

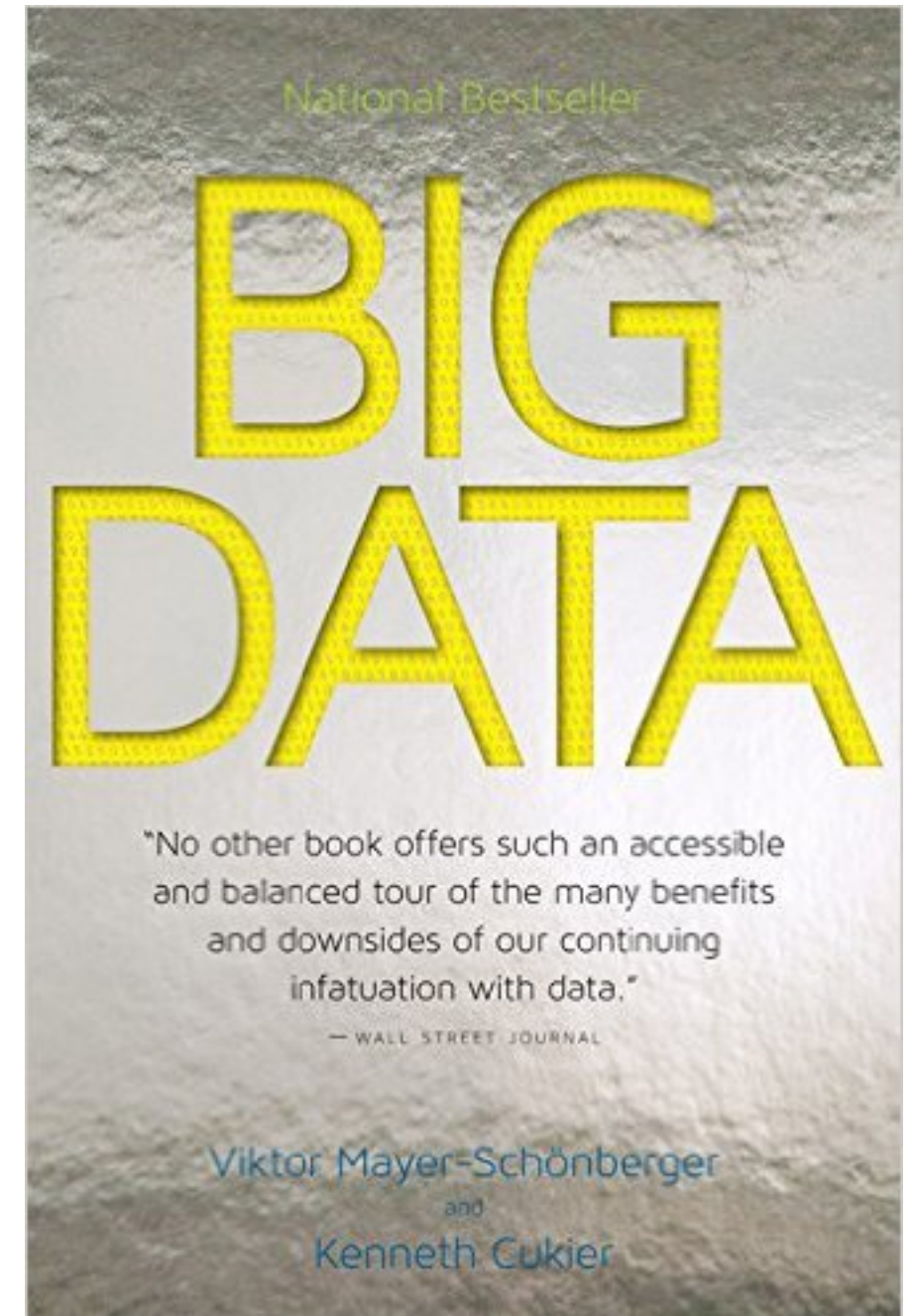
# Lecture 3: Big

Core 109S IDWT?, Spring 2017

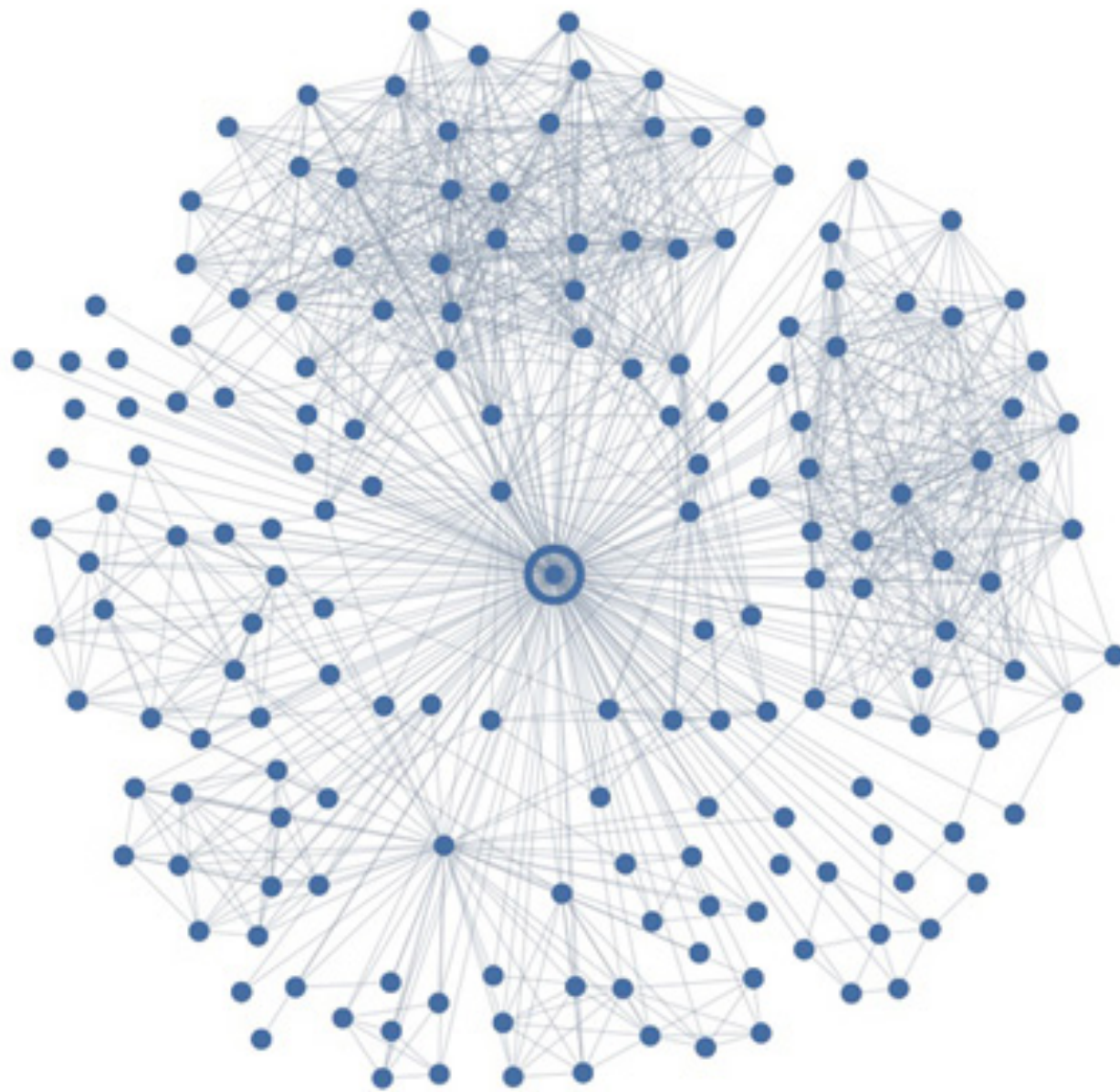
Michael Hay

# Discussion

- The book describes three shifts in mindset in Ch. 2-5.
- What is the shift described in Ch. 2 "More"? (Find good quotes to support your answers!)
- Example from chapter?
- What did you find interesting?
- Connection to Facebook study?



# Social networks



*"Thus, instead of embeddedness, we propose that the link between an individual  $u$  and his or her partner  $v$  should display a 'dispersed' structure: the mutual neighbors of  $u$  and  $v$  are not well-connected to one another, and hence  $u$  and  $v$  act jointly as the only intermediaries between these different parts of the network."*

<https://arxiv.org/pdf/1310.6753v1.pdf>

# Databases and SQL

- Over the next couple of weeks, we'll focus on relational databases and the *structured query language* (SQL).
- Why?
  - Practical skill: go beyond the spreadsheet!
  - Formal programming languages are an essential ingredient to computer science
  - Want you to experience challenge of translating from "plain English" to *formal specification*

# Invention of relational databases

- Invented by E.F. Codd in 1970
- Radical idea: describe data in a logical way that is *independent* of how it is actually stored on disk

“Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation)... Activities of users at terminals and most application programs should remain unaffected when the internal representation of data is changed...”

- Invented in 1970 yet did not see widespread adoption for another 10 years... why?

