lecturer

Violates 2NF

Second Normal Form (2NF)

- › Problem: redundancy.

In the example above, the teacher position would have to be repeated for all units of study that teacher teaches!

- › More formally,

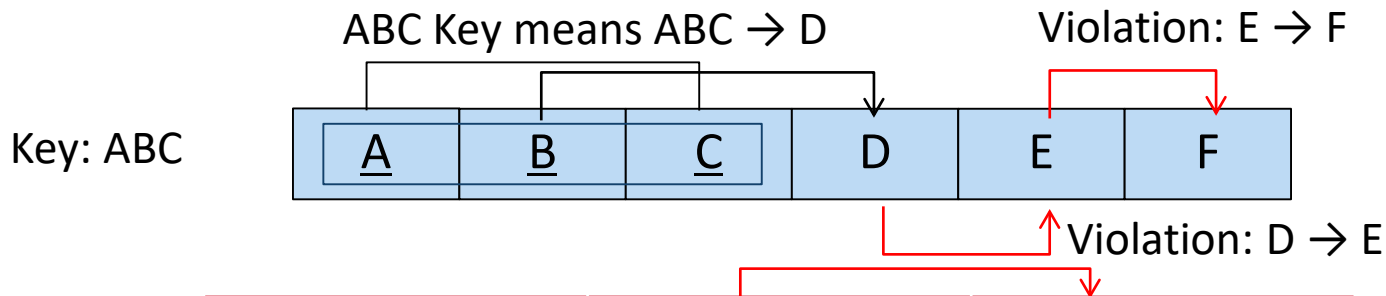
Definition: a relation is in the 2NF if the closure F^+ contains *no* functional dependency of the form: $X \rightarrow Y$ where Y is nonprime (not part of a candidate key) and X is a proper subset of a candidate key. In this case, Y is said to be *fully functionally dependent* on the key.

- › What is the solution?

- Split up the relation in two relations such that X and Y form one relation and X is in the remaining relation. Check that the resulting relations are in 2NF. The relation (X,Y) is now in 2NF. Check the remaining relation if it is in 2NF.
-

› Third Normal Form (3NF)

- › Formal Definition: a relation R is in 3NF if for each dependency $X \rightarrow Y$ in F^+ , *at least one of the following holds*:
 - $X \rightarrow Y$ is a trivial FD ($Y \subseteq X$)
 - X is a superkey for R
 - $Y \subset$ (is a proper subset of) a candidate key for R
 - › Another definition of 3NF is : a relation is in 3NF iff it is in 2NF and every non-key attribute is *not* transitively functionally dependent on a key.
-



<u>Employee</u>	Department	Location
Jim	Research	Brisbane
John	Manufacturing	Melbourne
Mary	Research	Brisbane
Sue	Manufacturing	Melbourne

Violates 3NF

- › Note that there is a functional dependency between *Department* and *Location* (thus transitive dependency).

› Problem: redundancy

The department's location is repeated with every employee working in that department.

› Solution: split up the relation into two relations:

<u>Employee</u>	Department
-----------------	------------

<u>Department</u>	Location
-------------------	----------

- › Boyce-Codd Normal Form (BCNF) is a stronger form of 3NF

Problem: 3NF allows functionally dependencies between *non-prime attributes to prime attributes*.

- › A relation R is in **BCNF** if **all** non-trivial FDs that hold over R **have** a superkey of R on the LHS. In another words, a relation is in BCNF iff every attribute *is dependent on the key, the whole key and nothing but the key*.
 - Formal: For all *non-trivial* $X \rightarrow A$ for R : X is a superkey for R
 - Informal: All dependency arrows come out of candidate keys

Diagram illustrating a table with three columns: Teacher_id, UnitOfStudy, and Address. The table contains two rows of data:

<u>Teacher_id</u>	<u>UnitOfStudy</u>	<u>Address</u>
Mary	COMP9120	One Street
Mary	COMP5313	One Street

Key: AB

Allowed: $AB \rightarrow C, D$

Violation: $C \rightarrow B$

Violation: $\text{Address} \rightarrow \text{Teacher-Id}$

› Redundancy

- Update/Insertion/deletion anomalies

› Functional Dependencies and Normal Forms

- Functional dependencies
- Attribute closure, candidate keys
- 1NF, 2NF, 3NF, BCNF
- Multivalued dependencies and 4NF

› Schema Decomposition

- Dependency-preservation and Lossless decomposition into BCNF

MultiValued Dependency (MVD)

