

COMS W4111: Introduction to Databases Spring 2024, Sections 002/V02

Homework 1 *Introduction to Core Concepts, ER Modeling, Relational Algebra, SQL*

Introduction

This notebook contains Homework 1. **Both Programming and Nonprogramming tracks should complete this homework.**

Submission Instructions

- You will submit **PDF and ZIP files** for this assignment. Gradescope will have two separate assignments for these.
 - For the PDF:
 - The most reliable way to save as PDF is to go to your browser's menu bar and click **File -> Print**. Switch the orientation to landscape mode, and hit save.
 - **MAKE SURE ALL YOUR WORK (CODE AND SCREENSHOTS) IS VISIBLE ON THE PDF. YOU WILL NOT GET CREDIT IF ANYTHING IS CUT OFF.** Reach out for troubleshooting.
 - **MAKE SURE YOU DON'T SUBMIT A SINGLE PAGE PDF.** Your PDF should have multiple pages.
 - For the ZIP:
 - Zip a folder containing this notebook and any screenshots.
 - You may delete any unnecessary files, such as caches.
-

Add Student Information

In [1]: `# Print your name, uni, and track below`

```
name = "Lexa Huang"
uni = "sh4350"
track = "Programming Track"

print(name)
print(uni)
print(track)
```

Lexa Huang
sh4350
Programming Track

Setup

SQL Magic

The `sql` extension was installed in HW0. Double check that if this cell doesn't work.

In [6]: `%load_ext sql`

The `sql` extension is already loaded. To reload it, use:
`%reload_ext sql`

You may need to change the password below.

In [7]: `%sql mysql+pymysql://root:dbuserdbuser@localhost`

In [8]: `%sql SELECT * FROM db_book.student WHERE ID = 12345`

* mysql+pymysql://root:***@localhost
1 rows affected.

Out[8]:

| ID | name | dept_name | tot_cred |
|-------|---------|------------|----------|
| 12345 | Shankar | Comp. Sci. | 32 |

Python Libraries

In [5]: `from IPython.display import Image`
`import pandas`

Written Questions

Chapter 1 from the recommended textbook [Database System Concepts, Seventh Edition](#) covers general information and concepts about databases and database management systems. Lecturing on the general and background information is not a good use of precious class time. To be more efficient with class time, the chapter 1 information is a reading assignment.

Answering the written questions in HW 1, Part 1 does not require purchasing the textbook and reading the chapter. The [chapter 1 slides](#) provided by the textbook authors provide the necessary information. In some cases, students may also have to search the web or other sources to "read" the necessary information.

When answering the written questions, do not "bloat". The quantity of words does not correlate with the quality of the answer. We will deduct points if you are not succinct. The answers to the questions require less than five sentences or bullet points.

"If you can't explain something in a few words, try fewer."

You may use external resources, but you should cite your sources.

W1

What is a database management system and how do relational databases organize data?

A database management system is a system to organize data not built directly on top of file systems. It addresses the problems of: data redundancy and inconsistency, difficulty in accessing data, data isolation, integrity problems, atomicity of updates, concurrent access by multiple users, and security problems.

Relational databases organize data by storing it in tables, and has core concepts of relation, tuple (row), and column (Attribute).

W2

Columbia University uses several applications that use databases to run the university. Examples are SSOL and CourseWorks. An alternate approach could be letting students, faculty, administrators, etc. use shared Google Sheets to create, retrieve, update, and delete information. What are some problems with the shared spread sheet approach and what functions do DMBS implement to solve the problems?

1. Security problems: It is hard to provide user access to some, but not all, data.

Access would need to be carefully controlled to prevent instances where an unauthorized user modifies information they should not be able to modify. For example, students choose to change their own grades or other students' grades to

whatever score they want. DBMS solves this problem by having security features like user authentication, role-based access control, and encryption.

2. Concurrent Access by multiple users: Uncontrolled concurrent access can lead to inconsistencies in data or even loss of data. DBMS use locking mechanisms and transaction management to ensure data integrity during concurrent access.
3. Scalability and performance: Spreadsheets may not be able to scale to the large number of students, faculty, and administrators in the school, and the large amount of data that needs to be tracked. DBMS are designed for scalability and handling large volumes of data and transactions efficiently, by using indexing, partitioning, and optimized query processing.

