# Introduction to Data Management
# CSE 344

## Lectures 18: Design Theory

# Announcements

- HW6 and WQ6 are out
    - Due on Monday, 11/21

- Back to Tues/Wed cycle for HW7-8 and WQ7


- Today and next lecture:
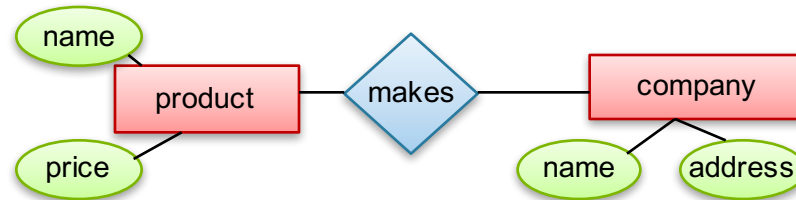    - Design theory (3.1-3.4)

# Where are we?

- ## First half of 344:

  - ### Data models: instance, schema, languages

    - Relational and NoSQL

  - ### Query processing

- ## Second half of 344: Using DBMSs effectively

  - ### Conceptual design

  - ### Transactions

  - ### Parallel databases

# What is this class about?

- **Focus: Using DBMSs**
- Relational Data Model
  - SQL, Relational Algebra, Relational Calculus, datalog
- Semistructured Data Model
  - JSon, CouchDB (NoSQL)
- Conceptual design
  - E/R diagrams, Views, and Database normalization
- Transactions
- Parallel databases, MapReduce, and Spark
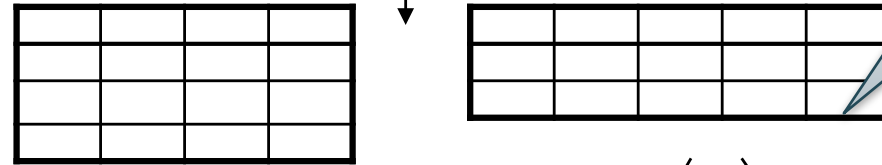- Data integration and data cleaning
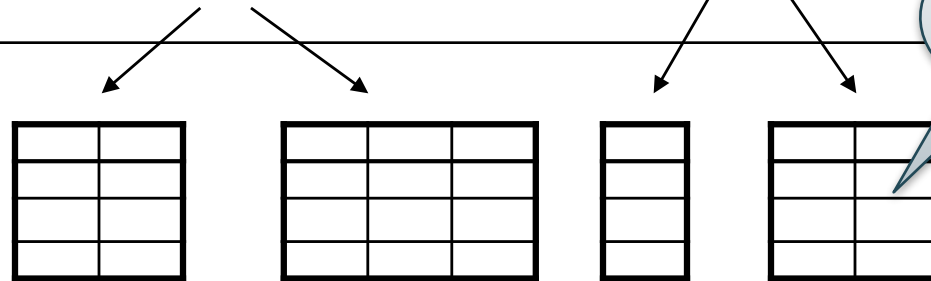
# Database Design Process

Conceptual Model:

Lec 17

name
product — makes — company
price
name    address

Relational Model:
Tables + constraints
And also functional dep.

