

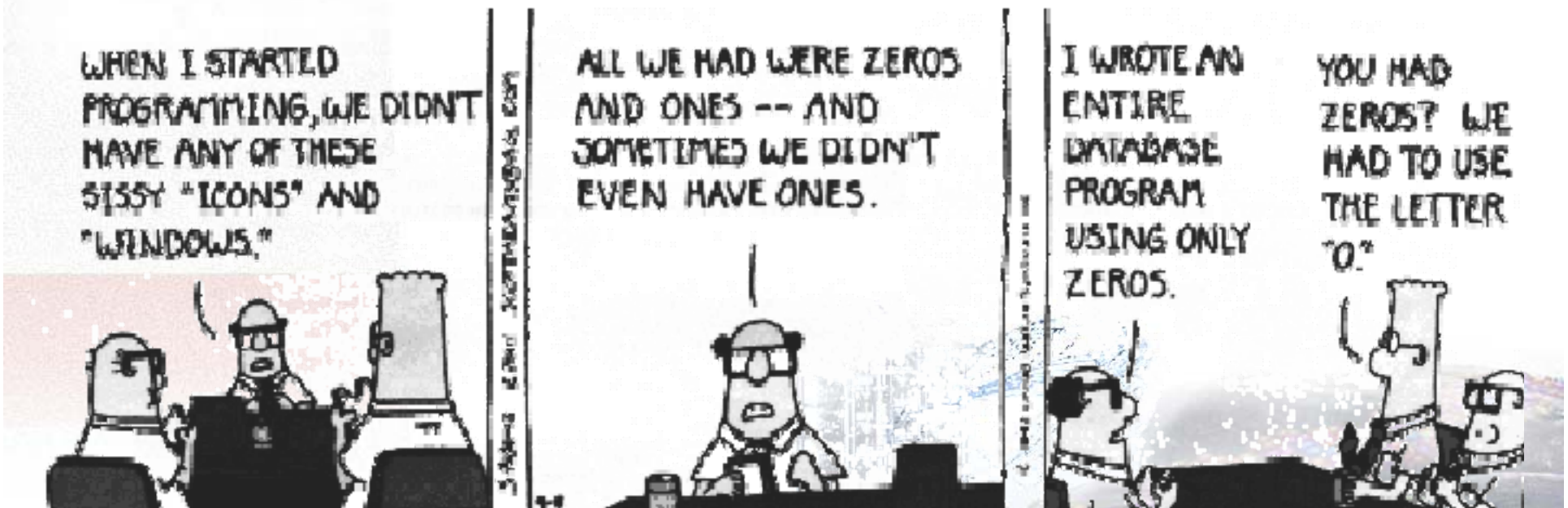


Introduction to Informatics

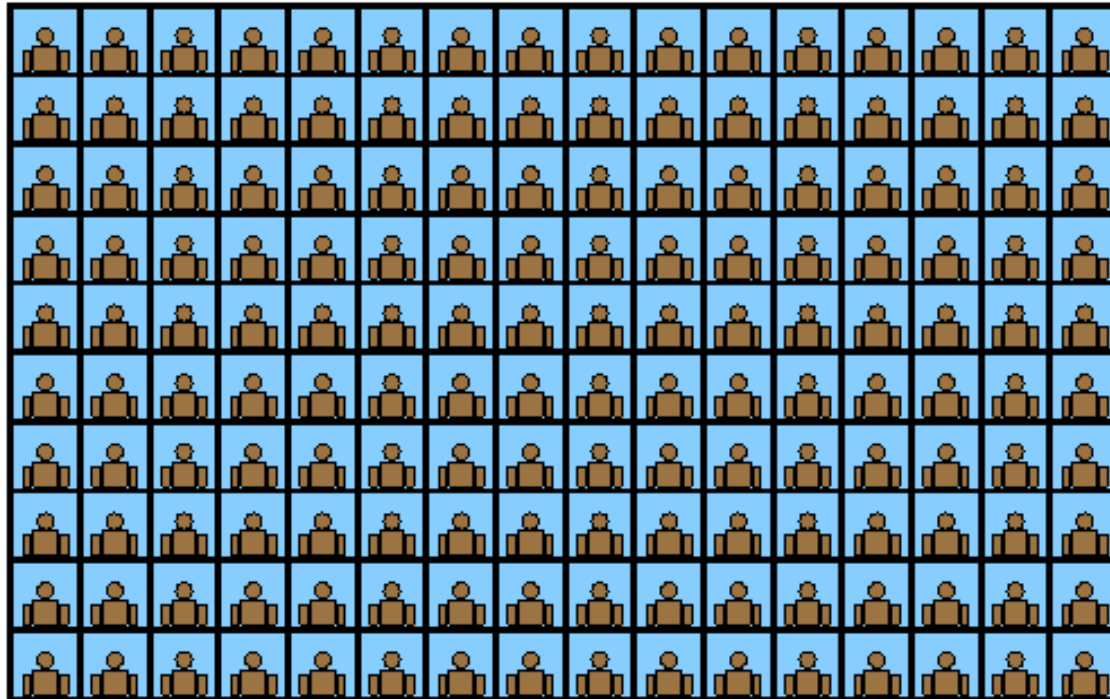
Lecture 26:

Information Technology in the Real World

Databases



NO MORE LABS !!!



Exam Schedule

- 11595
 - Midterm
 - March 1st (Thursday)
 - Regular Class time
 - Final Exam
 - May 3rd (Thursday)
 - 7:15-9:15 p.m.

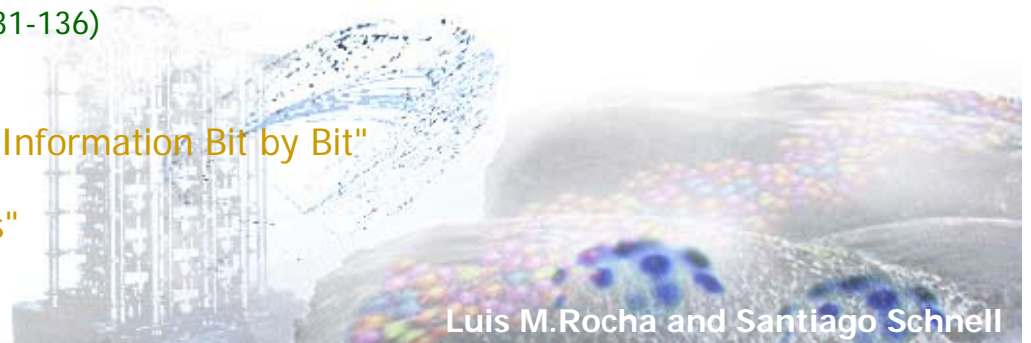
OH NO! OH NO!



Luis M.Rocha and Santiago Schnell

Readings until now

- Lecture notes
 - Posted online
 - <http://informatics.indiana.edu/rocha/i101>
 - *The Nature of Information*
 - *Technology*
 - *Modeling the World*
 - @ infoport
 - <http://infoport.blogspot.com>
 - From course package
 - Von Baeyer, H.C. [2004]. *Information: The New Language of Science*. Harvard University Press.
 - Chapters 1, 4 (pages 1-12)
 - Chapter 10 (pages 13-17)
 - From Andy Clark's book "*Natural-Born Cyborgs*"
 - Chapters 2 and 6 (pages 19 - 67)
 - From Irv Englander's book "*The Architecture of Computer Hardware and Systems Software*"
 - Chapter 3: Data Formats (pp. 70-86)
 - Klir, J.G., U. St. Clair, and B.Yuan [1997]. *Fuzzy Set Theory: foundations and Applications*. Prentice Hall
 - Chapter 2: Classical Logic (pp. 87-97)
 - Chapter 3: Classical Set Theory (pp. 98-103)
 - Norman, G.R. and D.L. Streinrt [2000]. *Biostatistics: The Bare Essentials*.
 - Chapters 1-3 (pages 105-129)
 - OPTIONAL: Chapter 4 (pages 131-136)
 - Chapter 13 (pages 147-155)
 - Chapter 5 (pages 141-144)
 - Igor Aleksander, "Understanding Information Bit by Bit"
 - Pages 157-166
 - Ellen Ullman, "Dining with Robots"
 - Pages 167-172