## Information Retrieval

# A Relational Model of Data for Large Shared Data Banks

E. F. CODD  
*IBM Research Laboratory, San Jose, California*

Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation). A prompting service which supplies such information is not a satisfactory solution. Activities of users at terminals and most application programs should remain

The relational view (or model) in Section 1 appears to be superior in a graph or network model [3, 4] to preexisting inferential systems. It provides a method with its natural structure only—without imposing any additional structure for many purposes. Accordingly, it provides a data language which will yield maximum between programs on the one hand and action and organization of data on the

A further advantage of the relational model is that it provides a sound basis for treating derivation and consistency of relations—these are the two main properties of a relational model on the other

# moving to external most internal

## INCONSISTENCIES IN PR

toward the goal of data independent facilitate changing certain characterisation stored in a data bank. However, data representation characteristics *without logically impairing some* are still quite limited. Further, the mod users interact is still cluttered with errors, particularly in regard to the sections of data (as opposed to individual principal kinds of data dependence to be removed are: ordering dependence, and access path dependence. These dependencies are not clearly separa

**KEY WORDS AND PHRASES:** data bank, data base, data structure, data organization, hierarchies of data, networks of data, relations, derivability, redundancy, consistency, composition, join, retrieval language, predicate calculus, security, data integrity

CR CATEGORIES: 3.70, 3.73, 3.75, 4.20, 4.22, 4.29

## 1. Relational Model and Normal Form

## 1.1. INTRODUCTION



# Relational Database: Definitions

- *Relational database*: a set of *relations*
- *Relation*: made up of 2 parts:
  - *Schema* : name of relation, plus name and type/domain of each column.
  - *Instance* : a *table*, with rows and columns.

Students(*sid*: string, *name*: string, *login*: string,  
*age*: integer, *gpa*: real).

# Types

Types will be important later: certain operations behave differently, depending on the type

- **Strings:** a string is a sequence of characters. Used to store names, addresses, documents, etc.
  - char (a string of characters)
  - varchar (variable length character string)
- **Int:** stores an integer {..., -2, -1, 0, 1, 2, ...}
- **Real:** stores (approximation of) real-valued number (3.14159...)
  - Also known as float, double

# Relational instance: a table

Students

column,  
attribute,  
field

sid	name	login	age	gpa
53666	Jones	jones@cs	18	3.4
53688	Smith	smith@eecs	18	3.2
53650	Smith	smith@math	19	3.8

row, tuple

Attribute value

The diagram illustrates a relational instance as a table. The table has five columns: 'sid', 'name', 'login', 'age', and 'gpa'. The first three columns are highlighted with a light gray background. The table contains three rows of data. The third row is highlighted with a blue border. The 'name' cell in the third row is further highlighted with a dark blue border. An arrow points from the text 'Attribute value' to this cell. A red box highlights the 'age' and 'gpa' columns. The text 'column, attribute, field' is positioned above this box. The text 'row, tuple' is positioned to the left of the third row. The title 'Students' is centered above the table.



# Example Database

STUDENT

sid	name
1	Jill
2	Bo
3	Maya

Takes

sid	cid
1	445
1	483
3	435

COURSE

cid	title	sem
445	DB	F12
483	AI	S14
435	Arch	F12

PROFESSOR

fid	name
1	Diao
2	Saul
8	Weems

Teaches

fid	cid
1	445
2	483
8	435

# Real Schemas Are Complex

Wikipedia Schema

# Relation instance: dimensions

- Cardinality
  - number of rows
- Arity/Degree
  - number of attributes
  - unary (1), binary (2), ternary (3), ...

# Example Database

STUDENT

sid	name
1	Jill
2	F
3	

binary table w/  
cardinality 3

Takes

sid	cid
1	445
1	483
3	435

COURSE

cid	title	sem
445	DP	F12
483		14
435		F12

ternary table w/  
cardinality 3

PROFESSOR

fid	name
	Diao
	Saul

binary table w/  
cardinality 2

Teaches

fid	cid
1	445
2	483
8	435

# Keys

- A **key** is a set of one or more attributes whose values are *guaranteed* to identify tuples in the relation uniquely
- Book distinguishes between superkey, candidate key, primary key, foreign key
  - You don't need to know this level of detail.
  - "key"  $\approx$  record identifier
  - "foreign key"  $\approx$  reference to a record of some other (foreign) table.

# Example Database

STUDENT

sid	name
1	Jill
2	Bo
3	Maya

Takes

sid	cid
1	445
1	483
3	435

COURSE

cid	title	sem
445	DB	F12
483	AI	S14
435	Arch	F12

PROFESSOR

fid	name
1	Diao
2	Saul
8	Weems

Teaches

fid	cid
1	445
2	483
8	435

What are the keys  
of these relations?

Are there any  
foreign keys?