- Because the 1NF does not allow more than one value for each attribute, and in case there are several of such attributes *that are independent* from each other, there is a need to repeat all combinations of these attributes *so not to infer* a relationship.
- More formally, let $R(X, Y, Z)$ be a relation and X, Y and Z are sets of attributes of R , there is a *multivalued dependency*, noted:

$$X \twoheadrightarrow Y$$

- if for all tuples t_1 and t_2 that agree in X , i.e., $t_1[X] = t_2[X]$, there exist tuples t_3 and t_4 such that:

1. $t_1[X] = t_2[X] = t_3[X] = t_4[X]$ and
2. $t_3[Y] = t_1[Y]$ and $t_4[Y] = t_2[Y]$ and
3. $t_3[Z] = t_2[Z]$ and $t_4[Z] = t_1[Z]$

tuples	X	Y	Z
t_1	$a_1 \dots a_i$	$a_{i+1} \dots a_j$	$a_{j+1} \dots a_n$
t_2	$a_1 \dots a_i$	$b_{i+1} \dots b_j$	$b_{j+1} \dots b_n$
t_3	$a_1 \dots a_i$	$a_{i+1} \dots a_j$	$b_{j+1} \dots b_n$
t_4	$a_1 \dots a_i$	$b_{i+1} \dots b_j$	$a_{j+1} \dots a_n$

- In simple English: if there is a multivalued dependency between X and Y , then this means that *no relationship can be inferred* between Y and Z , i.e., they are *independent* of each other.

MultiValued Dependency (MVD)

Example

Name \twoheadrightarrow *Profession*? No, this is not an MVD

The values suggest that whenever

- John is electrician he speaks French.
- John is plumber he speaks Swahili.

NAME	PROFESSION	LANGUAGE
JOHN	ELECTRICIAN	FRENCH
MARY	DOCTOR	SPANISH
JOHN	PLUMBER	SWAHILI
MARY	AUTHOR	CHINESE

- This of course does not make much sense. For instance, John speaks French and Swahili regardless of his professions. Therefore, there should be a multivalued dependency:

we add a few tuples to reflect this information.

NAME	PROFESSION	LANGUAGE
JOHN	ELECTRICIAN	FRENCH
MARY	DOCTOR	SPANISH
JOHN	PLUMBER	SWAHILI
MARY	AUTHOR	CHINESE
JOHN	PLUMBER	FRENCH
MARY	AUTHOR	SPANISH
JOHN	ELECTRICIAN	SWAHILI
MARY	DOCTOR	CHINESE

Fourth Normal Form (4NF)

Redundancy problem in MVDs: should list all professions for every language a person speaks.

Deals with *multivalued dependencies*.

Formally:

R is in 4NF if for all MVDs of the form $X \twoheadrightarrow Y$ in F^+ , at least one of the following holds:

- $X \twoheadrightarrow Y$ is a trivial MVD (i.e., either $Y \subseteq X$ or $X \cup Y = R$)
- X is a superkey for R

Fourth Normal Form (4NF)

Assuming the only key to the following relation is the set
(Employee, Project-id, Personal-phone#):

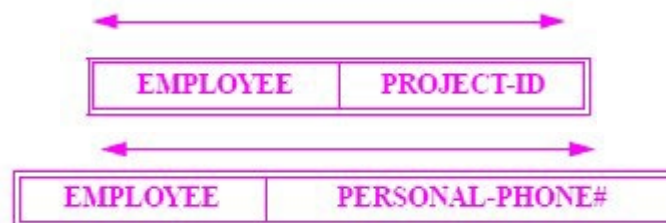
Employee	Project-id	Personal-phone#
Bob	P1	202-222-2222
Bob	P3	703-333-3333
Bob	P1	703-333-3333
Bob	P3	202-222-2222
Janet	P1	421-534-4452
Nguyen	P7	703-666-6666

