

CS 4440 A

# Emerging Database Technologies

---

Lecture 3  
01/21/26

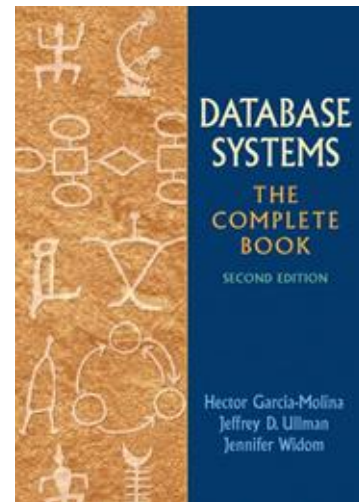
# Announcements

- Assignment 1
  - Preference survey: We will use this to determine technology presentation (assignment 3) grouping
- Office hours
  - Instructor: Thursdays 1:00-2:00, KACB 3322
  - Yihao: Tuesdays 2:00-3:00, in the common area next to Klaus 3324
  - Tianji: Fridays 2:00-3:00, in the common area next to Klaus 3324

# Reading Materials

Database Systems: The Complete Book (2nd edition)

- Chapter 3: Design Theory for Relational Databases  
(3.1 – 3.3)



Acknowledgement: The following slides have been adapted from EE477 (Database and Big Data Systems) taught by Steven Whang and CS145 (Intro to Big Data Systems) taught by Peter Bailis

# Agenda

1. Normal forms & functional dependencies
2. Finding functional dependencies
3. Closures, superkeys & keys

# 1. Normal forms & functional dependencies

# Normal Forms

- 1<sup>st</sup> Normal Form (1NF) = All tables are flat
- 2<sup>nd</sup> Normal Form = disused
- Boyce-Codd Normal Form (BCNF)
- 3<sup>rd</sup> Normal Form (3NF)
- 4<sup>th</sup> and 5<sup>th</sup> Normal Forms = see textbooks

# 1<sup>st</sup> Normal Form (1NF)

Student	Courses
Mary	{CS4440,CS6422}
Joe	{CS4440,CS6400}
...	...

Violates 1NF.

Student	Courses
Mary	CS4440
Mary	CS6422
Joe	CS4440
Joe	CS6400

In 1<sup>st</sup> NF

1NF Constraint: Types must be atomic!

# Normal Forms

- 1<sup>st</sup> Normal Form (1NF) = All tables are flat
- 2<sup>nd</sup> Normal Form = disused
- Boyce-Codd Normal Form (BCNF)
- 3<sup>rd</sup> Normal Form (3NF)
- 4<sup>th</sup> and 5<sup>th</sup> Normal Forms = see textbooks

DB designs based on functional dependencies, intended to prevent data anomalies

Our focus in this lecture + next one



# Data Anomalies

A poorly designed database causes *anomalies*:

Student	Course	Room
Mary	CS6400	B01
Joe	CS6400	B01
Sam	CS6400	B01
..	..	..

If every course is in only one room, contains redundant information!

# Data Anomalies

A poorly designed database causes *anomalies*:

Student	Course	Room
Mary	CS6400	B01
Joe	CS6400	C12
Sam	CS6400	B01
..	..	..

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

# Data Anomalies

A poorly designed database causes *anomalies*:

Student	Course	Room
..	..	..

If everyone drops the class, we lose  
what room the class is in!  
= a delete anomaly

# Data Anomalies

A poorly designed database causes *anomalies*:

...	CS6422	C12



Student	Course	Room
Mary	CS6400	B01
Joe	CS6400	B01
Sam	CS6400	B01
..	..	..

Similarly, we can't reserve a room without students = an insert anomaly

# Data Anomalies

Student	Course
Mary	CS6400
Joe	CS6422
Sam	CS6400
..	..