Assignment Situation

■ Labs

■ Past

- Lab 1: Blogs
 - Closed (Friday, January 19): Grades Posted
- Lab 2: Basic HTML
 - Closed (Wednesday, January 31): Grades Posted
- Lab 3: Advanced HTML: Cascading Style Sheets
 - Closed (Friday, February 2): Grades Posted
- Lab 4: More HTML and CSS
 - Closed (Friday, February 9): Grades Posted
- Lab 5: Introduction to Operating Systems: Unix
 - Closed (Friday, February 16): Grades Posted
- Lab 6: More Unix and FTP
 - Closed (Friday, February 23): Grades Posted
- Lab 7: Logic Gates
 - Closed (Friday, March 9): Grades Posted
- Lab 8: Intro to Statistical Analysis using Excel
 - Closed (Friday, March 30): Grades Posted
- Lab 9: Data analysis with Excel (linear regression)
 - Closed (Friday, April 6): Grades Posted
- Lab 10: Simple programming in Excel and Measuring Uncertainty
 - April 12 and 13, Due April 20



■ Assignments

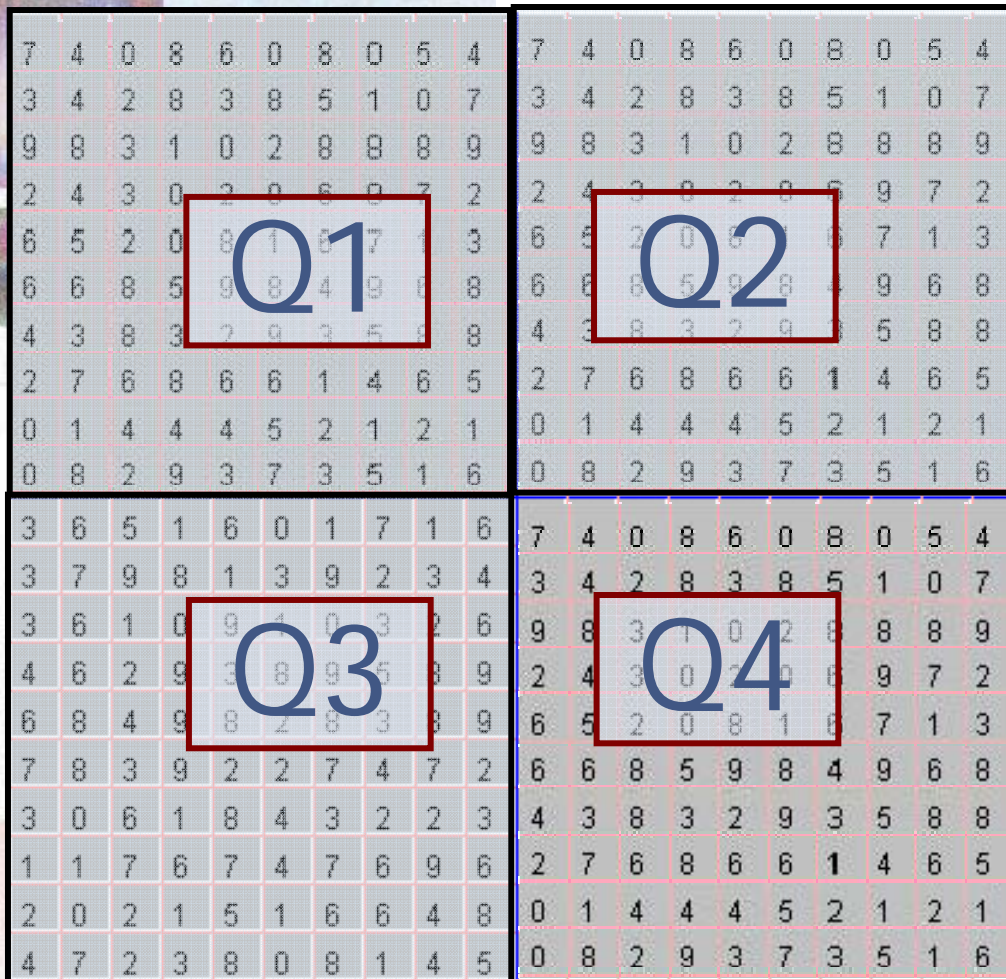
■ Individual

- First installment
 - Closed: February 9: Grades Posted
- Second Installment
 - Past: March 2: Grades Posted
- Third installment
 - Past: Grades Posted
- Fourth Installment
 - Presented April 10th, Due April 20th

■ Group

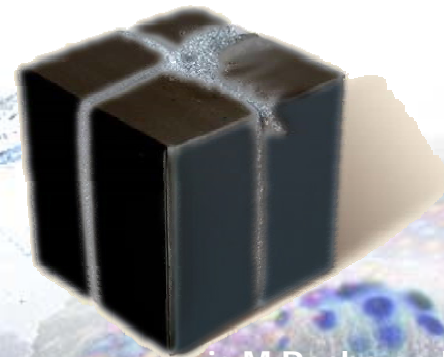
- First Installment
 - Past: March 9th, graded
- Second Installment
 - Past: April 6th Graded
- Third Installment
 - Presented Thursday, April 12; Due Friday, April 27

Individual Assignment – Part IV



Cycles = 1

- Step by step analysis of “dying” squares
 - 4th Installment
 - Presented: April 10th
 - Due: April 20th
- Use inductive and deductive reasoning
 - To uncover the algorithm in each quadrant
 - Build from inductive knowledge accumulated so far



Summary of Black Box

■ Quadrant 1

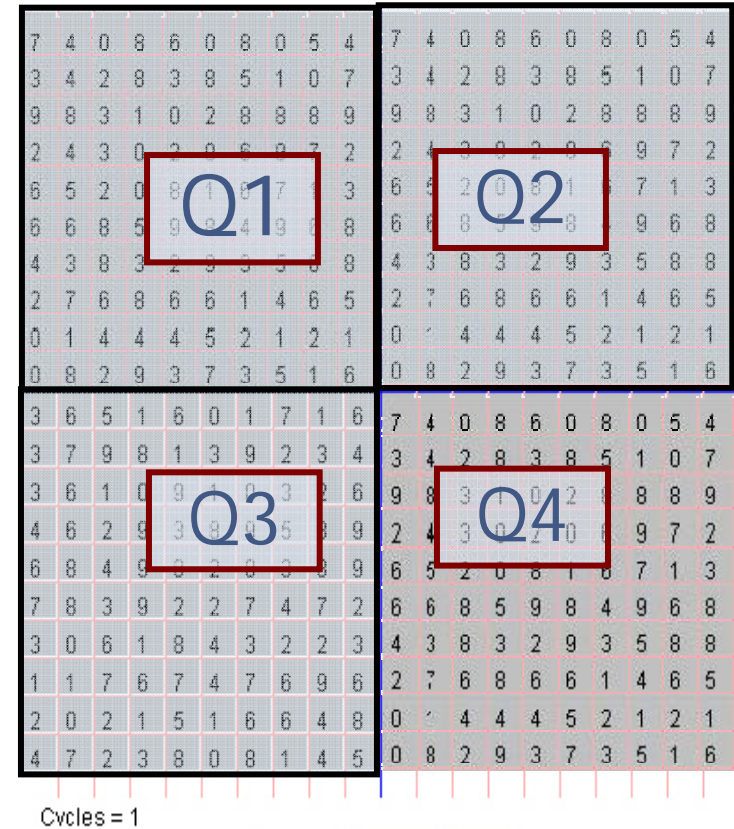
■ At the random initial state

- All numbers have equal probability of being initially present
- But the probability of changes are different

■ In Any State

- Any number changes depending on its neighbors
- It 'gravitates' towards the smallest number that it 'sees' most often.
- Odd and Even numbers do not show different behavior

■ What is the Algorithm?