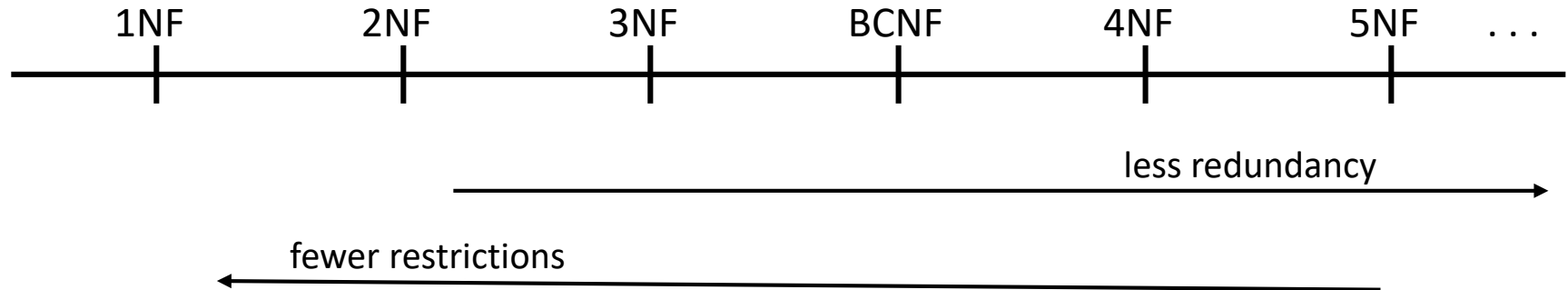
Is this relation in 4NF?

No: There is at least *one non-trivial multivalued dependency*

Employee \twoheadrightarrow Project-id, and Employee is *not* a superkey.

Solution: split the above relation into two relations:





› Higher normal forms

- 5NF, 6NF/DKNF
- They also exploit other types of dependencies
 - Join dependencies
 - Inclusion dependencies

› Redundancy

- Update/Insertion/deletion anomalies

› Functional Dependencies and Normal Forms

- Functional dependencies
- Attribute closure, candidate keys
- 1NF, 2NF, 3NF, BCNF
- Multivalued dependencies and 4NF

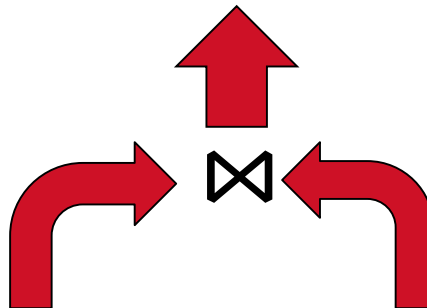
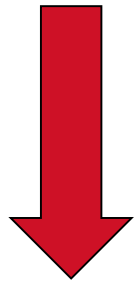
› Schema Decomposition

- Dependency-preservation and Lossless decomposition into BCNF

Decomposition Example 1

$uosCode \rightarrow uosName$

Student	UoSCode	UoSName	Grade
Alice	COMP5138	Database Management Systems	CR
Alice	COMP5338	Advanced Data Models	D
Bob	COMP5138	Database Management Systems	P
Clare	COMP5338	Advanced Data Models	HD
David	COMP5338	Advanced Data Models	CR



Student	UoSCode	Grade
Alice	COMP5138	CR
Alice	COMP5338	D
Bob	COMP5138	P
Clare	COMP5338	HD
David	COMP5338	CR

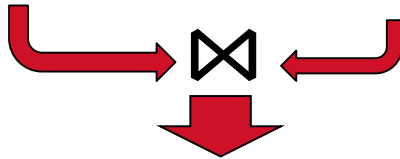
UoSCode	UoSName
COMP5138	Database Management Systems
COMP5338	Advanced Data Models



Decomposition Example 2

uosCode → uosName

Student	UoSCode	UoSCode	UoSName	Grade
Alice	COMP5138	COMP5138	Database Management Systems	CR
Alice	COMP5338	COMP5138	Database Management Systems	P
Bob	COMP5138	COMP5338	Advanced Data Models	CR
Clare	COMP5338	COMP5338	Advanced Data Models	D
David	COMP5338	COMP5338	Advanced Data Models	HD



Student	UoSCode	UoSName	Grade
Alice	COMP5138	Database Management Systems	CR
Alice	COMP5138	Database Management Systems	P
Alice	COMP5338	Advanced Data Models	CR
Alice	COMP5338	Advanced Data Models	D
Alice	COMP5338	Advanced Data Models	HD
Bob	COMP5138	Database Management Systems	CR
Bob	COMP5138	Database Management Systems	P
Clare	COMP5338	Advanced Data Models	CR
Clare	COMP5338	Advanced Data Models	D
Clare	COMP5338	Advanced Data Models	HD
David	COMP5338	Advanced Data Models	CR
David	COMP5338	Advanced Data Models	D
David	COMP5338	Advanced Data Models	HD

- › A decomposition of R replaces R by two or more distinct relations:
 - Each new relation schema contains a subset of the attributes of R (and no new attributes that do not appear in R), and
 - Every attribute of R appears as an attribute of at least one of the new relations.
 - Many possible decompositions – however, not all equally good
 - **Lossless-join:** Re-joining a decomposition of R should give us back R!
 - **Dependency preserving:** No FDs are lost in the decomposition

› Dependency preservation

When decomposing a relation, we may require that dependencies be preserved in the resulting component relations because dependencies are used to express the constraints on the database. If they are not preserved, a solution is to perform costly joins to check whether the dependencies still hold.