Course	Room
CS6400	B01
CS6422	C12

Eliminate anomalies by  
**decomposing relations.**

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

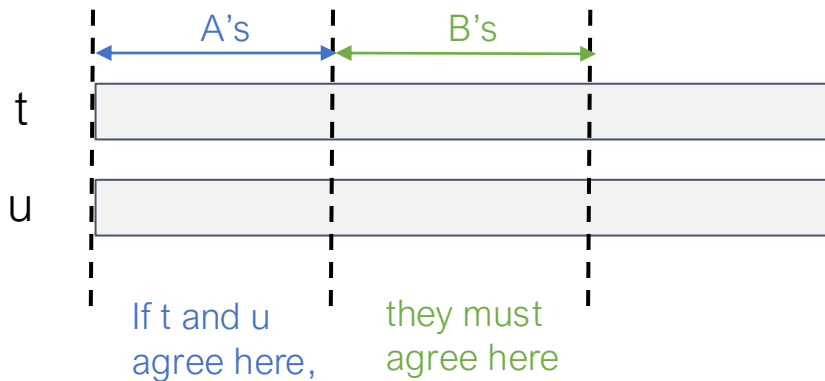
Goal: develop theory to understand why this design may be better and how to find this decomposition...

# Functional Dependencies

# Functional dependency (FD)

**Definition:** if two tuples of R agree on all the attributes  $A_1, A_2, \dots, A_n$ , they must also agree on (or functionally determine)  $B_1, B_2, \dots, B_m$

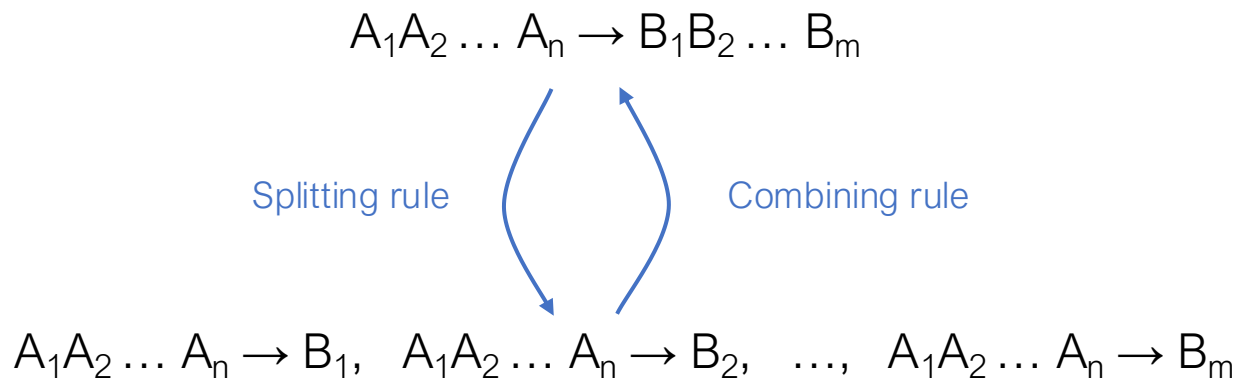
- Denoted as  $A_1 A_2 \dots A_n \rightarrow B_1 B_2 \dots B_m$



$A \rightarrow B$  means that  
“whenever two tuples agree on  
A then they agree on B.”

# Splitting/combining rule

- Splitting/combining can be applied to the **right sides** of FD's





# Splitting/combining rule

- For example,

title year  $\rightarrow$  length genre studioName



title year  $\rightarrow$  length

title year  $\rightarrow$  genre

title year  $\rightarrow$  studioName

# Splitting rule

- Splitting rule does not apply to the left sides of FD's

title year  $\rightarrow$  length



title  $\rightarrow$  length  
year  $\rightarrow$  length

# Functional Dependencies as Constraints

A functional dependency is a form of constraint

- Holds on some instances (but not others) – can check whether there are violations
- Part of the schema, helps define a valid instance

Student	Course	Room
Mary	CS6400	B01
Joe	CS6400	B01
Sam	CS6400	B01
..	..	..

Recall: an instance of a schema is a multiset of tuples conforming to that schema, i.e. a table

Note: The FD  
 $\{\text{Course}\} \rightarrow \{\text{Room}\}$   
holds on this instance

# Functional Dependencies as Constraints

Note that:

- You can check if an FD is **violated** by examining a single instance;
- However, you **cannot prove** that an FD is part of the schema by examining a single instance.
  - This would require checking every valid instance

Student	Course	Room
Mary	CS6400	B01
Joe	CS6400	B01
Sam	CS6400	B01
..	..	..