1 Restart Go

Summary of Black Box

■ Quadrant 3

■ At the random initial state

- All numbers have equal probability of being initially present
- But the probability of changes are different

■ In Any State

- 0 can only change to 0
- 5 can only change to 5 or 0
- Even digits always change to even digits
- Odd digits could change to any other digit

■ What is the Algorithm?

	$n(i)$	$p(i)$
0	27	0.27
1	4	0.04
2	12	0.12
3	4	0.04
4	12	0.12
5	9	0.09
6	12	0.12
7	4	0.04
8	12	0.12
9	4	0.04

1. $0 \rightarrow 0$
2. $\{5\} \rightarrow \{0, 5\}$
3. $\{2, 4, 6, 8\} \rightarrow \{0, 2, 4, 6, 8\}$
4. $\{1, 3, 7, 9\} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$

Summary of Black Box

■ Quadrant 2

■ At the random initial state

- All numbers have equal probability of being initially present
- But the probability of changes are different

■ In Any State

- 0 can only change to 0
- 5 can only change to 5 or 0
- Even digits always change to even digits
- Odd digits could change to any other digit

■ What is the Algorithm?

1. $0 \rightarrow 0$
2. $\{5\} \rightarrow \{0, 5\}$
3. $\{2, 4, 6, 8\} \rightarrow \{0, 2, 4, 6, 8\}$
4. $\{1, 3, 7, 9\} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$

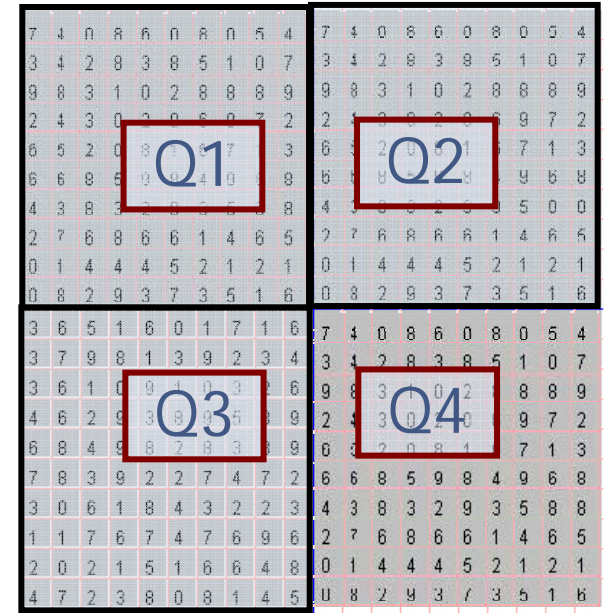
Possible Operations Q2 and Q3

Operator	Meaning	Excel	Example
()	Brackets, grouping	()	$y = (a + b) * (c + d)$
*	Multiplication	*	$i = j * k$
+	Add	+	$i = i + 1$
-	Subtract	-	$i = j - 3.2$
/	Real division	/	$i = 8 / 5 = 1.6$
div	Integer division	Quotient (a,b)	$i = 8 / 5 = 1$
Mod, %	remainder	Mod (a, b)	$i = 8 \text{ mod } 5 = 3$
ROUND	Rounds	ROUND (a, d)	$i = \text{ROUND}(3.67, 0) = 4$
INT	Integer Part	INT	$i = \text{INT}(3.67) = 3$
rand	Random number	Rand() RandBetween(a,b)	$i = \text{rand}(n)$

Tip for Individual Assignment

■ Quadrant Q

- There are 100 cells in each 10x10 quadrant
 - $C = 1..100$
- Each cell can take one of 10 colors
 - $V(C) = 0..9$
 - is the value of the cell
 - This is the state cell C is in
- Random initialization of quadrant Q at cycle 1
 - For $c=1$ to 100 do
 - $V(C) \leftarrow \text{randbetween}(0,9)$ {random number 0 to 9}
 - EndFor
 - Cycle $\leftarrow 1$
- Run for Number of cycles
 - $n \leftarrow$ Input dialog
 - For $k=1$ to n do
 - Cycle \leftarrow cycle+1
 - {Pick random cell}
 - $C \leftarrow \text{randbetween}(1,100)$
 - {Update the value of the cell (NOT THE REAL THING)}
 - $V(C) \leftarrow ((V(C) * \text{randbetween}(0,9)) \text{ div } 2) - 5 * X$
 - EndFor
- X may be a hidden variable
 - $X \leftarrow ???$



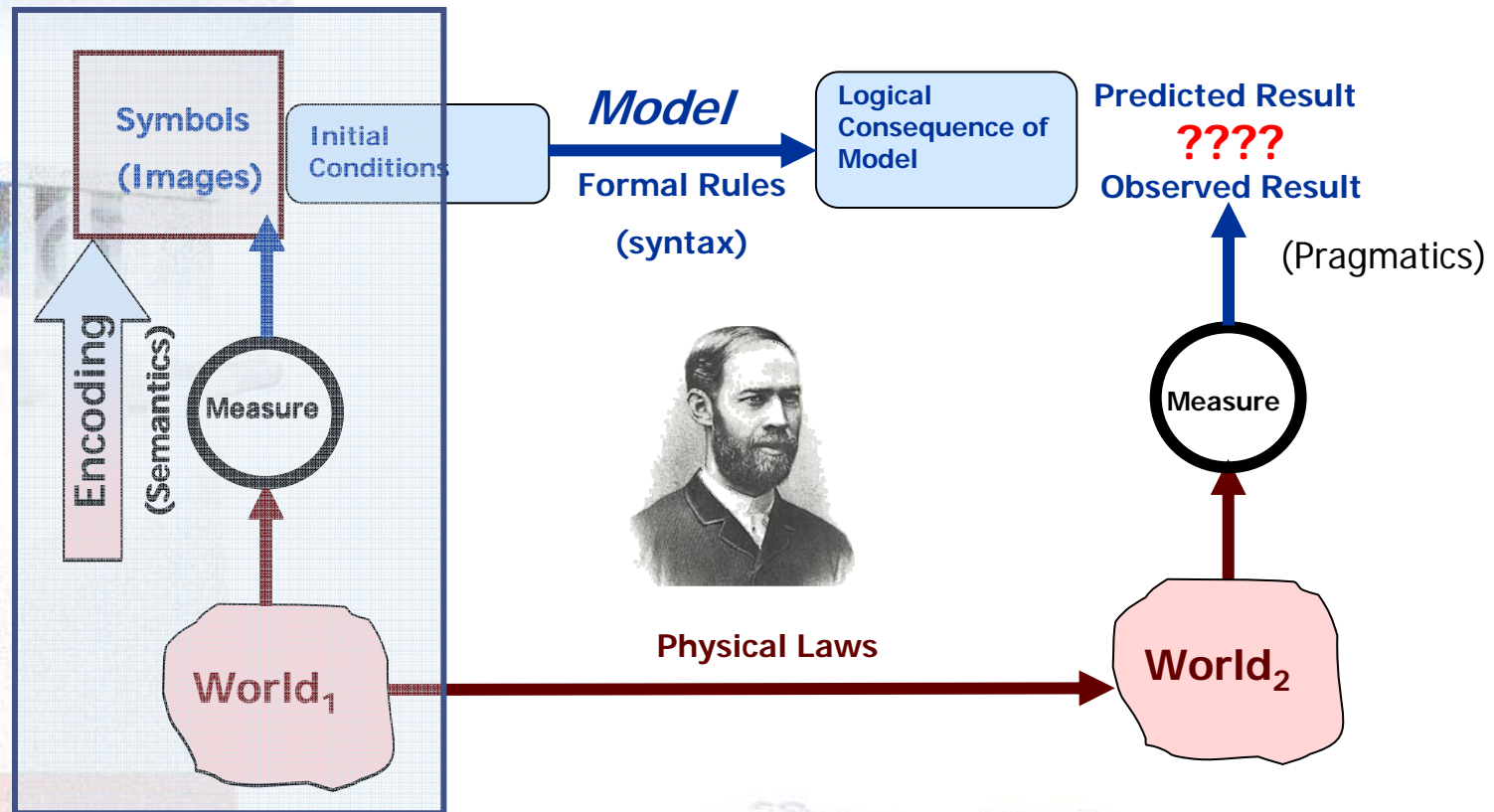
Cycles = 1

Restart Go



The Modeling Relation

Hertz' Modeling Paradigm



- Organizing Data
 - After encoding
 - Modern Problems Require large storage capabilities

The Entity-Relationship Model

- **Conceptual *Data Model***
 - A kind of “pseudocode” for *models of data storage*
- **What should we consider?**
 - What are the interesting entities and relationships in our model of reality?
 - What information about these entities and relationships do we need to store?
 - What are the reality constraints and rules that must hold?