› Assuming R is decomposed into $R_1 R_2 \dots R_{(n-1)} R_n$

Idea: First check whether the "lost" dependency can be deduced across relations. If it cannot be deduced, we have no other choice but to perform joins to check whether a functional dependency still holds every time we have an update.

› How do we check?

Let $F' = F_1 \cup F_2 \cup \dots \cup F_{(n-1)} \cup F_n$ be the union of sets of FDs of the decomposed relations. F_i is that subset of F^+ that is applicable (i.e., projection) to R_i .

If $F' \neq F$, run a simple algorithm that checks whether $F'^+ = F^+$.

› Lossless join decomposition

The loss here refers *to the loss of information* and *not* to the loss of tuples!

When joining back the resulting relations, *we may be losing information by adding extra tuples as we saw before!*

- › A decomposition is not always lossless. We may have more than one choice to split up a relation in the course of normalization.

For example, if the decomposition is on a non-key attribute, the decomposition may be lossy.

To ensure lossless join decomposition use *Functional Dependencies*.

- › If R is a relation decomposed into relations $R_1, R_2, \dots, R_{n-1}, R_n$ and F is a set of functional dependencies. This decomposition is lossless (with respect to F) if for every relation R that satisfies F , the following equation is true:
- › $R = \Pi R_1(R) \bowtie \Pi R_2(R) \dots \bowtie \Pi R_{n-1}(R) \bowtie \Pi R_n(R)$ where Π and \bowtie are the symbols for *projection* and *natural join*, respectively.
- › More formally: Let R be a relation and R_1 and R_2 be a decomposition of R . Let F be the set of functional dependencies. This decomposition is a lossless-join decomposition if at least one of the following FDs are in F^+ .
 - $R_1 \cap R_2 \rightarrow R_1$
 - $R_1 \cap R_2 \rightarrow R_2$
- › In plain English, this means that if the intersection of the set of attributes between R_1 and R_2 functionally determines either R_1 or R_2 (i.e., the intersection is a key to either of the resulting relations), then the composition is lossless.

› BCNF Lossless decomposition algorithm

result = {R}

done = false

compute F^+

while *not done* do

 if there is a relation R_i in result that is not in BCNF then

 {

 let $X \rightarrow Y$ be a nontrivial FD that holds on R_i such that $X \rightarrow R_i$ is not in F^+ and
 $X \cap Y = \emptyset$

 result = (result - R_i) \cup ($R_i - Y$) \cup (X, Y)

 }

 else

 done = true

If we did not require $X \cap Y = \emptyset$, then those attributes in $X \cap Y$ would not appear in the schema ($R_i - Y$) and the dependency $X \rightarrow Y$ (across the two new relations) would no longer hold.

Decomposing a Schema into BCNF

- › Suppose we have a schema R and a non-trivial dependency $X \rightarrow Y$ causes a violation of BCNF. We decompose R into:

$$R_1 = X \cup Y$$

$$R_2 = R - Y$$

- › Example schema that is *not* in BCNF:

$loan_info = (\underline{customer_id}, \underline{loan_number}, amount)$ with $loan_number \rightarrow amount$

but $loan_number$ is not a superkey

- › Here,

- $X = loan_number$
- $Y = amount$

So, $loan_info$ is replaced by

$$(X \cup Y) = (\underline{loan_number}, amount)$$

$$(R - Y) = (\underline{customer_id}, \underline{loan_number})$$

- › Is the following relation in BCNF? If not, give a lossless-join decomposition.

employee (*emp_ID*, *name*, *salary*, *dept_name*, *building_number*) with
dept_name → *building_number*

- › No, because *dept_name* is not a superkey
- › Here,
- $X = dept_name$
 - $Y = building_number$