However, cannot prove that the FD {Course}  $\rightarrow$  {Room} is part of the schema

# Trivial functional dependencies

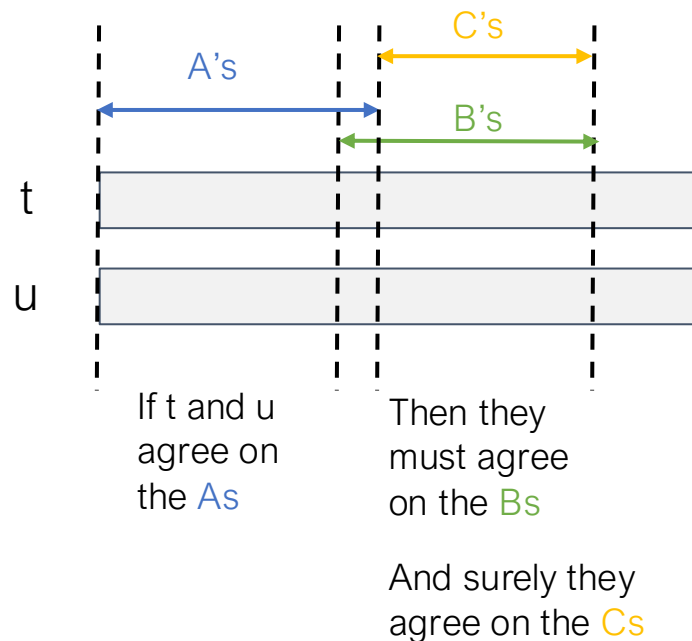
A constraint is *trivial* if it holds for every possible instance of the relation.

## Trivial FDs:

$A_1 A_2 \dots A_n \rightarrow B_1 B_2 \dots B_m$  such that  
 $\{B_1, B_2, \dots B_m\} \subseteq \{A_1, A_2, \dots, A_n\}$

## Trivial dependency rule:

$A_1 A_2 \dots A_n \rightarrow B_1 B_2 \dots B_m$  is equivalent  
to  $A_1 A_2 \dots A_n \rightarrow C_1 C_2 \dots C_k$ , where the  
C's are the B's that are not also A's



# In-class Exercise

Q1: Find an FD that holds on this instance

Q2: Find an FD that is violated on this instance

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	Lawyer

## 2. Finding functional dependencies

# FDs for Relational Schema Design

High-level idea: why do we care about FDs?

1. Start with some relational schema
2. Find out its functional dependencies (FDs)
3. Use these to design a better schema
  1. One which minimizes possibility of anomalies

This part can be tricky!



# Finding Functional Dependencies

There can be a large number of FDs...

Let's start with this problem:

Given a set of FDs,  $F = \{f_1, \dots, f_n\}$ , does an FD  $g$  hold?

Three simple rules called **Armstrong's Rules**.

1. Reflexivity,
2. Augmentation,
3. Transitivity

# Armstrong's axioms

You can derive any FDs that follows from a given set using these axioms:

## 1. Reflexivity:

If  $Y$  is a subset of  $X$ , then  $X \rightarrow Y$

This means that a set of attributes always determines a subset of itself

## 2. Augmentation:

If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$  for any  $Z$

This means we can add the same attributes to both sides of a functional dependency.

## 3. Transitivity:

If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$

This allows us to chain functional dependencies.

# Armstrong's axioms

- Does  $AB \rightarrow D$  follow from the FDs below?

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

1.  $AB \rightarrow C$  (given)
2.  $BC \rightarrow AD$  (given)

# Armstrong's axioms

- Does  $AB \rightarrow D$  follow from the FDs below?

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

1.  $AB \rightarrow C$  (given)
2.  $BC \rightarrow AD$  (given)
3.  $AB \rightarrow BC$  (Augmentation on 1)

# Armstrong's axioms

- Does  $AB \rightarrow D$  follow from the FDs below?

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

1.  $AB \rightarrow C$  (given)
2.  $BC \rightarrow AD$  (given)
3.  $AB \rightarrow BC$  (Augmentation on 1)
4.  $AB \rightarrow AD$  (Transitivity on 2,3)

# Armstrong's axioms

- Does  $AB \rightarrow D$  follow from the FDs below?

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

1.  $AB \rightarrow C$  (given)
2.  $BC \rightarrow AD$  (given)
3.  $AB \rightarrow BC$  (Augmentation on 1)
4.  $AB \rightarrow AD$  (Transitivity on 2,3)
5.  $AD \rightarrow D$  (Reflexivity)

# Armstrong's axioms

- Does  $AB \rightarrow D$  follow from the FDs below?