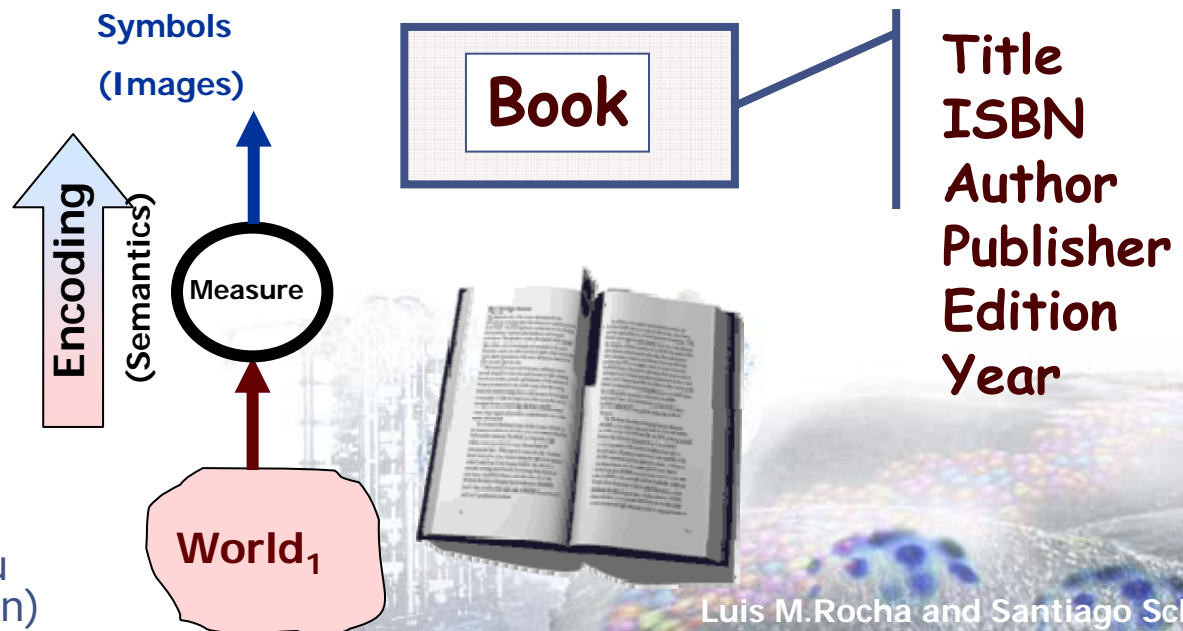
Peter Chen
(1976)

Adapted from Yuqing Melanie Wu
(I308: Information Representation)

Luis M.Rocha and Santiago Schnell

Entities in Data Modeling

- **Objects, people, places**
 - Basically *a noun*: a discrete object
 - Choose a meaningful name
 - Represented by a rectangle
 - Attributes
 - Describe the proprieties of an entity

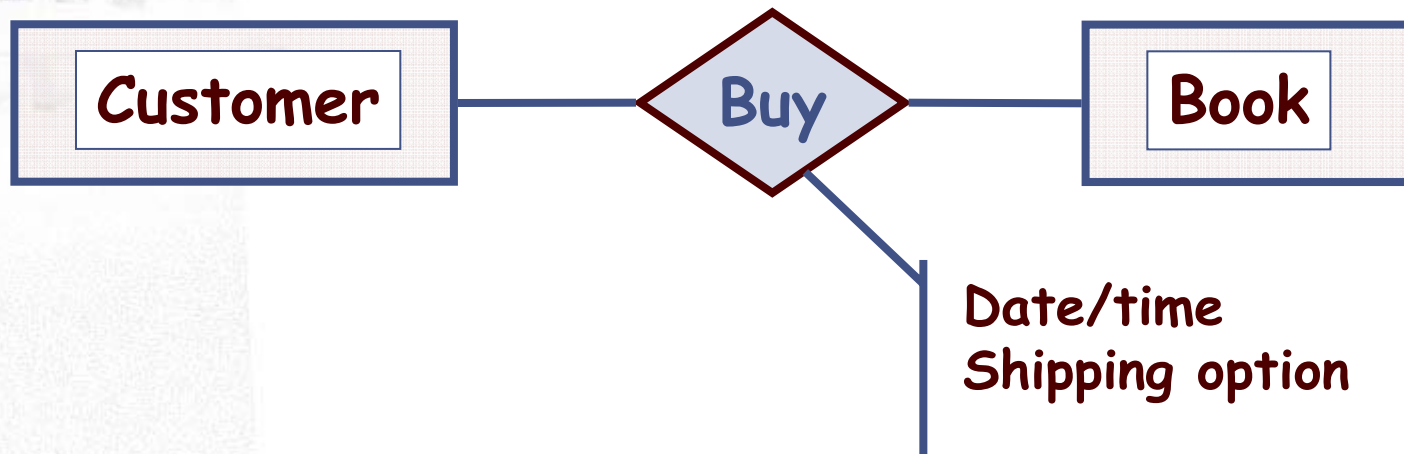


Adapted from Yuqing Melanie Wu
(I308: Information Representation)

Luis M.Rocha and Santiago Schnell

Relationships in Data Modeling

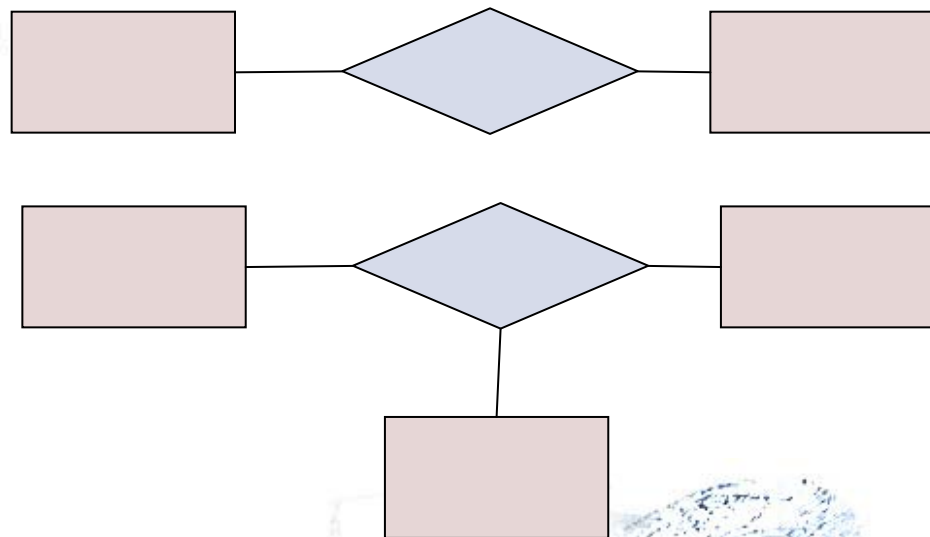
- **Relationship:**
 - An association among two or more entities.
 - Verbs



