


# Introduction

Introduction to Databases  
CompSci 316 Spring 2020



1


## Welcome to

### CompSci 316: Introduction to Database Systems!! Spring 2020

2


### About us: instructor

- Instructor: [Sudeepa Roy](#)
  - At Duke CS since Fall 2015
  - Member of “Duke Database Devils” a.k.a. the database research group
  - PhD. UPenn, Postdoc: U. of Washington
  - Research interests:
    - “data”
    - data management, database theory, data analysis, data science, causality and explanations, uncertain data, data provenance, crowdsourcing, ....



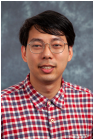
3

### Meet your grad TAs




[Tiangang Chen](#)

- Duke CS MS student
- Interested in computer systems
- Runs a half-marathon every year!



[Xiangchen Shen](#)

- Duke CS MS student
- Interested in data science and machine learning, currently working on image processing
- Loves cats!




[Yuchao Tao](#)

- Duke CS PhD student
- Interested in databases and privacy research
- Enjoys cooking!

\* All CompSci 516 veterans


4

### Meet your UTAs




[David Chen](#)

- Duke CS major
- Interested in applying computational models and algorithms to biological systems
- Loves to downhill ski!



[Jane Li](#)

- Duke CS major, minor in Visual Media studies
- Interested in UI/UX, front-end development, and project management
- Has two dogs and a hamster!



[Runxin \(Rebecca\) Wang](#)

- Duke CS major
- Interested in coding, machine learning, and theory
- Loves hiking and bubble tea!

\* All CompSci 316 veterans

5

### What are the goals of this course?

- Learn about “databases” or data management

6

## Why do we care about data? (easy)

How big data can help find new mineral deposits

... The three years of gathering and analyzing data culminated in what U.S. Sealing calls their "Bio Weather Playbook," a body of critical information about each of the seven courses only available to the U.S. team...

— FiveThirtyEight, "Will Data Help U.S. Ling Get Back On The Olympic Podium?" Aug 15, 2016

**Data = Money Information Power Fun in Science, Business, Politics, Security Sports, Education, ....**

7

Wait.. don't we need to take a Machine Learning or Stat course for those things?

Yes, but..

Pic: <https://www.technobuffalo.com/sites/technobuffalo.com/files/styles/large/public/wp/2012/05/confused-student.jpg>

8

... we also need to manage this (huge or not-so-huge) data!

9

### Also think about building a new App or website based on data from scratch

- E.g., your own version of book purchase platform (like a mini-Amazon)
- Large data! (think about all books in the world or even in English)

**• How do we start?**

\* You are going to do something similar in the course project!

10

## Who are the key people?

11

## Who are the key people?

12

What should the user be able to do?<sup>13</sup>

13

What should the user be able to do?<sup>14</sup>

14

What should the platform do?<sup>15</sup>

15

What should the platform do?<sup>16</sup>

16


What are the desired and necessary properties of the platform?<sup>17</sup>

17

What are the desired and necessary properties of the platform?<sup>18</sup>

18

That was the design phase  
(a basic one though)



How about C++, Java, or Python?  
On data stored in large files

https://it.wp.com/dynamiclandscapes.vita-learn.org/wp-content/uploads/2019/05/lets-code.jpg?resize=768%2C432&ssl=1

19

Sounds simple!

```
James Morgan#Durham, NC
....
A tale of two cities#Charles Dickens#3.50#7
To Kill a Mockingbird#Harper Lee#7.20#1
Les Miserables#Victor Hugo#12.80#2
....
```

- Text files – for books, customer, ...
- Books listed with title, author, price, and no. of copies
- Fields separated by #'s

20

Query by programming

```
James Morgan#Durham, NC
....
A tale of two cities#Charles Dickens#3.50#7
To Kill a Mockingbird#Harper Lee#7.20#1
Les Miserables#Victor Hugo#12.80#2
....
```

- James Morgan wants to buy "To Kill a Mockingbird"


What if he changes the "query" and wants to buy a book by Victor Hugo?