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

1.  $AB \rightarrow C$  (given)
2.  $BC \rightarrow AD$  (given)
3.  $AB \rightarrow BC$  (Augmentation on 1)
4.  $AB \rightarrow AD$  (Transitivity on 2,3)
5.  $AD \rightarrow D$  (Reflexivity)
6.  $AB \rightarrow D$  (Transitivity on 4,5)

Can we find an algorithmic way to do this?

# Closures



# Closure of attributes

Given a set of attributes  $A_1, \dots, A_n$  and a set of FDs  $F$ , the closure,  $\{A_1, \dots, A_n\}^+$  is the set of attributes  $B$  where  $\{A_1, \dots, A_n\} \rightarrow B$  follows from the FDs in  $F$

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

$\{A, B\}^+$

A, B

# Closure of attributes

Given a set of attributes  $A_1, \dots, A_n$  and a set of FDs  $F$ , the closure,  $\{A_1, \dots, A_n\}^+$  is the set of attributes  $B$  where  $\{A_1, \dots, A_n\} \rightarrow B$  follows from the FDs in  $F$

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

$\{A, B\}^+$

$A, B, C$

# Closure of attributes

Given a set of attributes  $A_1, \dots, A_n$  and a set of FDs  $F$ , the closure,  $\{A_1, \dots, A_n\}^+$  is the set of attributes  $B$  where  $\{A_1, \dots, A_n\} \rightarrow B$  follows from the FDs in  $F$

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

$\{A, B\}^+$

$A, B, C, D$

# Closure of attributes

Given a set of attributes  $A_1, \dots, A_n$  and a set of FDs  $F$ , the closure,  $\{A_1, \dots, A_n\}^+$  is the set of attributes  $B$  where  $\{A_1, \dots, A_n\} \rightarrow B$  follows from the FDs in  $F$

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

$\{A, B\}^+$

A, B, C, D, E

# Closure of attributes

Given a set of attributes  $A_1, \dots, A_n$  and a set of FDs  $F$ , the closure,  $\{A_1, \dots, A_n\}^+$  is the set of attributes  $B$  where  $\{A_1, \dots, A_n\} \rightarrow B$  follows from the FDs in  $F$

$AB \rightarrow C$

$BC \rightarrow AD$

$D \rightarrow E$

$CF \rightarrow B$

$\{A, B\}^+$

A, B, C, D, E

Cannot be expanded further, so this is a closure

# Closure algorithm

Start with  $X = \{A_1, \dots, A_n\}$  and set of FDs  $F$ .

**Repeat** until  $X$  doesn't change; **do**:

**if**  $\{B_1, \dots, B_n\} \rightarrow C$  is entailed by  $F$

        and  $\{B_1, \dots, B_n\} \subseteq X$

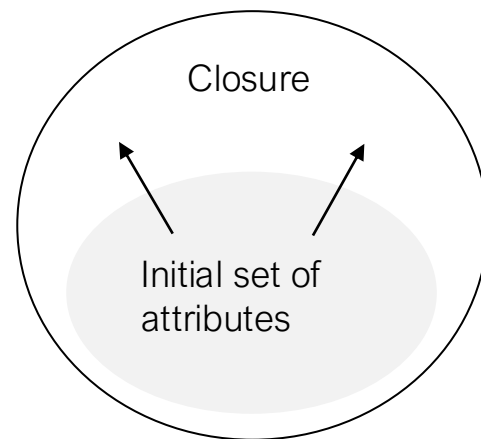
**then** add  $C$  to  $X$ .

**Return**  $X$  as  $X^+$

Helps to split the FD's of  $F$ , so each FD has a single attribute on the right

The algorithm (proof in book)

- only produces true FDs
- discovers all true FDs



# 3. Closures, Superkeys & Keys

# Why Do We Need the Closure?

With closure we can find all FD's easily

To check if  $X \rightarrow A$

1. Compute  $X^+$
2. Check if  $A \in X^+$

Note here that  $X$  is a set of attributes, but  $A$  is a single attribute. Why does considering FDs of this form suffice?