- **Attributes also describe relationships**

Arity of Relationships

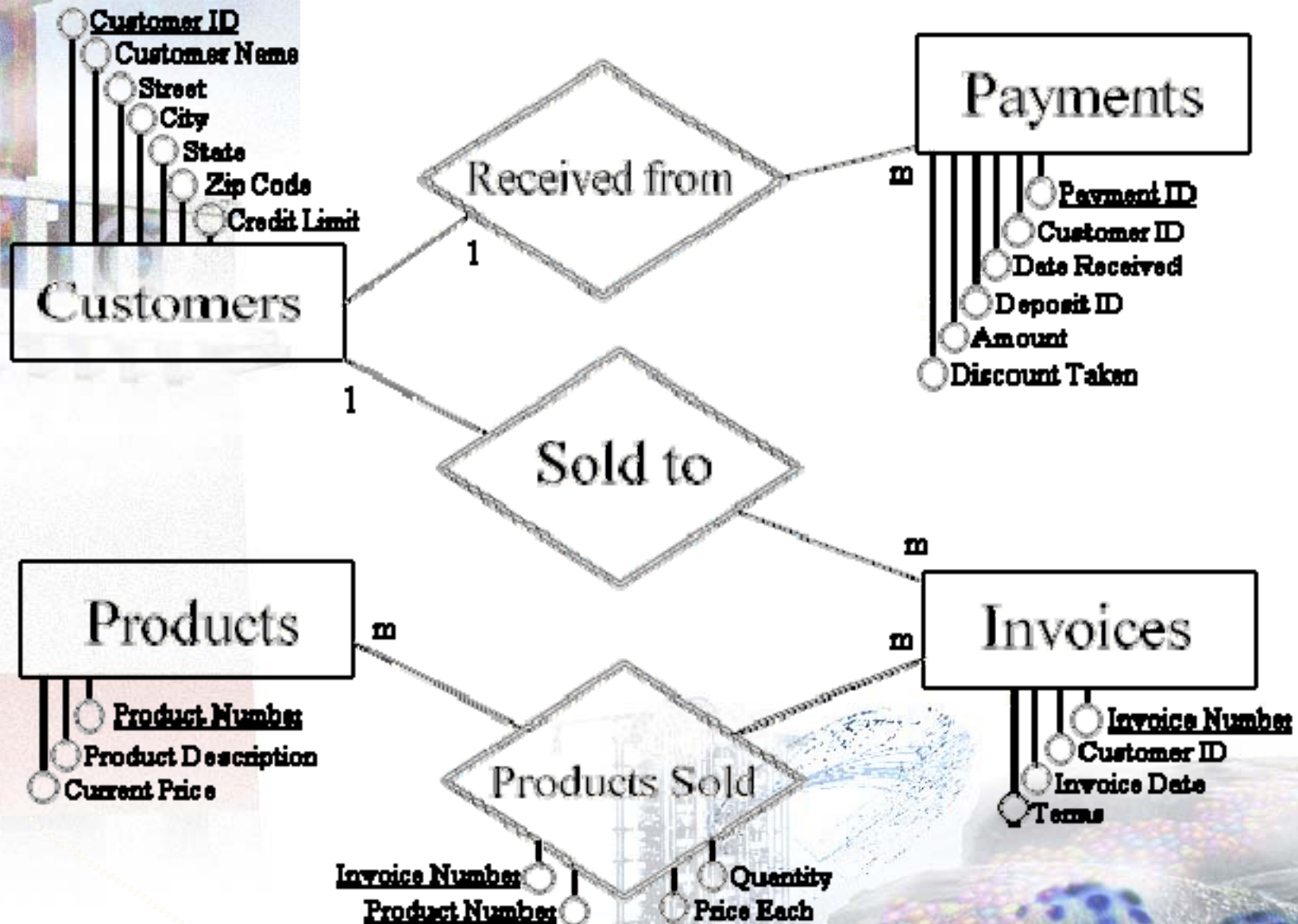
- The number of entities participate in a relationship
 - Binary, ternary, N-ary



Adapted from Yuqing Melanie Wu
(I308: Information Representation)