21


Revisit: What are the desired and necessary properties of the platform?

22

Solution?



- DBMS = Database Management System



23

A DBMS takes care of all of the following (and more):

In an easy-to-code, efficient, and robust way

Optimization

Recovery

Index

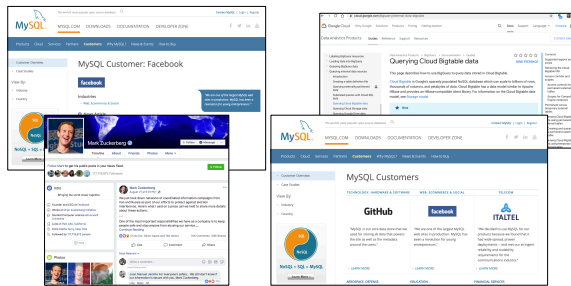
Declarative

Consistency

\* We will learn these in the course!

24

## DBMS helps the big ones!



Note: Not always the “standard” DBMS (called Relational DBMS), but we need to know pros and cons of all alternatives

25

## CompSci 316 gives an intro to DBMS

- How can a user use a DBMS (programmer’s/designer’s perspective)
  - Run queries, update data (SQL, Relational Algebra)
  - Design a good database (ER diagram, normalization)
  - Use different types of data (Relational, XML, JSON)
- How does a DBMS work (system’s or admin’s perspective)
  - Storage, index
  - Query processing, join algorithms, query optimizations
  - Transactions: recovery and concurrency control
- Glimpse of advance topics and other DBMS
  - NOSQL, Spark (big data)
  - Data mining
- Hands-on experience in class projects by building an end-to-end website or an app that runs on a database

26

## Misc. course info

- Website: <https://www2.cs.duke.edu/courses/spring20/compsci316/>
  - Course info; tentative schedule and reference sections in the book; lecture slides, assignments, help docs, ...
- Book: *Database Systems: The Complete Book*, by H. Garcia-Molina, J. D. Ullman, and J. Widom. 2<sup>nd</sup> Ed.
- **Programming**: VM required, need significant programming on different platforms and languages
- Prerequisite: **CompSci 230** (will need basic understanding of discrete maths, data structure, and algorithms) - or talk to us
- Q&A on **Piazza**
- Grades, sample solutions on **Sakai**
- Submissions on **Gradescope** and **Gradiance**
- Watch your email for announcements

27

## Important: Grading

Absolute but adjustable grading

Guarantees:

[90%, 100%] A- / A / A+

[80%, 90%) B- / B / B+

[70%, 80%) C- / C / C+

[60%, 70%) D

Class topper gets A+

- Scale will not go upwards but can get downwards (e.g., based on the class performance in the exams)

- We will give you a feedback on your approximate standing after the midterm.

28

## Duke Community Standard

- See course website for link
- Group discussion for assignments is okay (and encouraged), but
  - Acknowledge any help you receive from others
  - Make sure you “own” your solution
- All suspected cases of violation will be aggressively pursued

29

## Course load

- (See course webpage for full details)
- **Weekly (short) homework assignments (25%)**
  - Each homework has same weight
  - Released on Tuesdays and due next Tuesday night (mostly)
  - **Gradiance**: immediately and automatically graded
  - **Gradescope**: programming problems, immediate feedback, later also manual grading
  - **Gradescope**: written solution, manual grading
- **Midterm and final (20% each)**
  - Open book, open notes
  - No communication/Internet whatsoever
  - Final is comprehensive, but emphasizes the second half of the course

30

### Course load (contd.)

- **Course project (20%)**
  - Details to be given in the next 1-2 weeks
- **In-class quiz (5%)**
  - To review concepts right away in class – will be open for 5-10 mins
  - Will be announced at least one class in advance and on piazza
  - Each quiz: 50% for attempt on time and 50% for correct solution
  - Lowest score will be dropped (each quiz has same weight)
- **In-class labs (5%)**
  - Practice problems in class (both programming and conceptual) – each lab has the same weight
  - Will be announced at least one class in advance and on piazza
  - Due by the next day after class, 10% bonus points for finishing all problems in class correctly
  - TAs will be around to help you

31

### Tentative office hours schedules

- Locations: TBD.
- See the updated info on the webpage
- More office hours around Tuesday (hws due), but good to start early!