Recall the split/combine rule:  
 $X \rightarrow A_1, \dots, X \rightarrow A_n$   
implies  
 $X \rightarrow \{A_1, \dots, A_n\}$



# Using Closure to Infer ALL FDs

Example:

Given F =

$\{A,B\} \rightarrow C$

$\{A,D\} \rightarrow B$

$\{B\} \rightarrow D$

Step 1: Compute  $X^+$ , for every set of attributes X:

$\{A\}^+ = \{A\}$

$\{B\}^+ = \{B,D\}$

$\{C\}^+ = \{C\}$

$\{D\}^+ = \{D\}$

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

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

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

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

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

# Using Closure to Infer ALL FDs

Example:

Given F =

$\{A,B\} \rightarrow C$

$\{A,D\} \rightarrow B$

$\{B\} \rightarrow D$

Step 1: Compute  $X^+$ , for every set of attributes X:

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

Step 2: Enumerate all FDs  $X \rightarrow Y$ , s.t.  $Y \subseteq X^+$  and  $X \cap Y = \emptyset$ :

# Using Closure to Infer ALL FDs

Example:

Given F =

$\{A,B\} \rightarrow C$

$\{A,D\} \rightarrow B$

$\{B\} \rightarrow D$

Step 1: Compute  $X^+$ , for every set of attributes X:

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

*Y is in the  
closure of X*

Step 2: Enumerate all FDs  $X \rightarrow Y$ , s.t.  $Y \subseteq X^+$  and  $X \cap Y = \emptyset$ :

$\{A,B\} \rightarrow \{C,D\}$ ,  $\{A,D\} \rightarrow \{B,C\}$ ,  
 $\{A,B,C\} \rightarrow \{D\}$ ,  $\{A,B,D\} \rightarrow \{C\}$ ,  
 $\{A,C,D\} \rightarrow \{B\}$

*The FD  
 $X \rightarrow Y$  is  
non-trivial*

# Minimal basis

The full set of implied FDs can be large and redundant...

For the purpose of data normalization, it's often easier to work with the cleanest, smallest set of FDs.

A **minimal basis** (or minimal cover) for a set of FDs  $F$  is a simplified set of FDs  $G$  that satisfies the following conditions:

- No redundant/extraneous FDs
- RHS has a single attribute
- No extraneous attributes on the LHS

Given a set of FD's  $F$ , any set of FD's equivalent to  $F$  is a **basis** for  $F$

# Minimal basis generation

Input:  $F = \{A \rightarrow AB, AB \rightarrow C\}$

1. Split FD's so that they have singleton right sides

$G = \{A \rightarrow B, A \rightarrow A, AB \rightarrow C\}$

2. Remove trivial FDs

$G = \{A \rightarrow B, AB \rightarrow C\}$

3. Minimize the left sides of each FD

$G = \{A \rightarrow B, A \rightarrow C\}$

4. Remove redundant FDs

$G = \{A \rightarrow B, A \rightarrow C\}$

Step 3:

*For each FD  $X \rightarrow A$  in  $F$ :  
For each attribute  $B$  in  $X$ :  
If  $(X - \{B\})^+$  contains  $A$ ,  
remove  $B$  from  $X$ .*

# Why Do We Need the Closure?

With closure we can find **keys** and **superkeys** of a relation

For each set of attributes  $X$

1. Compute  $X^+$
2. If  $X^+ =$  set of all attributes then  $X$  is a **superkey**
3. If  $X$  is minimal, then it is a **key**

# Keys and Superkeys

A superkey is a set of attributes  $A_1, \dots, A_n$   
s.t.  
for any other attribute  $B$  in  $R$ ,  
we have  $\{A_1, \dots, A_n\} \rightarrow B$

i.e. all attributes are  
functionally determined  
by a superkey

A key is a minimal  
superkey

This means that no subset of a key  
is also a superkey  
(i.e., dropping any attribute from the  
key makes it no longer a superkey)

## Example of Finding Keys

Product(name, price, category, color)

{name, category} → price

{category} → color

What is a key?



# Example of Finding Keys

Product(name, price, category, color)

{name, category} → price

{category} → color

1. Is a key always guaranteed to exist?
2. Can we have more than one key?

$\{\text{name, category}\}^+ = \{\text{name, price, category, color}\}$

⇒ this is a **superkey**

⇒ this is a **key**, since neither **name** nor **category** alone is a superkey

# In-class Exercise

Given  $R(A, B, C, D)$  and FD's  $AB \rightarrow C$ ,  $C \rightarrow D$ ,  $D \rightarrow A$

- What are all keys of  $R$ ?