Sec 7 +
Lec 18

Normalization:
Eliminates anomalies

Lec 18
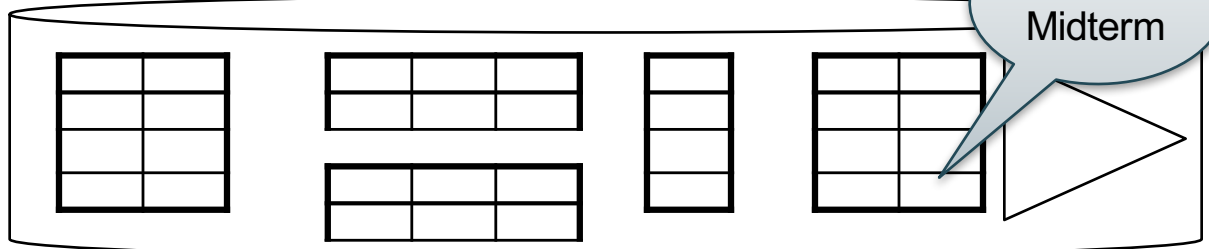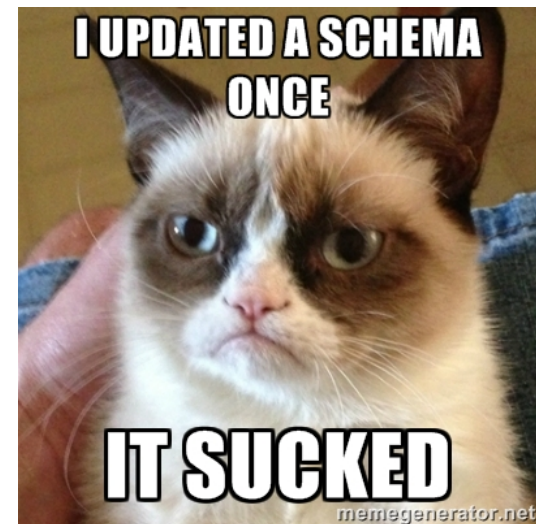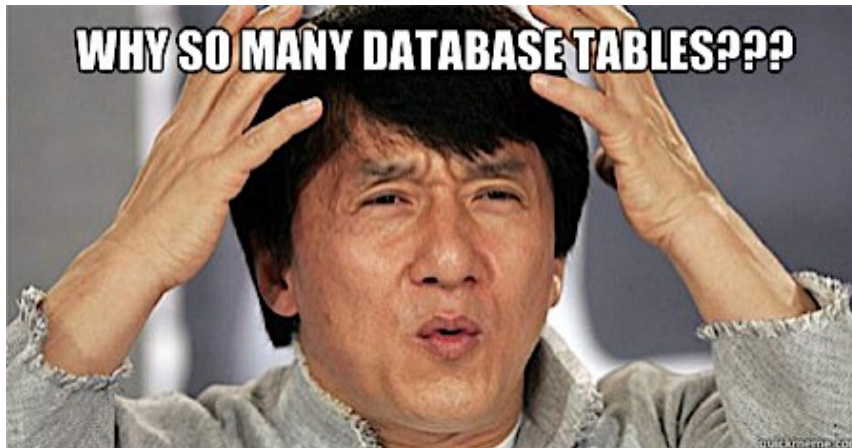Lec 19

Conceptual Schema

Physical storage details

Before
Midterm

Physical Schema

# What makes good schemas?

# Relational Schema Design

| Name | SSN | PhoneNumber | City |
|------|-----|-------------|------|
| Fred | 123-45-6789 | 206-555-1234 | Seattle |
| Fred | 123-45-6789 | 206-555-6543 | Seattle |
| Joe | 987-65-4321 | 908-555-2121 | Westfield |

One person may have multiple phones, but lives in only one city

Primary key is thus (SSN, PhoneNumber)

What is the problem with this schema?

# Relational Schema Design

| Name | SSN | PhoneNumber | City |
|------|-----|-------------|------|
| Fred | 123-45-6789 | 206-555-1234 | Seattle |
| Fred | 123-45-6789 | 206-555-6543 | Seattle |
| Joe | 987-65-4321 | 908-555-2121 | Westfield |

## Anomalies:

- Redundancy = repeat data
- Update anomalies = what if Fred moves to "Bellevue"?
- Deletion anomalies = what if Joe deletes his phone number?

# Relation Decomposition

**Break the relation into two:**

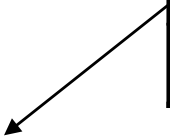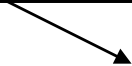| Name | SSN | PhoneNumber | City |
|------|-----|-------------|------|
| Fred | 123-45-6789 | 206-555-1234 | Seattle |
| Fred | 123-45-6789 | 206-555-6543 | Seattle |
| Joe | 987-65-4321 | 908-555-2121 | Westfield |

| Name | SSN | City |
|------|-----|------|
| Fred | 123-45-6789 | Seattle |
| Joe | 987-65-4321 | Westfield |

| SSN | PhoneNumber |
|-----|-------------|
| 123-45-6789 | 206-555-1234 |
| 123-45-6789 | 206-555-6543 |
| 987-65-4321 | 908-555-2121 |

## Anomalies have gone:

- No more repeated data
- Easy to move Fred to "Bellevue" (how ?)
- Easy to delete all Joe's phone numbers (how ?)

# Relational Schema Design (or Logical Design)

How do we do this systematically?

- Start with some relational schema

- Find out its ***functional dependencies*** (FDs)

- Use FDs to ***normalize*** the relational schema

# Functional Dependencies (FDs)

**<u>Definition</u>**

If two tuples agree on the attributes

$$A_1, A_2, \ldots, A_n$$

then they must also agree on the attributes

$$B_1, B_2, \ldots, B_m$$

Formally:

$$A_1, A_2, \ldots, A_n \rightarrow B_1, B_2, \ldots, B_m$$

$A_1 \ldots A_n$ **determines** $B_1 .. B_m$

# Functional Dependencies (FDs)

**Definition** $A_1, ..., A_m \rightarrow B_1, ..., B_n$ **holds** in R if:

$\forall t, t' \in R,$

$(t.A_1 = t'.A_1 \wedge ... \wedge t.A_m = t'.A_m \rightarrow t.B_1 = t'.B_1 \wedge ... \wedge t.B_n = t'.B_n )$

| R | | $A_1$ | ... | $A_m$ | | $B_1$ | ... | $B_n$ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | |
| t | | | | | | | | | | |
| | | | | | | | | | | |
| t' | | | | | | | | | | |
| | | | | | | | | | | |

if t, t' agree here then t, t' agree here

# Example

An FD <u>holds</u>, or <u>does not hold</u> on an instance:

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

EmpID  →   Name, Phone, Position

Position  →   Phone

but  not  Phone  →   Position

# Example

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

Position → Phone

# Example

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

But not Phone  →    Position

# Example

name → color
category → department
color, category → price

| name | category | color | department | price |
|------|----------|-------|------------|-------|
| Gizmo | Gadget | Green | Toys | 49 |
| Tweaker | Gadget | Green | Toys | 99 |

Do all the FDs hold on this instance?

# Example

name → color
category → department
color, category → price

| name | category | color | department | price |
|------|----------|-------|------------|-------|
| Gizmo | Gadget | Green | Toys | 49 |
| Tweaker | Gadget | Green | Toys | 49 |
| Gizmo | Stationary | Green | Office-supp. | 59 |

What about this one ?

# Terminology

- FD **holds** or **does not hold** on an instance

- If we can be sure that *every instance of R* will be one in which a given FD is true, then we say that **R satisfies the FD**

- If we say that R satisfies an FD F, we are **stating a constraint on R**
  - Recall constraints from lec 17 and sec 7

# Why bother with FDs?

| Name | SSN | PhoneNumber | City |
|------|-----|-------------|------|
| Fred | 123-45-6789 | 206-555-1234 | Seattle |
| Fred | 123-45-6789 | 206-555-6543 | Seattle |
| Joe | 987-65-4321 | 908-555-2121 | Westfield |

## Anomalies:
- Redundancy          = repeat data
- Update anomalies  = what if Fred moves to "Bellevue"?
- Deletion anomalies = what if Joe deletes his phone number?

# An Interesting Observation

If all these FDs are true:

> name → color
> category → department
> color, category → price

Then this FD also holds:

> name, category → price

If we find out from application domain that a relation satisfies some FDs, it doesn't mean that we found all the FDs that it satisfies!
There could be more FDs implied by the ones we have.