32

### Projects from past years

- RA: next-generation relational algebra interpreter
  - You may get to try it out for Homework #1!
- *Managing tent shifts and schedules!*
- *Tutor-tutee matching*
- *What's in my fridge and what can I cook?*
- *Hearsay: manage your own musics*
- *Dining at Duke (and deliver meals to students)*
- *National Parklopedia: a website to find information about national parks*

• *More examples later - but we expect you to be creative with a new idea!*

33

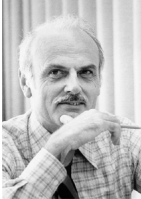
Let's get started!

## Relational Data Model

What is a good model to store data?

34

### Edgar F. Codd (1923-2003)



- Pilot in the Royal Air Force in WW2
- Inventor of the relational model and algebra while at IBM
- Turing Award, 1981

RDBMS = Relational DBMS

[http://en.wikipedia.org/wiki/File:Edgar\\_F\\_Codd.jpg](http://en.wikipedia.org/wiki/File:Edgar_F_Codd.jpg)

35

### The famous "Beers" database

36

See online database for more tuples

## "Beers" as a Relational Database

**Bar**

name	address
The Edge	108 Morris Street
Satisfaction	905 W. Main Street

**Beer**

Name	brewer
Budweiser	Anheuser-Busch Inc.
Corona	Grupo Modelo
Dixie	Dixie Brewing

**Drinker**

name	address
Amy	100 W. Main Street
Ben	101 W. Main Street
Dan	300 N. Duke Street

**Serves**

bar	beer	price
The Edge	Budweiser	2.50
The Edge	Corona	3.00
Satisfaction	Budweiser	2.25

**Frequents**

drinker	bar	times_a_week
Ben	Satisfaction	2
Dan	The Edge	1
Dan	Satisfaction	2

**Likes**

drinker	beer
Amy	Corona
Dan	Budweiser
Dan	Corona
Ben	Budweiser

37

## Relational data model

- A database is a collection of **relations** (or **tables**)
- Each relation has a set of **attributes** (or **columns**)
- Each attribute has a name and a **domain** (or **type**)
  - Set-valued attributes are not allowed
- Each relation contains a **"set"** of **tuples** (or **rows**)
  - Each tuple has a value for each attribute of the relation
  - Duplicate tuples are not allowed (Two tuples are duplicates if they agree on all attributes)
  - Ordering of rows doesn't matter (even though output is always in some order)
- However, SQL supports **"bag"** or duplicate tuples (why?)
  - Simplicity is a virtue**
    - not a weakness!

**Serves**

bar	beer	price
The Edge	Budweiser	2.50
The Edge	Corona	3.00
Satisfaction	Budweiser	2.25

38

## Schema vs. instance

- Schema**
  - Beer (name string, brewer string)
  - Serves (bar string, beer string, price float)
  - Frequents (drinker string, bar string, times\_a\_week int)
- Instance**
  - Actual tuples or records

☞ Compare to **types** vs. collections of **objects of these types** in a programming language

**Beer**

Name	brewer
Budweiser	Anheuser-Busch Inc.
Corona	Grupo Modelo
Dixie	Dixie Brewing

**Serves**

bar	beer	price
The Edge	Budweiser	2.50
The Edge	Corona	3.00
Satisfaction	Budweiser	2.25

**Frequents**

drinker	bar	times_a_week
Ben	Satisfaction	2
Dan	The Edge	1
Dan	Satisfaction	2

39

## SQL: Querying a RDBMS

- SQL: **Structured Query Language**
  - Pronounced "S-Q-L" or "sequel"
  - The standard query language supported by most DBMS
  - First developed at IBM System R
  - Follows ANSI standards

