

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

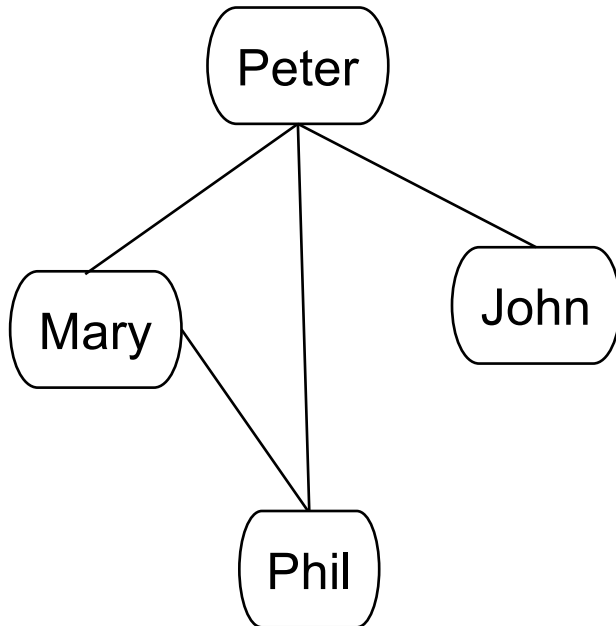


Unit 2



Unit 3

Example: storing FB friends



Or

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

As a graph

As a relation

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

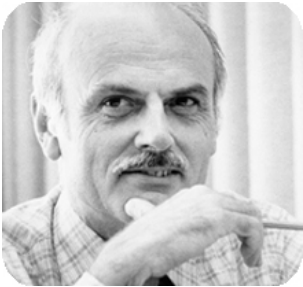
3 Elements of Data Models

- 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



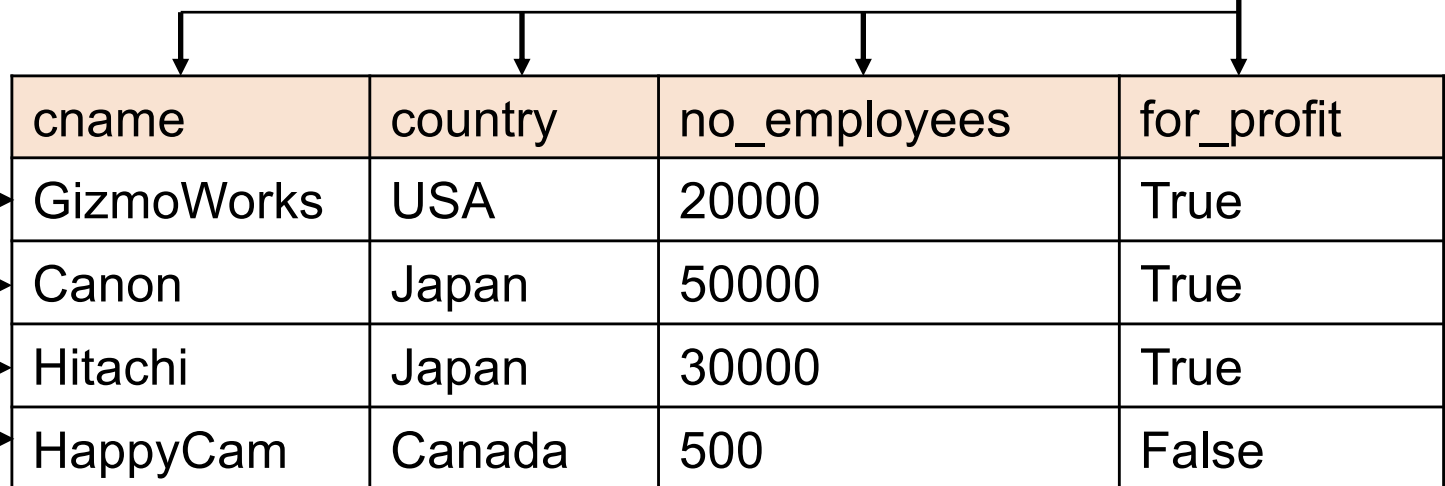
Michael Stonebraker, 2014
INGRES and Postgres