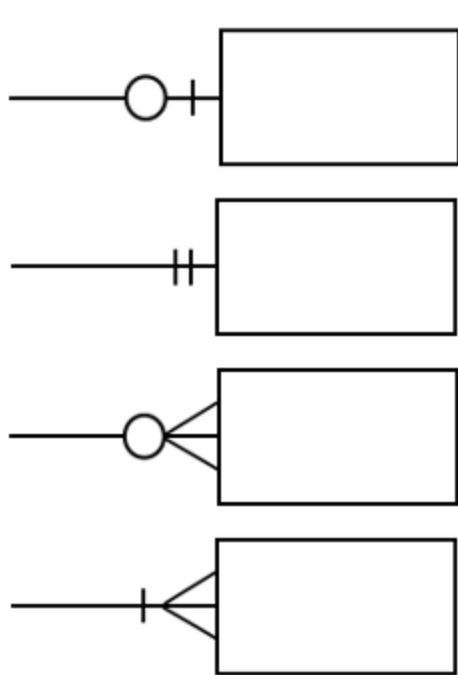
W3

Explain the differences between SQL, MySQL Server and DataGrip.

SQL is a standardized programming language used for managing and manipulating relational databases, while MySQL Server is an example of an RDBMS (Relational Database Management System) that implements SQL for its operations. On the other hand, DataGrip is an integrated development environment (IDE) developed by JetBrains, designed to provide a user-friendly interface for database administrators and developers. It supports various database systems, including MySQL, enabling efficient interaction with databases through features like advanced SQL coding assistance, query execution, and data management tools.

W4

Crow's Foot Notation has four endings for relationship lines. Briefly explain the meaning of each ending.



1. One or Zero: optional relationship, there might be no connection between the entities, or there might be one.
2. One and only One: Not an optional relationship. No more than one relationship is possible between instances.
3. Zero or Many: optional relationship, there might be no connection between the entities, or there might be one or more.
4. One or Many: Not an optional relationship There might be more than one relationship between entities.

Source: <https://www.gleek.io/blog/crows-foot-notation>

W5

What is a primary key and why is it important?

The primary key is a composite key that indicates that the value occurs at most once. It is important as it acts as the unique identifier for each member in the set.

W6

The relational algebra is closed under the operators. Explain what this means and give an example.

This means that whenever you apply a relational algebra operation to a relation in a database, the result is always another relation. This ensures that operations can be

composed and nested.

W7

Some of the Columbia University databases/applications represent the year/semester attribute of a section in the form "2023_2". The first four characters are the academic year, and the last character is the semester (1, 2, or 3). The data type for this attribute might be CHAR(6). Using this example, explain the concepts of domain and atomic domain. How is domain different from type?

Domain refers to the set of all values that the year/semester attribute can take on, which in this case includes all the possible combinations of valid years and valid semester numbers (1, 2, or 3). It is constrained by format (YYYY_S), and valid year numbers and semester numbers.

Atomic domain refers to the fact that the values in the domain are indivisible - they are treated as single units and cannot be decomposed into smaller components. The database treats "2023_2" as a single, indivisible entity.

Type refers to "CHAR(6)", and is a technical specification of data storage format, whereas domain refers to the logical and business rules governing the attribute.

W8

Briefly explain the difference between a database schema and database instance.

Database schema is the logical structure of the database, whereas database instance is a snapshot of the data in the database at a given instant in time.

W9

Briefly explain the concepts of data definition language and data manipulation language.

Data definition language is the specific notation for defining database schema. The DDL compiler generates a set of table templates stored in a data dictionary containing metadata such as database schema, integrity constraints, and authorization.

Data manipulation language is the language for accessing and updating the data organized by the appropriate data model, also known as query language. There are two classes of languages, pure and commercial.

W10

What is physical data independence?

It is the ability to modify the physical schema without changing the logical schema.

Entity-Relationship Modeling

Overview

The ability to understand a general description of a requested data model and to transform into a more precise, specified *logical model* is one of the most important skills for using databases. SW and data engineers build applications and data models for end-users. The end-users, product managers and business managers are not SW or data modeling experts. They will express their *intent* in imprecise, text and words.