So, *employee* is replaced by

$(X \cup Y) = (dept_name, building_number)$

$(R - Y) = (emp_ID, name, salary, dept_name)$

Is the following relation in BCNF? If not, give a lossless-join decomposition.

contracts (contractID, supplierID, projectID, deptID, itemID, quantity, value)

Functional dependencies:

contractID \rightarrow *supplierID, projectID, deptID, itemID, quantity, value*

supplierID, deptID \rightarrow *itemID*

projectID \rightarrow *supplierID*

contractID, supplierID, projectID, deptID, itemID, quantity, value

supplierID, deptID \rightarrow *itemID*

supplierID, deptID, itemID

supplierID, deptID, contractID, projectID, quantity, value

projectID \rightarrow *supplierID*

projectID, supplierID

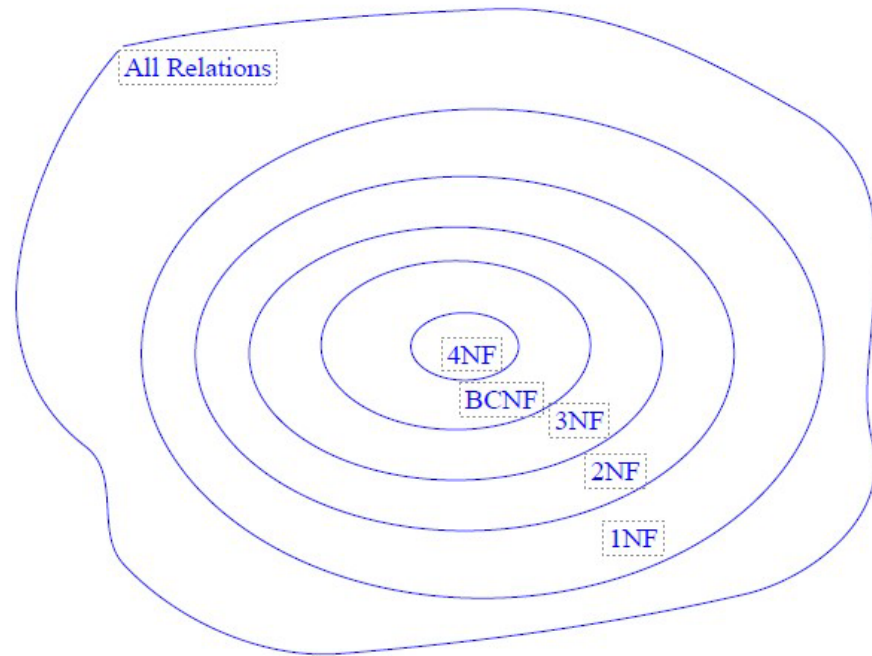
projectID, deptID, contractID, quantity, value



› Summary: Properties of Normal Forms and their decompositions

Property	1NF	2NF	3NF	BCNF	4NF
Eliminates redundancy due to FDs	<i>NO</i>	<i>Some</i>	<i>Most</i>	<i>All</i>	<i>All</i>
Eliminates redundancy due to MVDs	<i>NO</i>	<i>NO</i>	<i>NO</i>	<i>NO</i>	<i>All</i>
Preserves FDs	<i>N/A</i>	<i>Yes</i>	<i>Yes</i>	<i>Maybe</i>	<i>Maybe</i>

Relationships between Normal forms:



Remember that normalization is

- geared towards update
- not necessarily good for retrieval

This means that it is not always good to normalize up to the 4NF/5NF/6NF normal form (e.g., archives, historical databases, etc).

You should be able to perform the following tasks

- › Identify and interpret functional dependencies for a database schema
- › Identify the candidate keys and normal form (up to 4NF) of a relation schema using its functional and multivalued dependencies
 - Identify whether a set of attributes forms a minimal superkey
 - Determine all possible candidate keys
- › Decompose a relational instance into a set of BCNF relational instances
 - Determine a lossless join decomposition of a relational schema
 - Correctly determine the decomposed relation instances

- › Storage and Indexing
 - Storing data in a database
 - Retrieving records from a database
 - B⁺ Tree index
- › Kifer/Bernstein/Lewis
 - Chapter 9 (9.1-9.5)
- › Ramakrishnan/Gehrke
 - Chapter 8
- › Ullman/Widom
 - Chapter 8 (8.3 onwards)
- › Silberschatz/Korth/Sudarshan (5th ed)
 - Chapter 11 and 12

See you next week!



THE UNIVERSITY OF
SYDNEY