Luis M.Rocha and Santiago Schnell

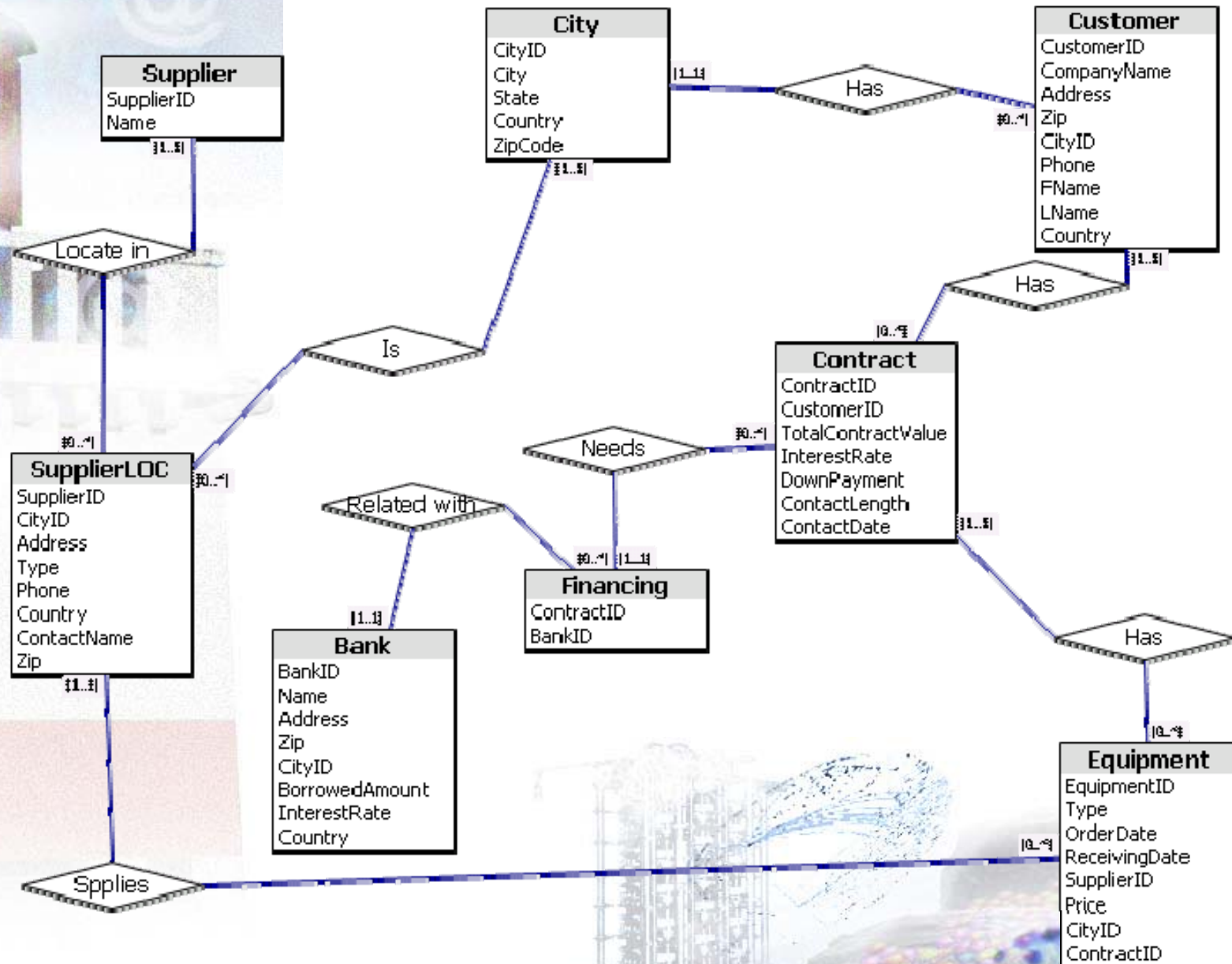
ER Data Modeling Example



From Carol Brown

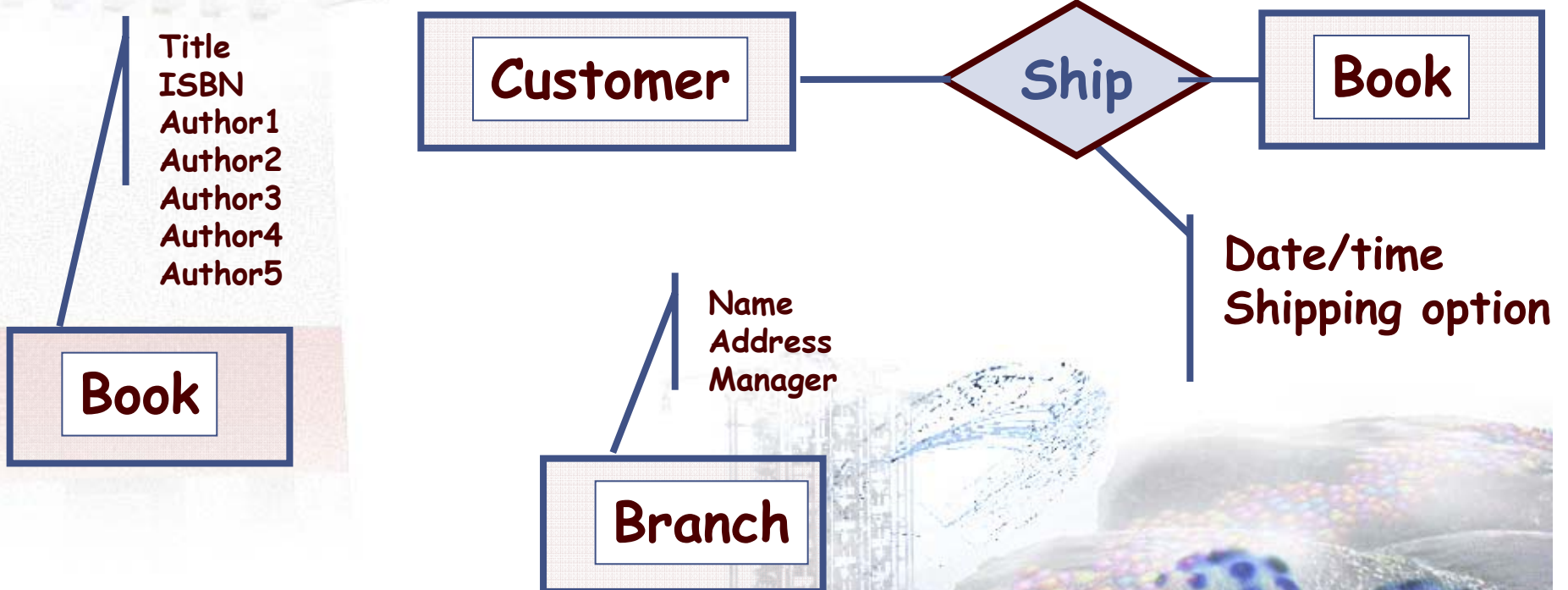
Luis M.Rocha and Santiago Schnell

ER Data Modeling Example 2



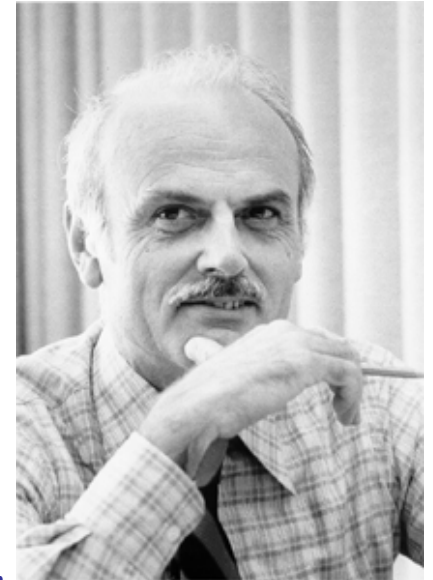
Try this at home...

- How to represent the following?
 - A book can have no more than 5 authors
 - A customer has to specify the shipping option
 - Each branch has only one manager.

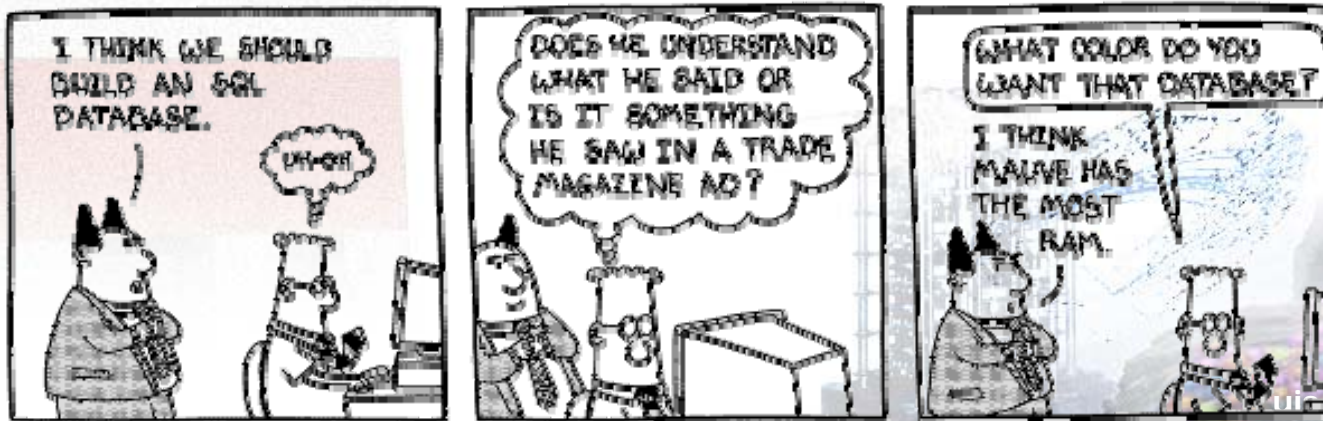


The Relational Database Model

- Relational database management system (RDBMS)
 - Most popular commercial database type.
 - a data model based on *logic* and *set theory*.
- invented by Ted Codd in 1970
 - Oxford, IBM, U. Michigan, IBM
- System R
 - IBM's San Jose research center
 - Structured English Query Language ("SEQUEL")
 - Data Manipulation Language (DML)
 - SEQUEL was later condensed to SQL due to a trademark dispute
 - In 1979, Relational Software, Inc. (now Oracle Corporation) introduced the first commercially available implementation of SQL