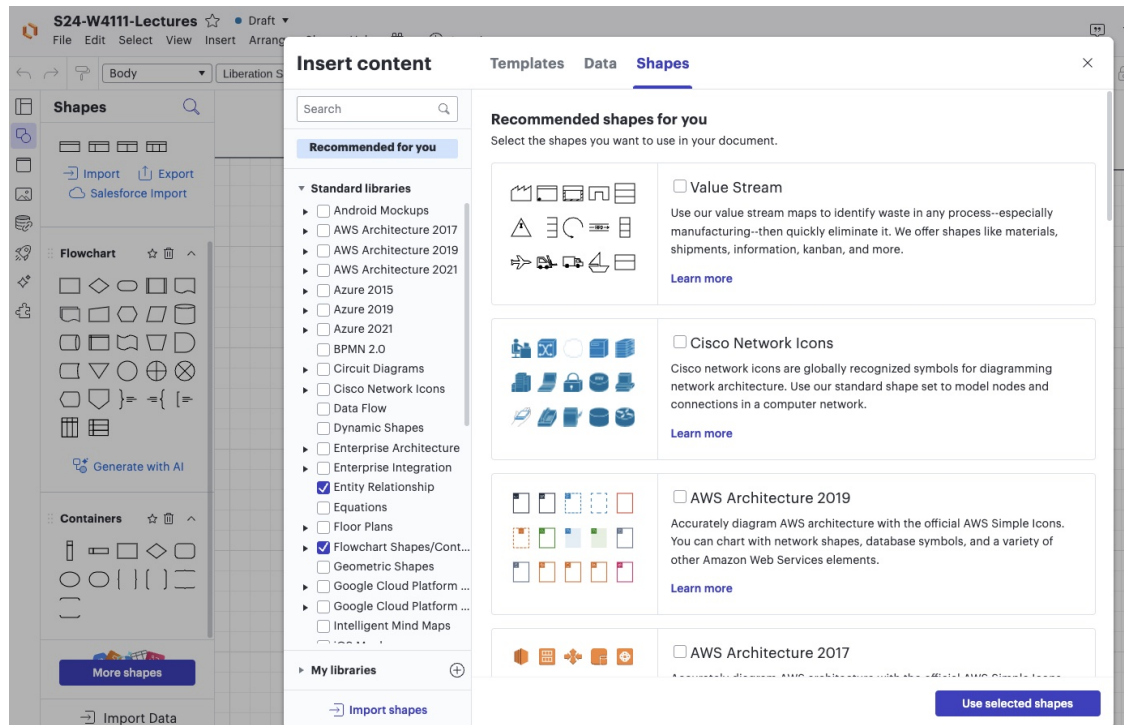
The users and business stakeholder often can understand and interact using a *conceptual model* but details like keys, foreign keys, ... are outside their scope.

In this problem, you will:

- Understand a short written description of a requested data model.
- Produce a *conceptual data model diagram* using Lucidchart.
- Produce a *logical data model diagram* using Lucidchart.

You can sign up for a free [Lucidchart account](#). The free account provides the capabilities you will need for this course.

To draw the diagrams, you need to add the *entity relationship* shapes. Lecture 2 demonstrated how to add the shapes.



Adding Entity Relationship Shapes

We provide a simple [Lucidchart document](#) from Lecture 2 that helps you get started. You need a Lucidchart account to access the document and diagrams.

Data Model Description

The data model represents banks, customers, employees and accounts. The model has the following entity types/sets:

1. *Customer*
2. *Employee* of the banking company
3. *Branch*, which is a location of one of the banks offices
4. *Savings Account*
5. *Checking Account*
6. *Loan*
7. *Portfolio*

Customer has the following properties:

- *customerID*
- *lastName*
- *firstName*
- *email*
- *dateOfBirth*

Employee has the following properties:

- *employeeID*
- *lastName*
- *firstName*
- *jobTitle*

Branch has the following properties:

- *branchID*
- *zipCode*

Savings Account has the following properties:

- *accountID*
- *balance*
- *interestRate*

Checking Account has the following properties:

- *accountID*
- *balance*

Loan has the following properties.

- *loanID*
- *balance*
- *interestRate*

Portfolio has the following properties:

- *portfolioID*
- *createdDate*

The data model has the following relationships:

- *Customer Branch* connects a customer and a branch. A *Customer* is connected to exactly one *Branch*. A *Branch* may have 0, 1 or many customers.
- *Employee Branch* connects an employee and a branch. An *Employee* is connected to exactly one *Branch*. A *Branch* may have 0, 1 or many associated employees.
- *Savings Account Branch*, *Checking Account Branch*, and *Loan Branch* all have the same pattern.
 - An account/loan has exactly one branch.
 - A *Branch* may have 0, 1 or many accounts/loans.
- *Savings Customer*, *Checking Customer*, *Loan Customer*, and *Portfolio Customer* follow the same pattern.
 - The account/loan has exactly one customer.

