

# Database Systems I, CSCI-GA.2433-001

## New York University, Fall 2019

instructor: Dennis Shasha  
shasha@cs.nyu.edu  
212-998-3086  
Courant Institute  
New York University  
251 Mercer Street  
NY, NY 10012 USA

Office Hours: On Mondays by appointment before class in Warren Weaver lobby

April 11, 2019

## 1 Goals

The course is divided into three major parts:

1. To understand the design of data intensive systems given an underlying database management systems (i.e. entity-relationship conceptual design, normalization theory, query languages).
2. To understand something about internals (indexes, query processing).
3. Some advanced topics (analytical processing, data cleaning and preparation, introduction to concurrency control and recovery and to distribution).

Most of the course notes we use were written by Professor Zvi Kedem, though I've modified a little (and at times skip some slides). They are all here on the website. Here is the order of presentation:

- 01\_Goals\_Of\_The\_Course.pptx

- 05\_SQL\_As\_Data\_Manipulation\_Language.pptx
- 06\_SQL\_As\_Data\_Definition\_And\_Control\_Language.pptx
- 02\_Modeling\_Enterprise\_With\_ER\_Diagrams.pptx
- 03\_From\_ER\_Diagrams\_To\_Relational\_Databases.pptx
- 04\_Relational\_Algebra\_With\_SQL\_Equivalents.pptx
- 07\_Logical\_Design\_With\_Normalization.pptx
- 08\_Physical\_Design\_And\_Query\_Execution\_Concepts.pptx
- 11\_Online\_Analytical\_Processing.pptx
- Some introduction to data cleaning, transaction processing, recovery, distribution etc.

## **2 Mechanics**

YOU MUST BE ENROLLED IN THIS CLASS TO SIT IN ON THE LECTURES.

### **2.1 Texts and Notes**

There are no textbooks, but you will look up manuals (notably for mySQL and for data structure libraries) on the web.

### **2.2 Prerequisites**

Fundamental Algorithms I is a co-requisite.

### **2.3 Course Requirements**

three problem sets (40%), project (60%). There are no exams in this course.

LATE HOMEWORKS OR PROJECTS WILL NOT BE ACCEPTED without a note from your physician or from your employer. (We will discuss

the solutions on the day you hand in the assignment. That's why I don't want any late homeworks. As for the project, this is a question of fairness.)

On the other hand, collaboration on the problem sets IS allowed. You may work together with one other student and sign both of your names to a single submitted homework. Both of you will receive the grade that the homework merits. So, you may work alone or in a team of two, but no team larger than two. Please choose your partner carefully.

## 2.4 How to Set Up MySQL

You will also need to set up a database using mysql as follows:

1. Go to this website: <https://cims.nyu.edu/webapps/databases>
2. Log in with your CIMS username (netID) and password
3. create a database. Take note of the password, which will be automatically assigned (you should see it at the top of the webpage). In my case, I created a database that I called db1 which then became shasha\_db1. We'll call the password PASSWD
4. Next ssh into access.cims.nyu.edu
5. Now, I typed (but your name will be different than shasha\_db1)

```
mysql -h warehouse.cims.nyu.edu -u shasha -p shasha_db1
```

and then typed in the password I was assigned.

6. Then I was good to go. I could create tables etc. Here are some commands:

```
show databases;
show tables;
create table friends(name1 varchar(20), name2 varchar(20));
insert into friends values("bob", "alice");
select * from friends;
```

7. Now suppose that you have a file foo.csv with vertical bar delimiters on your local machine (in my case in /dbcourse1.d/foo.csv)

```
carol|ted
susan|paul
tyler|cloe
```

You can then import that data into friends as follows:

```
LOAD DATA LOCAL INFILE '~/dbcourse1.d/foo.csv' INTO TABLE friends FIELDS TERMINATED BY '|'
```

## 2.5 Project: A miniature relational database with order

This project is due on Monday Dec 2, 2019 at 4:30 PM.

Given ordered tables (array-tables) whose rows consist of strings and integers, you are to write a program which will

- Perform the basic operations of relational algebra: selection, projection, join, group by, and count, sum and avg aggregates. The comparators for select and join will be  $=$ ,  $<$ ,  $>$ ,  $!=$ ,  $\geq$ ,  $\leq$
- Because the array-tables are potentially ordered, you can sort an array-table by one or more columns, and running moving sums and average aggregates on a column of an array-table.
- Import a vertical bar delimited file into an array-table (in the same order), export from an array-table to a file preserving its order, and assign the result of a query to an array-table.
- Each operation will be on a single line. Each time you execute a line, you should print the time it took to execute.
- You will support in memory B-trees and hash structures. You are welcome to take those implementations from wherever you can find them, but you must say where.
- Your program should be written in python or java. You will hand in clean and well structured source code in which each function has a header that says (i) what the function does, (ii) what its inputs are and what they mean (iii) what the outputs are and mean (iv) any side effects to globals.
- You must ensure that your software runs on the Courant Institute (cims) machine `crunchy5.cims.nyu.edu`