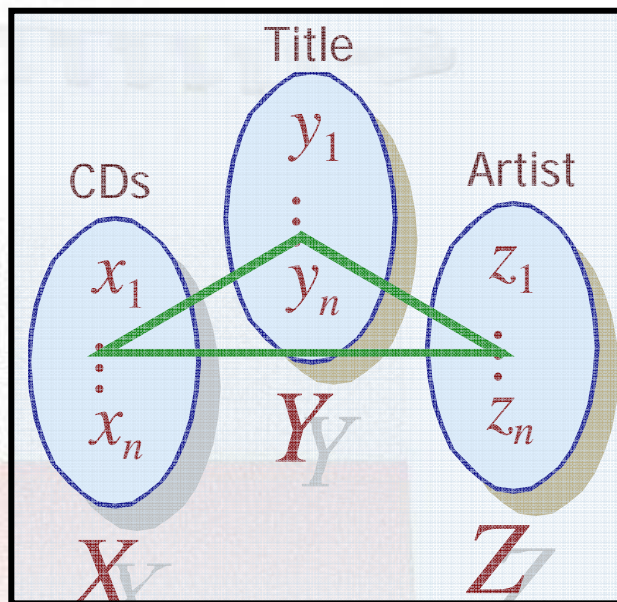


Ted Codd

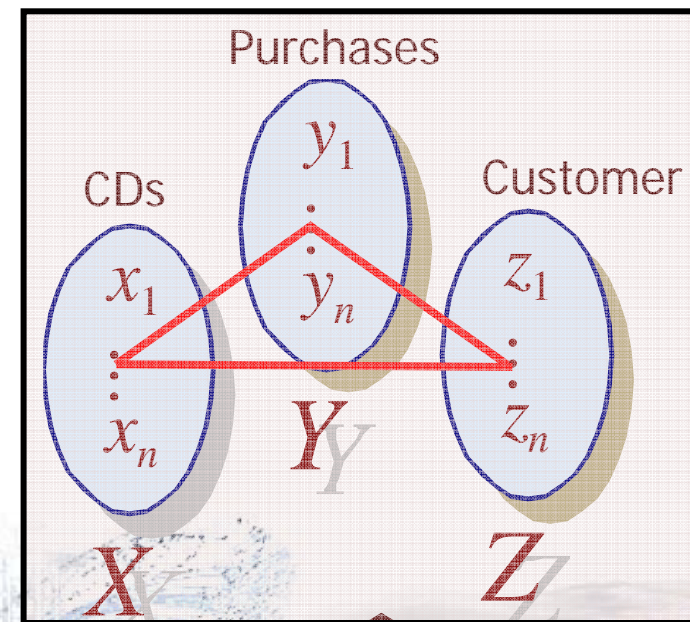


The Relational Model

- All data are represented as mathematical relations
 - Represent the presence of association, interaction or interconnectedness between the elements of two or more sets.
 - A relation associates the elements of 2 or more sets
 - Set of books with sets of attributes (*entities*)
 - Set of purchases with sets of attributes (*relationships*)
 - Tables store relations



Title
Artist

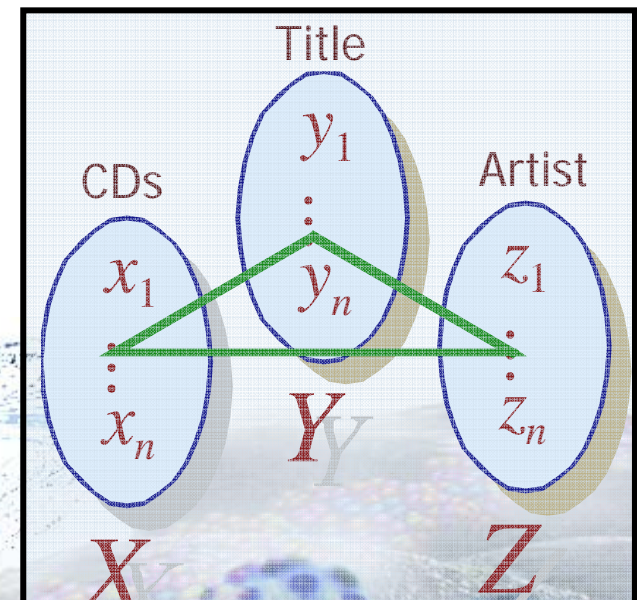


The Relational Database Model

- A relational database is a collection of tables
 - 2-dimensional
- Each table has a unique name in the database.
- Tables define Relations
 - Columns (number of sets)
 - Attributes plus key (*primary set*)
 - Row (number of relation instances)
 - A table is a **set** of rows: tuples

CDs

ID	Title	Artist
3592	Yes I am a Witch	Yoko Ono
2678	Big	Macy Gray
0623	Sound of Silver	LCD Soundsystem
0321	Welcome to Planet Sexor	Tiga
8854	Transparent Things	Fujiya & Miyagi



Example

Relation (table)

Attributes (columns)

Customer

Tuples
(rows)

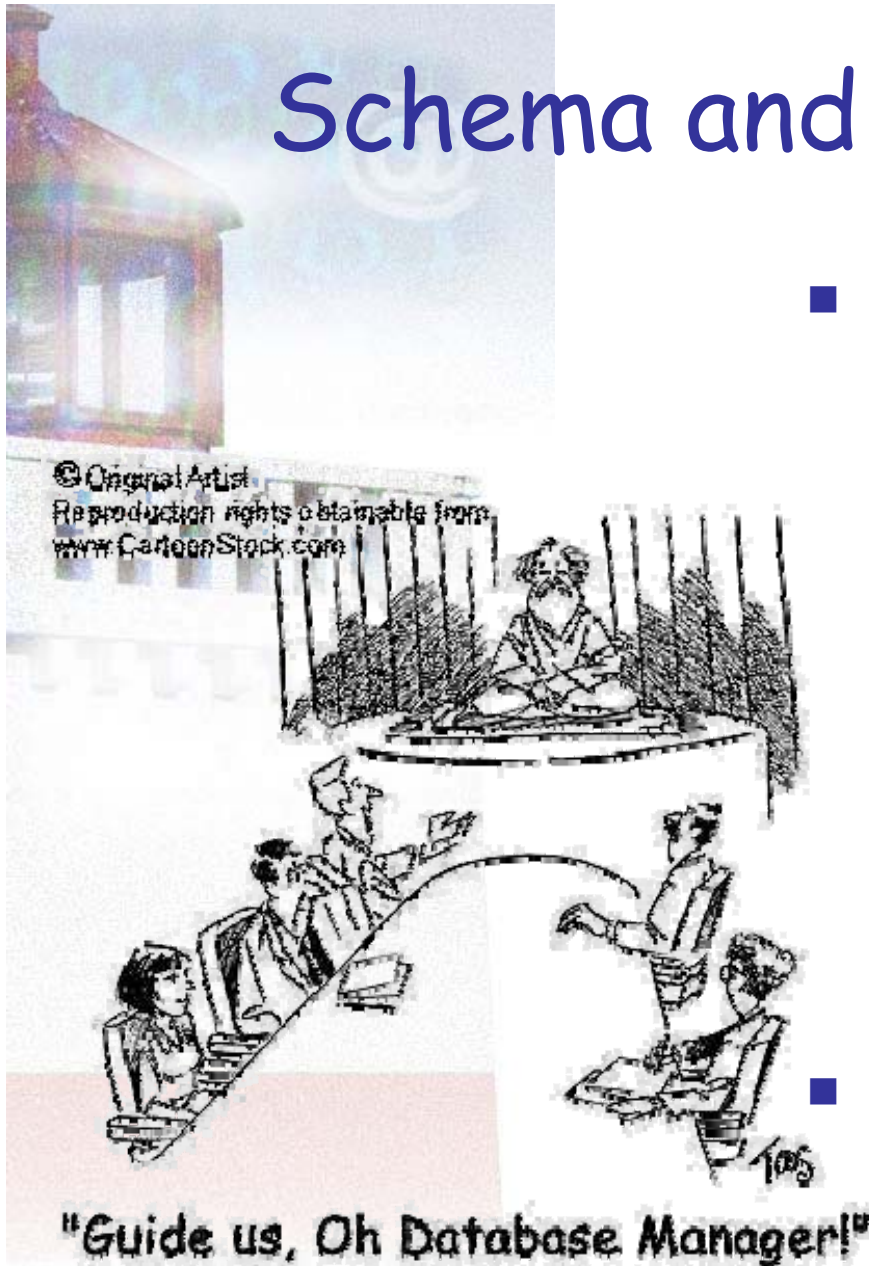
Phone	Name	Address
812-123-4567	Tom	408 3 rd st. Bloomington, IN
812-304-2378	Bill	#113, Redbud Hall, Bloomington, IN
812-856-1190	Kate	1205, Maritime ct. Bloomington, IN
812-754-9567	Mary	#901 10 th St. Bloomington, IN
317-897-4536	Pam	2400 Rd135, Greenwood, IN
812-906-2486	Jeff	#208 Union Ave. Bloomington, IN

Degree = 3

Cardinality = 6

Schema and Instance

- Database schema
 - Metadata or Model
 - The logical design of a database
 - E.g. using the *entity-relationship model*
 - *Entity* → *Table*
 - *Attribute* → *Columns*
 - *Relationship* → *Table*
 - Specifies names of tables/relations (*entities and relationships*), plus names and types of each column (*attributes*)
 - Database instance
 - A snapshot of the data in the database at a given instant in time.



Adapted from Yuqing Melanie Wu
(I308: Information Representation)

Luis M.Rocha and Santiago Schnell

Primary Key

Customer

Key

Phone	Name	Address
812-123-4567	Tom	408 3 rd st. Bloomington, IN
812-304-2378	Bill	#113, Redbud Hall, Bloomington, IN
.....