- The customer may have 0 or 1 of each type of account.
- A *Portfolio* is related to exactly one *Customer*, exactly one *Savings Account*, exactly one *Checking Account*, and exactly one *Loan*.
- *Portfolio Advisor* relates a *Portfolio* and *Employee*. An *Employee* may be the advisor for 0, 1 or many *Portfolios*. A *Portfolio* may have at most one *Employee* advisor.

Answer

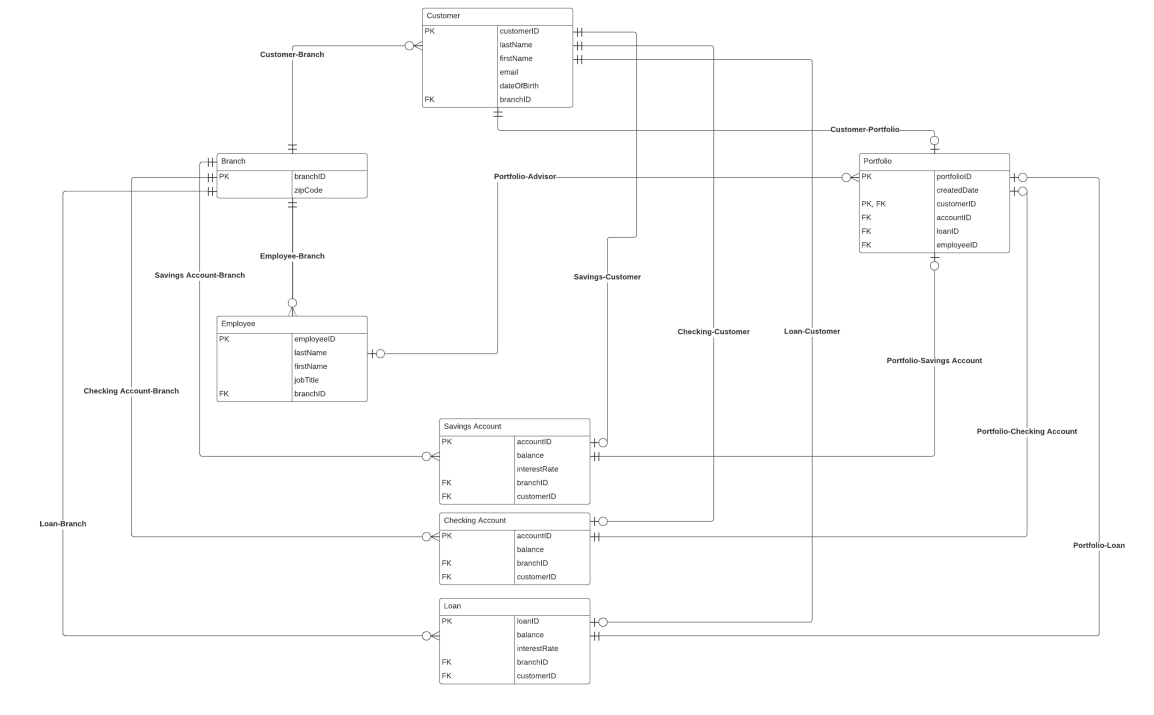
1. Place your Logical Model diagram below.
2. You *may* have to add attributes to entities to implement the model.
3. You *may* make reasonable assumptions. Please document your assumptions below.
You may add comments/notes to your diagram for clarity.

Assumptions:

1. There may be more than one branch located in the same zip code area. Hence, zipCode is not a primary key.
2. The portfolio may have at most one Employee advisor, so it is optional and is a zero-or-one relationship in that the portfolio may have zero advisors.
3. Since the customer may only have 0 or 1 of Portfolio, customerID in portfolio is both a foreign key and a primary key
4. Since the customer can only have maximum 1 savings account, checking account or loan, and maximum 1 portfolio, and each portfolio is related to one and only one savings/checking account or loan, and portfolios are optional for the customer, we can assume that savings account, checking account, and loan are all associated with one-or-zero of portfolio, assuming savings/checking account and loans do not have to be part of a portfolio if the customer doesn't have one.