# Closure of a set of Attributes

**Given** a set of attributes $A_1, \ldots, A_n$

The **closure** is the set of attributes B, notated $\{A_1, \ldots, A_n\}^+$,
$$\text{s.t. } A_1, \ldots, A_n \rightarrow B$$

Example:

1. name → color
2. category → department
3. color, category → price

Closures:

name$^+$ = {name, color}

{name, category}$^+$ = {name, category, color, department, price}

color$^+$ = {color}

# Closure Algorithm

X={A1, …, An}.

**Repeat until** X doesn't change **do**:
  **if**    $B_1, …, B_n \rightarrow C$  is a FD **and**
        $B_1, …, B_n$  are all in X
  **then**  add C to X.

Example:

1. name → color
2. category → department
3. color, category → price

{name, category}$^+$ =
    { name, category, color, department, price }

Hence: name, category → color, department, price

# Example

In class:

R(A,B,C,D,E,F)

A, B → C
A, D → E
B   → D
A, F → B

Compute {A,B}⁺   X = {A, B,                    }

Compute {A, F}⁺   X = {A, F,                    }

# Example

In class:

R(A,B,C,D,E,F)

$$A, B \rightarrow C$$
$$A, D \rightarrow E$$
$$B \rightarrow D$$
$$A, F \rightarrow B$$

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

Compute $\{A, F\}^+$    X = {A, F,                    }

# Example

In class:

R(A,B,C,D,E,F)

A, B → C
A, D → E
B → D
A, F → B

Compute {A,B}⁺   X = {A, B, C, D, E }

Compute {A, F}⁺   X = {A, F, B, C, D, E }

# Example

In class:

R(A,B,C,D,E,F)

A, B → C
A, D → E
B → D
A, F → B

Compute {A,B}+    X = {A, B, C, D, E }

Compute {A, F}+   X = {A, F, B, C, D, E }

What is the key of R?

# Practice at Home

Find all FD's implied by:

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

# Practice at Home

Find all FD's implied by:

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

Step 1: Compute $X^+$, for every X:

A+ = A,   B+ = BD,   C+ = C,   D+ = D

AB+ =ABCD, AC+=AC, AD+=ABCD,
             BC+=BCD,  BD+=BD,  CD+=CD

ABC+ = ABD+ = ACD$^+$ = ABCD (no need to compute– why ?)

BCD$^+$ = BCD,    ABCD+ = ABCD

Step 2: Enumerate all FD's X $\rightarrow$ Y, s.t. Y $\subseteq$ X$^+$ and X $\cap$ Y = $\varnothing$ :

AB $\rightarrow$ CD, AD$\rightarrow$BC,  ABC $\rightarrow$ D, ABD $\rightarrow$ C, ACD $\rightarrow$ B

# Keys

- A **superkey** is a set of attributes $A_1, ..., A_n$ s.t. for any other attribute B, we have $A_1, ..., A_n \rightarrow B$

- A **key** is a minimal superkey
  - A superkey and for which no subset is a superkey

# Computing (Super)Keys

- For all sets X, compute $X^+$

- If $X^+$ = [all attributes], then X is a superkey

- Try reducing to the minimal X's to get the key

# Example

Product(name, price, category, color)

name, category → price
category → color

What is the key ?

# Example

Product(name, price, category, color)

> name, category → price
> category → color

What is the key ?

(name, category) + = { name, category, price, color }

Hence (name, category) is a key

# Key or Keys ?

Can we have more than one key ?

Given R(A,B,C) define FD's s.t. there are two or more
    distinct keys

# Key or Keys ?

Can we have more than one key ?

Given R(A,B,C) define FD's s.t. there are two or more
   distinct keys

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

or

$$AB \rightarrow C$$
$$BC \rightarrow A$$

or

$$A \rightarrow BC$$
$$B \rightarrow AC$$

what are the keys here ?

# Eliminating Anomalies

Main idea:

- $X \rightarrow A$ is OK if X is a (super)key

- $X \rightarrow A$ is not OK otherwise
  - Need to decompose the table, but how?

## Boyce-Codd Normal Form