40

## Basic queries: SFW statement

- SELECT**  $A_1, A_2, \dots, A_n$   
**FROM**  $R_1, R_2, \dots, R_m$   
**WHERE** *condition*
- SELECT, FROM, WHERE are often referred to as SELECT, FROM, WHERE **"clauses"**

41

## Example: reading a table

- SELECT \***  
**FROM** Serves
- Single-table query
- WHERE** clause is optional
- \* is a short hand for "all columns"

**Serves**

bar	beer	price
The Edge	Budweiser	2.50
The Edge	Corona	3.00
Satisfaction	Budweiser	2.25

42

## Example: selecting few rows

- `SELECT beer AS mybeer`  
`FROM Serves`  
`WHERE price < 2.75`

Serves		
bar	beer	price
The Edge	Budweiser	2.50
The Edge	Corona	3.00
Satisfaction	Budweiser	2.25

- `SELECT beer`  
`FROM Serves`  
`WHERE bar = 'The Edge'`

What does these return?

- `SELECT` list can contain expressions  
Can also use built-in functions such as `SUBSTR`, `ABS`, etc.
- String literals (case sensitive) are enclosed in [single quotes](#)
- “AS” is optional
- Do not want duplicates? Write `SELECT DISTINCT beer ...`

43

## Example: Join

- Find addresses of all bars that ‘Dan’ frequents

- Which tables do we need?

44

## Example: Join

- Find addresses of all bars that ‘Dan’ frequents

Bar	
name	address
The Edge	108 Morris Street
Satisfaction	905 W. Main Street

bar	beer	price
The Edge	Budweiser	2.50
The Edge	Corona	3.00
Satisfaction	Budweiser	2.25

Name	brewer
Budweiser	Anheuser-Busch Inc.
Corona	Grupo Modelo
Dixie	Dixie Brewing

drinker	bar	times_a_week
Ben	Satisfaction	2
Dan	The Edge	1
Dan	Satisfaction	2

Which tables do we need?

How do we combine them?

name	address
Amy	100 W. Main Street.
Ben	101 W. Main Street
Dan	300 N. Duke Street

drinker	beer
Amy	Corona
Dan	Budweiser
Dan	Corona
Ben	Budweiser

Likes

45

## Example: Join

- Find addresses of all bars that ‘Dan’ frequents

- `SELECT B.address`  
`FROM Bar B, Frequents F`  
`WHERE B.name = F.bar`  
`AND F.drinker = 'Dan'`

Bar	
name	address
The Edge	108 Morris Street
Satisfaction	905 W. Main Street

- Okay to omit `table_name` in `table_name.column_name` if `column_name` is unique

- Can use “Aliases” for convenience

- “Bar as B” or “Bar B”

drinker	bar	times_a_week
Ben	Satisfaction	2
Dan	The Edge	1
Dan	Satisfaction	2

Frequents

46

## Let's try SQL in class!

(See how to access the pgweb interface for a small “Beers” database on the slides posted on the course website)

Next: semantics of SFW statements in SQL

47

## Announcements (Tue, 01/09)

- You should be on Sakai, Piazza, Gradescope
  - If you are not there or recently enrolled, please contact the instructor
- You will receive instructions on installing the VM
  - Please follow Piazza posts, all notifications will be posted there and you should receive emails right away
- First homework to be released on next class  
Tuesday 01/14, due in a week
  - No in-class quiz or labs unless explicitly announced in the class before (and posted on Piazza)

48



### Semantics of SFW

- `SELECT E1, E2, ..., En`  
`FROM R1, R2, ..., Rm`  
`WHERE condition`
- For each  $t_1$  in  $R_1$ :  
 For each  $t_2$  in  $R_2$ : ... .. 1. Apply "FROM"  
 For each  $t_m$  in  $R_m$ : Form cross-product of  $R_1, \dots, R_m$

If condition is true over  $t_1, t_2, \dots, t_m$ :

2. Apply "WHERE"  
 Only consider satisfying rows

Compute and output  $E_1, E_2, \dots, E_n$  as a row

3. Apply "SELECT"  
 Output the desired columns

49

### Step 1: Illustration of Semantics of SFW

- NOTE: This is "NOT HOW" the DBMS outputs the result, but "WHAT" is outputs!