ER Diagram:

Save your diagram to an image, place in the same directory as your notebook and change the file name in the HTML `img` tag in this Markdown cell.



Logical ER Diagram

Relational Algebra

R-1

The following is the SQL DDL for the `db_book.classroom` table.

```
CREATE TABLE IF NOT EXISTS db_book.classroom
(
    building    VARCHAR(15) NOT NULL,
    room_number VARCHAR(7)  NOT NULL,
    capacity    DECIMAL(4)  NULL,
    PRIMARY KEY (building, room_number)
);
```

Using the notation from the lecture slides, provide the corresponding relation schema definition.

`classroom(building, room_number, capacity)`

Answer Format

For the answers to the relational algebra questions, you will use the [RelaX calculator](#) with the schema associated with the book. Your answer should include the algebra statement

in as text and a screenshot of the execution result. Question **R0** below shows a sample of that the answer will look like.

R0

Write a relational algebra statement that produces a table of the following form:

- ID is the instructor ID
- name is the instructor name
- course_id, sec_id, semester, year of a section
- building, room_number

Note:

1. You will have to use the instructor, teaches and section relations
2. Your answer should only include sections taught in **Comp. Sci.** in **2009**

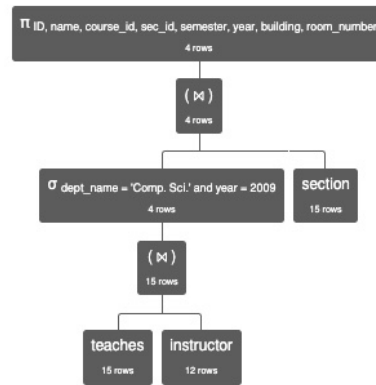
Algebra statement:

```

$$\pi_{ID, name, course\_id, sec\_id, semester, year, building, room\_number}(\sigma_{dept\_name='Comp. Sci.' \wedge year=2009}(\text{teaches} \bowtie \text{instructor})) \bowtie \text{section})$$

```

Execution:



$\pi_{ID, name, course_id, sec_id, semester, year, building, room_number} ((\sigma_{dept_name = 'Comp. Sci.' \text{ and } year = 2009} (teaches \bowtie instructor)) \bowtie section)$

Execution time: 1 ms

| teaches.ID | instructor.name | teaches.course_id | teaches.sec_id | teaches.semester | teaches.year | section.building | section.room_number |
|------------|-----------------|-------------------|----------------|------------------|--------------|------------------|---------------------|
| 10101 | 'Srinivasan' | 'CS-101' | 1 | 'Fall' | 2009 | 'Packard' | 101 |
| 10101 | 'Srinivasan' | 'CS-347' | 1 | 'Fall' | 2009 | 'Taylor' | 3128 |
| 83821 | 'Brandt' | 'CS-190' | 1 | 'Spring' | 2009 | 'Taylor' | 3128 |
| 83821 | 'Brandt' | 'CS-190' | 2 | 'Spring' | 2009 | 'Taylor' | 3128 |

RO Execution Result

R1

Write a relational algebra statement that produces a relation with the columns:

- student.name
- student.dept_name
- student.tot_cred
- instructor.name (the instructor that advises the student)
- instructor.dept_name

Only keep students who have earned more than 90 credits.

Note:

1. You will have to use the student, instructor, and advisor relations.
2. You should only include students that have an advisor, i.e., instructor.name and instructor.dept_name should be non-null for all rows.

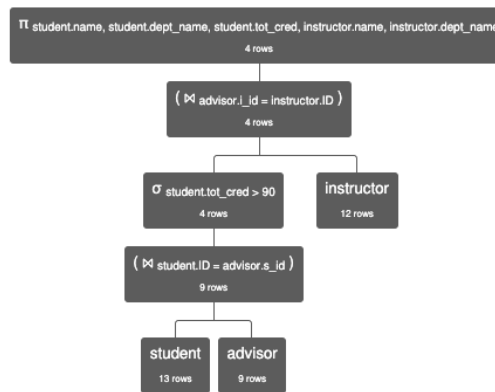
Algebra statement:

```

 $\pi$  student.name, student.dept_name, student.tot_cred,
instructor.name, instructor.dept_name (
  (  $\sigma$  student.tot_cred > 90
    ( student
       $\bowtie$  student.ID = advisor.s_id
        advisor )
    )
   $\bowtie$  advisor.i_id = instructor.ID
    instructor )