Customer(Phone, Name, Address)

- The identifying labels for the elements of the primary set of a table
- Every instance (row) in the database must have a distinct primary key
- Every instance in the database must have a particular (non-null) value for the primary key.

Book

Key

Customer

Key

ISBN	Title	Publisher	Phone	Name	Address
12345	Java	MIT press	812-123-4567	Tom
49082	Snow White	812-304-2378	Bill
72936	Honeymoon	812-856-1190	Kate

Book(ISBN, Title, Publisher)

Customer(Phone, Name, Address)

Sale

Key

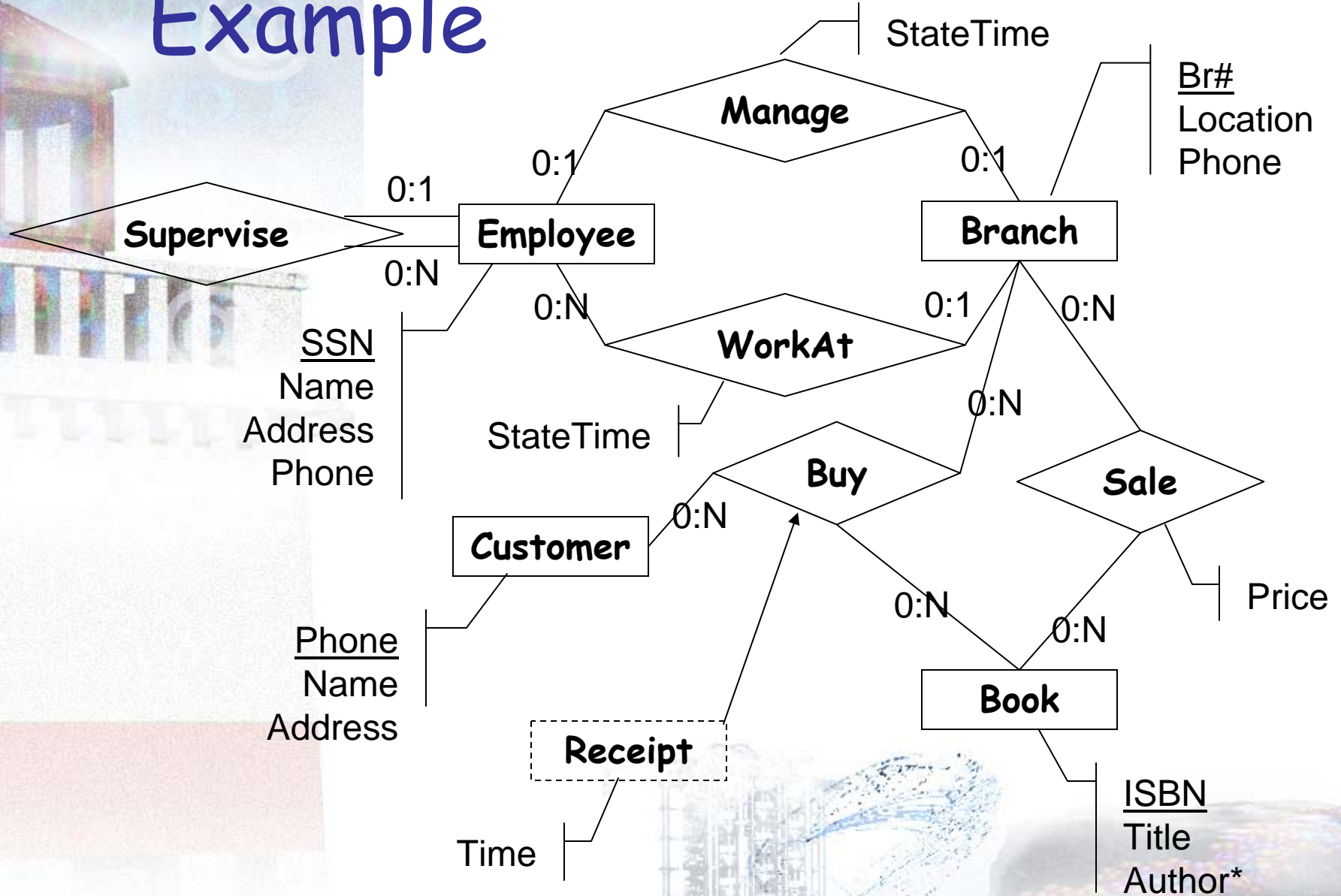
Come from primary keys in entities

ISBN	Phone	Price	Time
12345	812-123-4567	\$20	Feb 2, 05
49082	812-123-4567	\$25	Dec 20, 04
12345	812-856-1190	\$19

Primary keys in Relationships

Sale (ISBN, Phone, time, price)

Example



From Yuqing Melanie Wu (I308: Information Representation)

Luis M.Rocha and Santiago Schnell

Structured Query Language (SQL)

- The most popular computer language used to create, modify and retrieve data from relational database management systems. (Wikipedia)
- Three subsets of SQL
 - Data Definition Language (DDL)
 - Data Manipulation Language (DML)
 - Data Control Language (DCL) (for authorization)

Data Definition Language

- Used to create, alter, and delete databases and tables.
- Statements
 - Create Table
 - `CREATE TABLE table_name (column_name1 data_type primary key, column_name2 data_type);`
 - Some other operations
 - "alter" and "drop"

Data Manipulation Language

- Used to retrieve, insert, delete and update data in a database
- Statements
 - Select
 - Selects rows (records) according to attribute criteria
 - E.g. Select CDs published in YEAR=x
 - Some other operations
 - "insert", "update", "delete", and "truncate"

Select Statement

- **Select**
 - Selects rows (records) according to attribute criteria
 - E.g. papers published in YEAR=x
 - **SELECT** * FROM *list-of-relations* WHERE *condition*
 - **SELECT** * FROM *CITATION_TABLE* WHERE *PUBLISHED_YEAR='1995'*;
 - * Denotes ALL
 - **SELECT** * FROM *T*;
 - Returns all elements of all the rows of the table T

	MUID	Journal	Volume	Pages	Year
Paper1					
Paper2					
Paper3					
Paper4					

Select

Projection Operation

- Project

- Extracts columns

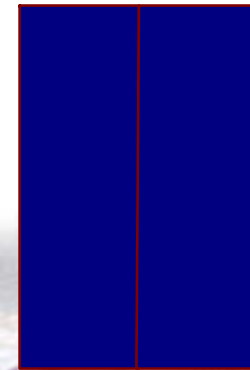
- E.g. projects a set of papers into a reduced set of attributes.

- **SELECT C1,C7 FROM T;**





Project



Join Operation

■ Join

- Merges records that contain matching values for specified attributes
 - given a key value join records from both tables
- **SELECT** * FROM *employee, department*;
- **SELECT** * FROM *citation-table, author-table* WHERE *citation-table.MUID = author-table.MUID*;

	MUID	Journal	Volume	Pages	Year
Paper1					
Paper2					
Paper3					
Paper4					

	MUID	Author
Author1,1		
Author1,2		
Author1,3		



MUID	Journal	Volume	Pages	Year	Author

Betty

By Delaney & Rasmussen



© 2001 by NEA, Inc. www.comix.com 11-78



Group Assignment

■ Third Installment

- Given any text such as the *library of babylon* or *Funes, the memorious*
 - Create a **database model** and a **relational database instance** using *Microsoft Access* to store the data and conclusions from previous installments
 - Use the entity-relationship model
 - Examples of items that should appear
 - Title, author, language, publication date
 - Frequency/probability of each letter
 - Conditional probabilities for letters 'e' and 'u' (as produced in installment 2)
 - Positively and negatively dependent letters
 - Use at least 4 texts
- Due on April 27th, 2005
- Upload to Oncourse





Next Class!

- Topics of next classes
 - Databases and SQL
 - Individual Assignment
 - Review
- Readings for Next week
 - @ *infoport*
 - course package
- No More Labs!!!!!!!