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

rows /
tuples /
records

- 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

Uniqueness should hold true for current AND all future tuples/rows that might be added

Keys

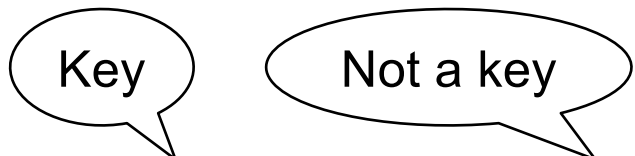
- Key = one (or multiple) attributes that uniquely identify a record

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



<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

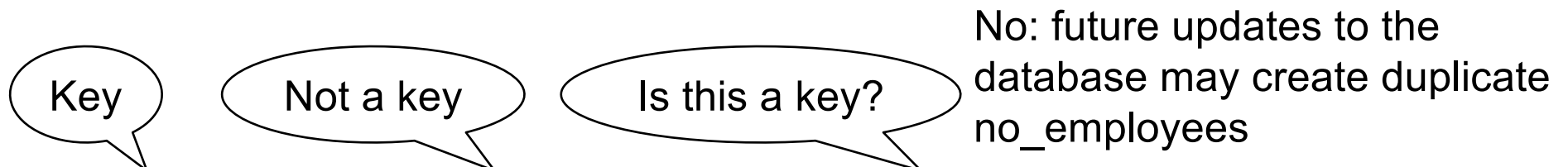
- Key = one (or multiple) attributes that uniquely identify a record



<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




The diagram shows a table with four columns: cname, country, no_employees, and for_profit. Three callouts are present: a speech bubble labeled 'Key' pointing to the cname column; an oval labeled 'Not a key' pointing to the country column; and an oval labeled 'Is this a key?' pointing to the no_employees column. To the right of the table, text states: '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

Multi-attribute Key

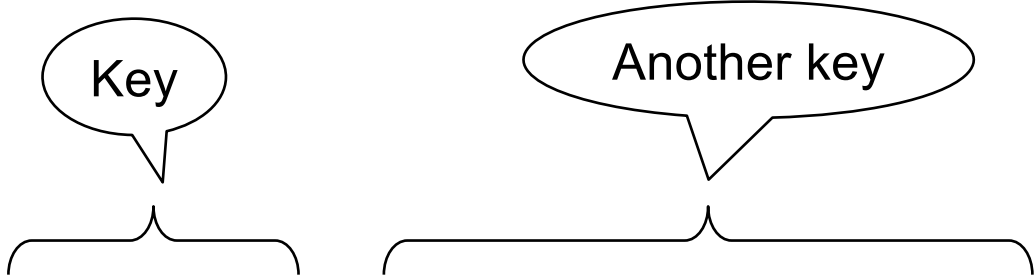
Key = fName, lName
(what does this mean?)

sometimes it takes more than a single attribute to uniquely define a tuple/row



<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



<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

Why have secondary keys?

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

every value/row in first table in our country columns will correspond to a unique row in our second table.

i.e. A foreign key is a column whose entries are keys in another table/relation

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*

Instances = actual data (rows of our data table)

Schema = what our variables/fields of the table are and what type they are

Query language = how we interact with this data

Query Language

- SQL
 - **Structured Query Language**
 - Developed by IBM in the 70s
 - Most widely used language to query relational data declarative programming language
“I want this piece of data, send me however you think is best?”
- Other relational query languages
 - Datalog, relational algebra

Our First DBMS

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

Demo 1

SHORT LIST OF COMMANDS

```
create table table_name(attribute1 type1, attribute2 type2, ...);
```

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><tr><th><u>pname</u></th><th>price</th><th>category</th></tr><tr><td>SingleTouch</td><td>149.99</td><td>Photography</td></tr><tr><td>Gadget</td><td>200</td><td>Toy</td></tr></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><tr><th><u>pname</u></th><th>price</th><th>category</th></tr><tr><td>AC</td><td>300</td><td>Appliance</td></tr></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	<u>pname</u>	price	category
				SingleTouch	149.99	Photography
				Gadget	200	Toy
Hitachi	Japan	30000	Y	<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