```

Execution:



```

 $\pi$  student.name, student.dept_name, student.tot_cred, instructor.name, instructor.dept_name ( (  $\sigma$  student.tot_cred > 90 ( student  $\bowtie$  student.ID =
  advisor.s_id advisor ) )  $\bowtie$  advisor.i_id = instructor.ID instructor )
Execution time: 1 ms

```

| student.name | student.dept_name | student.tot_cred | instructor.name | instructor.dept_name |
|--------------|-------------------|------------------|-----------------|----------------------|
| 'Zhang' | 'Comp. Sci.' | 102 | 'Katz' | 'Comp. Sci.' |
| 'Chavez' | 'Finance' | 110 | 'Singh' | 'Finance' |
| 'Bourikas' | 'Elec. Eng.' | 98 | 'Kim' | 'Elec. Eng.' |
| 'Tanaka' | 'Biology' | 120 | 'Crick' | 'Biology' |

R1 Execution Result

R2

Write a relational algebra statement that produces a relation with the columns:

- course_id
- title
- prereq_course_id
- prereq_course_title

This relation represents courses and their prereqs.

Note:

1. This query requires the `course` and `prereq` tables.
2. Your answer should only include courses in the `Comp. Sci.` department.
3. If a course has no prereqs, `prereq_course_id` and `prereq_course_title` should both be *null*.
4. You *may* have to use table and column renaming.

Algebra statement:

```
temp = (π course.course_id, course.title, prereq_id (
  ( σ course.dept_name = 'Comp. Sci.'
    ( course
      ⋈ course.course_id = prereq.course_id
        prereq)
    )
  ))
```

```
renamed_temp = ρ course_id2←course.course_id,
title2←course.title, prereq_course_id←prereq.prereq_id (π
course.course_id, course.title, prereq.prereq_id ( temp ) )
```

```
final = π course_id2, title2, prereq_course_id, course.title
(
  (renamed_temp
    ⋈ prereq.prereq_course_id = course.course_id
      course
    )
  )
```

```
ρ course_id←course.course_id2, title←course.title2,
prereq_course_id←prereq.prereq_course_id,
prereq_course_title←course.title (π course.course_id2,
course.title2, prereq.prereq_course_id, course.title ( final
) )
```

Execution:



| course_id | course_title | prereq_id | prereq_title |
|-----------|--------------------------|-----------|----------------------------|
| CS-190 | Game Design | CS-101 | Intro. to Computer Science |
| CS-315 | Robotics | CS-101 | Intro. to Computer Science |
| CS-319 | Image Processing | CS-101 | Intro. to Computer Science |
| CS-347 | Database System Concepts | CS-101 | Intro. to Computer Science |

R2 Execution Result

SQL

New Database

MySQL Tutorial is a good site with information that complements and extends the core material in our course. Much of the material the site covers is applicable to other SQL products. MySQL Tutorial uses an interesting dataset that is more complex than the simple "db_book" database. This is the Classic Models Dataset. The complexity allows us to better appreciate more complex SQL concepts.

You learned how to run a SQL script/file as part of HW0. Use the same approach to load and create the Classic Models Database . The file is classic-models-database.sql and is in the HW1 folder.

To test loading the data, you can use the cell below.

```
In [12]: %sql show tables;

* mysql+pymysql://root:***@localhost
8 rows affected.
```


Out [12]: **Tables_in_classicmodels**

| |
|--------------|
| customers |
| employees |
| offices |
| orderdetails |
| orders |
| payments |
| productlines |
| products |

In [11]: `%sql USE classicmodels;`

* mysql+pymysql://root:***@localhost
0 rows affected.

Out [11]: []

SQL 1

This query uses `customers` and `employees`.

Write and execute a SQL query that produces a table with the following columns:

- `customerContactName`
- `customerPhone`
- `salesRepName`

Only keep customers from France. Order your output by `customerContactName`.

Notes:

- The names of your columns must match exactly with what is specified.
- `customerContactName` can be formed by combining `customers.contactFirstName` and `customers.contactLastName`.
- `salesRepName` can be formed by combining `employees.firstName` and `employees.lastName`.