- You may NOT use any relational algebra or SQL library or system (e.g. no SQLite, no mySQL, no other relational database system, no Pandas ). Stick pretty much to the standard stuff (e.g. in Python: numpy, core language features, string manipulation, random number generators, and data structure support for in memory B-trees and hash structures). You may not use anyone else's code (other than for the data structure implementation). Doing so will constitute plagiarism.

We will run your programs on test cases of our choosing. For ease of parsing there will be one operation per line. Comments begin with // and go to the end of the line. For example,

```
R := inputfromfile(sales1) // import vertical bar delimited foo, first line
    // has column headers.
    // Suppose they are saleid|itemid|customerid|storeid|time|qty|pricerange
R1 := select(R, (time > 50) or (qty < 30))
    // select * from R where time > 50 or qty < 30
R2 := project(R1, saleid, qty, pricerange) // select saleid, qty, pricerange
    // from R1
R3 := avg(R1, qty) // select avg(qty) from R1
R4 := sumgroup(R1, time, qty) // select qty, sum(time) from R1 group by qty
R5 := sumgroup(R1, qty, time, pricerange) // select sum(qty), time,
    // pricerange from R1 group by time, pricerange
R6 := avggroup(R1, qty, pricerange) // select avg(qty), pricerange
    // from R1 group by by pricerange
S := inputfromfile(sales2) // suppose column headers are
    // saleid|I|C|S|T|Q|P
T := join(R, S, R.customerid = S.C) // select * from R, S
    // where R.customerid = S.C
T1 := join(R1, S, R1.qty > S.Q) // select * from R1, S where R1.qty > S.Q
T2 := sort(T1, S_C) // sort T1 by S_C
T2prime := sort(T1, R_time, S_C) // sort T1 by R_itemid, S_C (in that order)
T3 := movavg(T2, R_qty, 3) // perform the three item moving average of T2
    // on column R_qty. This will be as long as R_qty with the three way
    // moving average of 4 8 9 7 being 4 6 7 8
T4 := movsum(T2, R_qty, 5) // perform the five item moving sum of T2
    // on column R_qty
Q1 := select(R, qty = 5) // select * from R where qty=5
Btree(R, qty) // create an index on R based on column qty
```

```

    // Equality selections and joins on R should use the index.
Q2 := select(R, qty = 5) // this should use the index
Q3 := select(R, itemid = 7) // select * from R where itemid = 7
Hash(R,itemid)
Q4 := select(R, itemid = 7) // this should use the hash index
Q5 := concat(Q4, Q2) // concatenate the two tables (must have the same schema)
    // Duplicate rows may result (though not with this example).
outputtofile(Q5, bar) // This should output the table Q4 with vertical bar separators

```

Our tests may operate on different files with different column headers. Our queries may use different parameter values (e.g. 14 way moving average). Our joins may be on different fields.

Some constraints to make your life easier:

- There will be no syntax errors in our tests. However, white space may vary from the above examples.
- The only aggregates are count, sum, and avg and the corresponding countgroup, sumgroup, and avggroup.
- The only moving aggregates are movsum and movavg. There is no group by for moving sums and averages.
- All data is in main memory.
- The joins are on single columns
- The selects will be all ors or all ands.

Here is an example of the first few lines of the sales1 file:

```

saleid|itemid|customerid|storeid|time|qty|pricerange
45|133|2|63|49|23|outrageous
658|75|2|89|46|43|outrageous
149|103|2|23|67|2|cheap
398|82|2|41|3|27|outrageous
147|81|2|4|92|11|outrageous
778|75|160|72|67|17|supercheap
829|112|2|70|63|43|supercheap

```

101|105|2|9|74|28|expensive  
940|62|2|90|67|39|outrageous  
864|119|12|38|67|49|outrageous  
288|46|2|95|67|26|outrageous  
875|83|59|56|59|20|outrageous  
783|86|180|29|67|46|outrageous  
289|16|2|95|92|2|cheap

Full example files can be found here:

<http://cs.nyu.edu/cs/faculty/shasha/papers/sales1>  
<http://cs.nyu.edu/cs/faculty/shasha/papers/sales2>

You will hand in your code using Rezip inside a Docker Virtual Machine so that it is reproducible across platforms.