Form Cross product of two relations

- `SELECT B.address`  
`FROM Bar B, Frequents F`  
`WHERE B.name = F.bar`  
`AND F.drinker LIKE 'Dan%'`

name	address	drinker	bar	times_a_w week
The Edge	108 Morris Street	Ben	Satisfaction	2
The Edge	108 Morris Street	Dan	The Edge	1
The Edge	108 Morris Street	Dan	Satisfaction	2
Satisfaction	905 W. Main Street	Ben	Satisfaction	2
Satisfaction	905 W. Main Street	Dan	The Edge	1
Satisfaction	905 W. Main Street	Dan	Satisfaction	2

Bar	name	address
The Edge	108 Morris Street	
Satisfaction	905 W. Main Street	

Frequents	drinker	bar	times_a_week
Ben	Satisfaction	2	
Dan	The Edge	1	
Dan	Satisfaction	2	

50

### Step 2: Illustration of Semantics of SFW

- NOTE: This is "NOT HOW" the DBMS outputs the result, but "WHAT" is outputs!

Discard rows that do not satisfy WHERE condition

- `SELECT B.address`  
`FROM Bar B, Frequents F`  
`WHERE B.name = F.bar`  
`AND F.drinker LIKE 'Dan%'`

name	address	drinker	bar	times_a_w week
The Edge	108 Morris Street	Ben	Satisfaction	2
The Edge	108 Morris Street	Dan	The Edge	1
The Edge	108 Morris Street	Dan	Satisfaction	2
Satisfaction	905 W. Main Street	Ben	Satisfaction	2
Satisfaction	905 W. Main Street	Dan	The Edge	1
Satisfaction	905 W. Main Street	Dan	Satisfaction	2

Bar	name	address
The Edge	108 Morris Street	
Satisfaction	905 W. Main Street	

Frequents	drinker	bar	times_a_week
Ben	Satisfaction	2	
Dan	The Edge	1	
Dan	Satisfaction	2	

51

### Step 3: Illustration of Semantics of SFW

- NOTE: This is "NOT HOW" the DBMS outputs the result, but "WHAT" is outputs!

Output the "address" output of rows that survived

- `SELECT B.address`  
`FROM Bar B, Frequents F`  
`WHERE B.name = F.bar`  
`AND F.drinker LIKE 'Dan%'`

name	address	drinker	bar	times_a_w week
The Edge	108 Morris Street	Ben	Satisfaction	2
The Edge	108 Morris Street	Dan	The Edge	1
The Edge	108 Morris Street	Dan	Satisfaction	2
Satisfaction	905 W. Main Street	Ben	Satisfaction	2
Satisfaction	905 W. Main Street	Dan	The Edge	1
Satisfaction	905 W. Main Street	Dan	Satisfaction	2

Bar	name	address
The Edge	108 Morris Street	
Satisfaction	905 W. Main Street	

Frequents	drinker	bar	times_a_week
Ben	Satisfaction	2	
Dan	The Edge	1	
Dan	Satisfaction	2	

52

### Final output: Illustration of Semantics of SFW

- NOTE: This is "NOT HOW" the DBMS outputs the result, but "WHAT" is outputs!

Output the "address" output of rows that survived

- `SELECT B.address`  
`FROM Bar B, Frequents F`  
`WHERE B.name = F.bar`  
`AND F.drinker LIKE 'Dan%'`

Bar	name	address
The Edge	108 Morris Street	
Satisfaction	905 W. Main Street	

address
108 Morris Street
905 W. Main Street

Frequents	drinker	bar	times_a_week
Ben	Satisfaction	2	
Dan	The Edge	1	
Dan	Satisfaction	2	

53