In [26]: `%sql`

```
SELECT
    CONCAT(customers.contactFirstName, ' ', customers.contactLastName) AS customerContactName,
    customers.phone AS customerPhone,
    CONCAT(employees.firstName, ' ', employees.lastName) AS salesRepName
FROM
    customers
JOIN
```

```

employees ON customers.salesRepEmployeeNumber = employees.employeeNumber
WHERE
    customers.country = 'France'
ORDER BY
    customerContactName;

```

* mysql+pymysql://root:***@localhost
12 rows affected.

Out [26]:

| customerContactName | customerPhone | salesRepName |
|---------------------|------------------|------------------|
| Annette Roulet | 61.77.6555 | Gerard Hernandez |
| Carine Schmitt | 40.32.2555 | Gerard Hernandez |
| Daniel Tonini | 30.59.8555 | Gerard Hernandez |
| Daniel Da Silva | +33 1 46 62 7555 | Loui Bondur |
| Dominique Perrier | (1) 47.55.6555 | Loui Bondur |
| Frédérique Citeaux | 88.60.1555 | Gerard Hernandez |
| Janine Labrune | 40.67.8555 | Gerard Hernandez |
| Laurence Lebihan | 91.24.4555 | Loui Bondur |
| Marie Bertrand | (1) 42.34.2555 | Loui Bondur |
| Martine Rancé | 20.16.1555 | Gerard Hernandez |
| Mary Saveley | 78.32.5555 | Loui Bondur |
| Paul Henriot | 26.47.1555 | Loui Bondur |

SQL 2

This query uses `employees`, `customers`, `orders`, `orderdetails`.

Write and execute a SQL query that produces a table showing the amount of money each sales rep has generated.

Your table should have the following columns:

- `salesRepName`
- `moneyGenerated`

Order your output from greatest to least `moneyGenerated`.

Notes:

- The names of your columns must match exactly with what is specified.
- `salesRepName` can be formed by combining `employees.firstName` and `employees.lastName`.
- To calculate `moneyGenerated`:

- Every order in `orders` is associated with multiple rows in `orderdetails`. The total amount of money spent on an order is the sum of `quantityOrdered * priceEach` for all the associated rows in `orderdetails`. **Only consider orders that are Shipped.**
- A customer can have multiple orders. The total amount of money a customer has spent is the sum of the money spent on all that customer's orders.
- A sales rep can have multiple customers. `moneyGenerated` is the sum of the money spent by all that sales rep's customers.
- You may find the [WITH keyword](#) to be useful for cleaner code.

In [32]: `%%sql`

```
WITH OrderAmounts AS (
    SELECT
        orders.customerNumber,
        SUM(orderdetails.quantityOrdered * orderdetails.priceEach) AS orderTotal
    FROM
        orderdetails
    JOIN
        orders ON orderdetails.orderNumber = orders.orderNumber
    WHERE
        orders.status = 'Shipped'
    GROUP BY
        orders.customerNumber
),
CustomerTotals AS (
    SELECT
        customers.salesRepEmployeeNumber,
        SUM(OrderAmounts.orderTotal) AS customerTotal
    FROM
        OrderAmounts
    JOIN
        customers ON OrderAmounts.customerNumber = customers.customerNumber
    GROUP BY
        customers.salesRepEmployeeNumber
)
SELECT
    CONCAT(employees.firstName, ' ', employees.lastName) AS salesRepName,
    SUM(CustomerTotals.customerTotal) AS moneyGenerated
FROM
    CustomerTotals
JOIN
    employees ON CustomerTotals.salesRepEmployeeNumber = employees.employeeNumber
GROUP BY
    salesRepName
ORDER BY
    moneyGenerated DESC;
```

* mysql+pymysql://root:***@localhost
15 rows affected.

Out [32]:

| salesRepName | moneyGenerated |
|------------------|----------------|
| Gerard Hernandez | 1065035.29 |
| Leslie Jennings | 1021661.89 |
| Pamela Castillo | 790297.44 |
| Larry Bott | 686653.25 |
| Barry Jones | 637672.65 |
| George Vanauf | 584406.80 |
| Loui Bondur | 569485.75 |
| Peter Marsh | 523860.78 |
| Andy Fixter | 509385.82 |
| Foon Yue Tseng | 488212.67 |
| Mami Nishi | 457110.07 |
| Steve Patterson | 449219.13 |
| Martin Gerard | 387477.47 |
| Julie Firrelli | 386663.20 |
| Leslie Thompson | 307952.43 |

In []: