

# Introduction to Databases

## CSE 414

### Lecture 2: Data Models

# Class Overview

- Unit 1: Intro
- Unit 2: Relational Data Models and Query Languages
  - Data models, SQL, Relational Algebra, Datalog
- Unit 3: Non-relational data
- Unit 4: RDMBS internals and query optimization
- Unit 5: Parallel query processing
- Unit 6: DBMS usability, conceptual design
- Unit 7: Transactions

# Review

- What is a database?
  - A collection of files storing related data
- What is a DBMS?
  - An application program that allows us to manage efficiently the collection of data files

# Data Models

- Recall our example: want to design a database of books:
  - author, title, publisher, pub date, price, etc
  - How should we describe this data?
- **Data model** = mathematical formalism (or conceptual way) for describing the data

# Data Models

- Relational
  - Data represented as relations
- Semi-structured (JSON)
  - Data represented as trees
- Key-value pairs
  - Used by NoSQL systems
- Graph
- Object-oriented

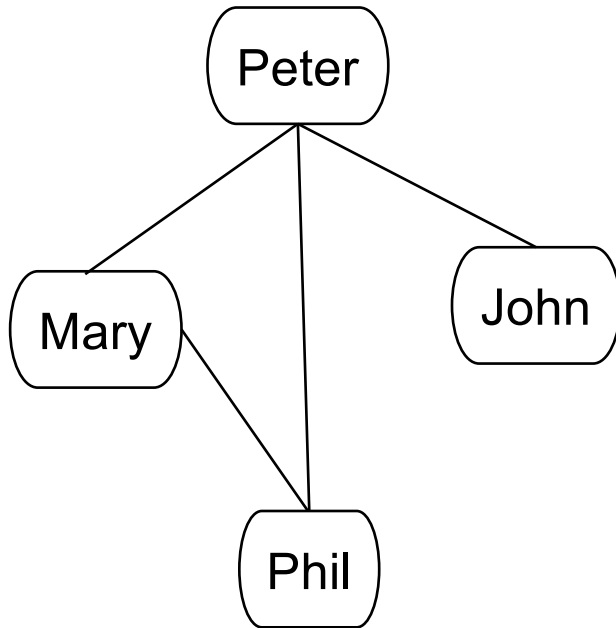


Unit 2



Unit 3

# Example: storing FB friends



As a graph

Or

Person1	Person2	is_friend
Peter	John	1
John	Mary	0
Mary	Phil	1
Phil	Peter	1
...	...	...

As a relation

We will learn the tradeoffs of different data models later this quarter

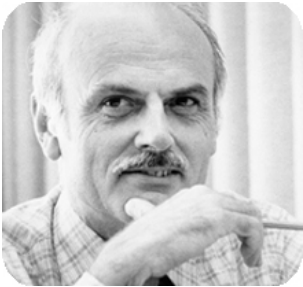
# 3 Elements of Data Models

- Instance
  - The actual data
- Schema
  - Describe what data is being stored
- Query language
  - How to retrieve and manipulate data

# Turing Awards in Data Management



Charles Bachman, 1973  
*IDS and CODASYL*



Ted Codd, 1981  
*Relational model*



Jim Gray, 1998  
*Transaction processing*



Michael Stonebraker, 2014  
*INGRES and Postgres*





# Relational Model

columns /  
attributes /  
fields

- Data is a collection of relations / tables:

cname	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

- mathematically, relation is a set of tuples
  - each tuple appears 0 or 1 times in the table
  - order of the rows is unspecified

# The Relational Data Model

- Degree (arity) of a relation = #attributes
- Each attribute has a type.
  - Examples types:
    - Strings: CHAR(20), VARCHAR(50), TEXT
    - Numbers: INT, SMALLINT, FLOAT
    - MONEY, DATETIME, ...
    - Few more that are vendor specific
  - Statically and strictly enforced

# Keys

- Key = one (or multiple) attributes that uniquely identify a record

# Keys

- Key = one (or multiple) attributes that uniquely identify a record

Key

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

# Keys

- Key = one (or multiple) attributes that uniquely identify a record

The diagram shows a table with four columns: cname, country, no\_employees, and for\_profit. A callout bubble labeled 'Key' points to the cname column. Another callout bubble labeled 'Not a key' points to the country column.

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

# Keys

- Key = one (or multiple) attributes that uniquely identify a record

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

# Keys

- Key = one (or multiple) attributes that uniquely identify a record

Key

Not a key

Is this a key?

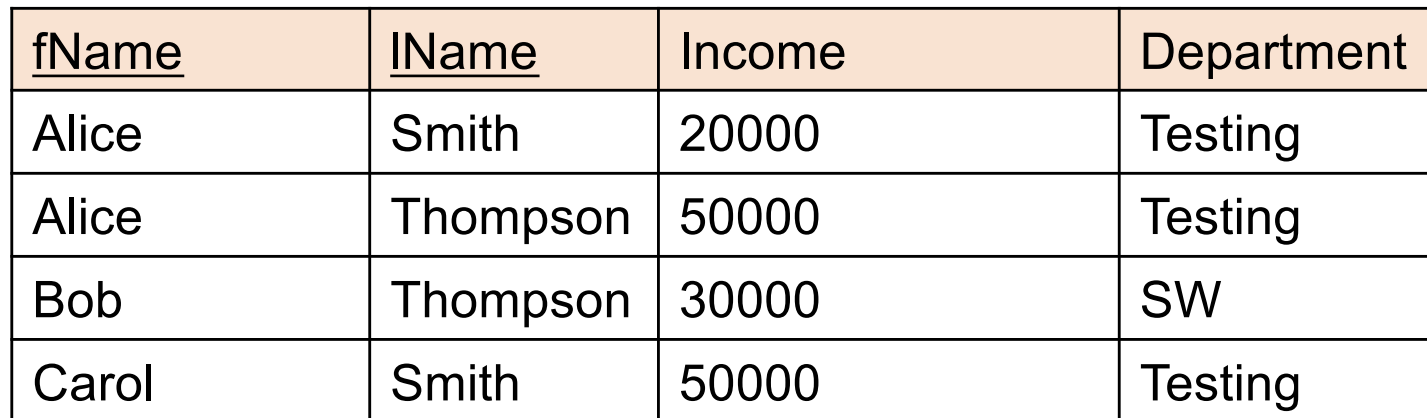
No: future updates to the database may create duplicate no\_employees

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

The diagram shows a table with four columns: cname, country, no\_employees, and for\_profit. Three callouts are present: 'Key' points to the cname column, 'Not a key' points to the country column, and 'Is this a key?' points to the no\_employees column. A text block to the right states: 'No: future updates to the database may create duplicate no\_employees'.

# Multi-attribute Key

Key = fName, lName  
(what does this mean?)



<u>fName</u>	<u>lName</u>	Income	Department
Alice	Smith	20000	Testing
Alice	Thompson	50000	Testing
Bob	Thompson	30000	SW
Carol	Smith	50000	Testing



# Multiple Keys

The diagram illustrates two keys for a table. A speech bubble labeled 'Key' points to the SSN column. Another speech bubble labeled 'Another key' points to the fName, IName, and Income columns.

<u>SSN</u>	fName	IName	Income	Department
111-22-3333	Alice	Smith	20000	Testing
222-33-4444	Alice	Thompson	50000	Testing
333-44-5555	Bob	Thompson	30000	SW
444-55-6666	Carol	Smith	50000	Testing

We can choose one key and designate it as primary key

E.g.: primary key = SSN

# Foreign Key

Company(cname, country, no\_employees, for\_profit)  
Country(name, population)

Company

<u>cname</u>	country	no_employees	for_profit
Canon	Japan	50000	Y
Hitachi	Japan	30000	Y

Foreign key to  
Country.name

Country

<u>name</u>	population
USA	320M
Japan	127M

# Keys: Summary

- Key = columns that uniquely identify tuple
  - Usually we underline
  - A relation can have many keys, but only one can be chosen as *primary key*
- Foreign key:
  - Attribute(s) whose value is a key of a record in some other relation
  - Foreign keys are sometimes called *semantic pointer*

# Query Language

- SQL
  - **Structured Query Language**
  - Developed by IBM in the 70s
  - Most widely used language to query relational data
- Other relational query languages
  - Datalog, relational algebra

# Our First DBMS

- SQL Lite
- Will switch to SQL Server later in the quarter

# Demo 1

# Discussion

- Tables are NOT ordered
  - they are sets or multisets (bags)
- Tables are FLAT
  - No nested attributes
- Tables DO NOT prescribe how they are implemented / stored on disk
  - This is called **physical data independence**

# Table Implementation

- How would you implement this?

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False



# Table Implementation

- How would you implement this?

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

Row major: as an array of objects

GizmoWorks	Canon	Hitachi	HappyCam
USA	Japan	Japan	Canada
20000	50000	30000	500
True	True	True	False

# Table Implementation

- How would you implement this?

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

Column major: as one array per attribute

GizmoWorks	Canon	Hitachi	HappyCam
USA	Japan	Japan	Canada
20000	50000	30000	500
True	True	True	False

# Table Implementation

- How would you implement this?

<u>cname</u>	country	no_employees	for_profit
GizmoWorks	USA	20000	True
Canon	Japan	50000	True
Hitachi	Japan	30000	True
HappyCam	Canada	500	False

## Physical data independence

The logical definition of the data remains unchanged, even when we make changes to the actual implementation

# First Normal Form

<u>cname</u>	country	no_employees	for_profit
Canon	Japan	50000	Y
Hitachi	Japan	30000	Y

- All relations must be flat: we say that the relation is in *first normal form*

# First Normal Form

<u>cname</u>	country	no_employees	for_profit
Canon	Japan	50000	Y
Hitachi	Japan	30000	Y

- All relations must be flat: we say that the relation is in *first normal form*
- E.g. we want to add products manufactured by each company:

# First Normal Form

<u>cname</u>	country	no_employees	for_profit
Canon	Japan	50000	Y
Hitachi	Japan	30000	Y

- All relations must be flat: we say that the relation is in *first normal form*
- E.g., we want to add products manufactured by each company:

<u>cname</u>	country	no_employees	for_profit	products									
Canon	Japan	50000	Y	<table border="1"> <thead> <tr> <th><u>pname</u></th> <th>price</th> <th>category</th> </tr> </thead> <tbody> <tr> <td>SingleTouch</td> <td>149.99</td> <td>Photography</td> </tr> <tr> <td>Gadget</td> <td>200</td> <td>Toy</td> </tr> </tbody> </table>	<u>pname</u>	price	category	SingleTouch	149.99	Photography	Gadget	200	Toy
<u>pname</u>	price	category											
SingleTouch	149.99	Photography											
Gadget	200	Toy											
Hitachi	Japan	30000	Y	<table border="1"> <thead> <tr> <th><u>pname</u></th> <th>price</th> <th>category</th> </tr> </thead> <tbody> <tr> <td>AC</td> <td>300</td> <td>Appliance</td> </tr> </tbody> </table>	<u>pname</u>	price	category	AC	300	Appliance			
<u>pname</u>	price	category											
AC	300	Appliance											

# First Normal Form

<u>cname</u>	country	no_employees	for_profit
Canon	Japan	50000	Y
Hitachi	Japan	30000	Y

- All relations must be flat: we say that the relation is in *first normal form*
- E.g., we want to add products manufactured by each company:

Non-1NF!

<u>cname</u>	country	no_employees	for_profit	products									
Canon	Japan	50000	Y	<table border="1"> <thead> <tr> <th><u>pname</u></th> <th>price</th> <th>category</th> </tr> </thead> <tbody> <tr> <td>SingleTouch</td> <td>149.99</td> <td>Photography</td> </tr> <tr> <td>Gadget</td> <td>200</td> <td>Toy</td> </tr> </tbody> </table>	<u>pname</u>	price	category	SingleTouch	149.99	Photography	Gadget	200	Toy
<u>pname</u>	price	category											
SingleTouch	149.99	Photography											
Gadget	200	Toy											
Hitachi	Japan	30000	Y	<table border="1"> <thead> <tr> <th><u>pname</u></th> <th>price</th> <th>category</th> </tr> </thead> <tbody> <tr> <td>AC</td> <td>300</td> <td>Appliance</td> </tr> </tbody> </table>	<u>pname</u>	price	category	AC	300	Appliance			
<u>pname</u>	price	category											
AC	300	Appliance											

# First Normal Form

Now it's in 1NF

## Company

<u>cname</u>	country	no_employees	for_profit
Canon	Japan	50000	Y
Hitachi	Japan	30000	Y

## Products

<u>pname</u>	price	category	manufacturer
SingleTouch	149.99	Photography	Canon
AC	300	Appliance	Hitachi
Gadget	200	Toy	Canon



# Demo 1 (cont'd)

# Data Models: Summary

- Schema + Instance + Query language
- Relational model:
  - Database = collection of tables
  - Each table is flat: “first normal form”
  - Key: may consists of multiple attributes
  - Foreign key: “semantic pointer”
  - Physical data independence