

COMS W 4111-002

W4111 - Introduction to Databases, Section 002, Fall 2021

Take Home Final

Exam Instructions

Overview

The Final Exam is worth 30 points out of the semester's total points. There are 10 questions of varying difficulty worth varying points. The amount of points is not necessarily indicative of difficulty/length of the question, but you can use it as a rough guide. The grade for the final is in the range 0-100. We map the score to final point by multiplying by $\frac{30}{100}$.

The Final Exam is open note, open book, open internet. You **may not collaborate** with other students. Posts on EdStem must be made private for you and the instructors only. Any common questions or clarifications will be made by the instructors on the Final Exam pinned thread. Students **are responsible** for monitoring the thread for corrections and clarifications.

You must cite any online sources in the comments Markdown cell for each questions.

Overview of Questions

1. Written — Core Databases Concepts (10 pts)
2. Relational Algebra (10 pts)
3. SQL Design and Query (10 pts)
4. Neo4j Design and Query Queries (10 pts)
5. MongoDB Design and Query (10 pts)
6. Implementation Scenario 1: Modeling and Implementing [RACI] in a Database (15 pts)
7. ~~Implementation Scenario 2: Data Model Comparisons (20 pts)~~
8. Impementation Scenario 3: Data Model Transformation (15 pts)

Note: I decided to drop the data model comparison to make the exam easier. So, everyone get's a free 20 points. Also, remember that **I never curve down**.

Submission Information

This exam is **due Sunday, December 19 at 11:59pm ET** to Gradescope. **You may not use Late Days.**

You submit a zip file containing the main Jupyter Notebook (this file), a PDF of this notebook, and several files in the folder. Each questions provides detailed instructions of

how to complete the question.

Your PDFs must be high enough resolution that the text is legible. It must be printed onto standard 8.5x11in pages. Any images that you embed **MUST** be visible in the PDF. Do not use HTML to embed your images or they will not be visible when you export to PDF.

Failure to meet these formatting specifications will result in a 0.

As always, respect for the individual is paramount. We will accommodate special circumstances, but we must be notified and discuss in advance.

Environment Setup and Test

Note: If you have already done the environment setup tests and succeeded, you only need to run the cells that:

1. Import `mysql_check`, `neo4j_check` and `mongodb_check`.
2. Run the cells that set the DB connection information (user ID, password, URL, ...) for the various databases.
3. You can go directly to the questions.

Instructions

This section tests your environment. You **MUST** completely follow and comply with the instructions.

Implementation Files

- Several of the questions requiring calling databases from Python code. The python code is simple and implements database queries and operations. This complies with the department's guidelines for *non-programming*.
- There is a section for testing access to each of MySQL, MongoDB and Neo4j. You **must** have installed or have access to the databases, and if locally installed the database must be running.

MySQL

- Download and load the [Classic Models \(https://www.mysqltutorial.org/mysql-sample-database.aspx\)](https://www.mysqltutorial.org/mysql-sample-database.aspx) database into MySQL. The download site provides installation instructions.
- The comments in the code snippets below provide instructions for completing and executing each cell.

```
In [1]: %load_ext autoreload
        %autoreload 2
```

```
In [2]: # Import the MySQL test and implementation template/helper functions
        # You do not need to modify this cell. You only need to implement it.
        #
        import mysql_check
```

```
In [3]: #
        # Call the function below to set the user, password and host for your
        # YOU MUST set the variables to the correct names for instance.
        #
        db_user = "root"
        db_password = "password"
        db_host = "127.0.0.1"

        mysql_check.set_connect_info(db_user, db_password, db_host)
```

```
In [4]: #
        # Execute the code below. Your answer should be the same as the example
        #
        df = mysql_check.test_pymysql()
        df
```

Out[4]:

Tables_in_classicmodels	
0	customers
1	employees
2	offices
3	orderdetails
4	orders
5	payments
6	productlines
7	products

```
In [5]: #
        # Execute the cell below. Your result should match the example.
        #
        result_df = mysql_check.test_sql_alchemy()
        result_df
```

Out[5]:

	customerNumber	customerName	country
0	103	Atelier graphique	France
1	119	La Rochelle Gifts	France
2	146	Saveley & Henriot, Co.	France
3	171	Daedalus Designs Imports	France
4	172	La Corne D'abondance, Co.	France
5	209	Mini Caravy	France

	customerNumber	customerName	country
6	242	Alpha Cognac	France
7	250	Lyon Souvenirs	France
8	256	Auto Associés & Cie.	France
9	350	Marseille Mini Autos	France
10	353	Reims Collectables	France

Neo4j

```
In [6]: #
# Run this cell.
#
import neo4j_check
```

```
In [7]: #
# Set the neo4j user and password for connecting to your database. Th
# You set the password when you created the project and graph.
#
db_user = "neo4j"
db_password = "password"

neo4j_check.set_neo4j_connect_info(db_user, db_password)
```

```
In [8]: #
# Your database MUST have the Movie DB installed. You had to do this i
# the sample output.
#
res = neo4j_check.get_people_in_matrix()
res
```

Out[8]:

	name	born
0	Emil Eifrem	1978
1	Hugo Weaving	1960
2	Laurence Fishburne	1961
3	Carrie-Anne Moss	1967
4	Keanu Reeves	1964
5	Laurence Fishburne	1961
6	Carrie-Anne Moss	1967
7	Emil Eifrem	1978
8	Keanu Reeves	1964
9	Hugo Weaving	1960

MongoDB

```
In [9]: # Import the MongoDB test and helper functions.
#
import mongodb_check
```

```
In [10]: #
# Set the connection URL to get to your instance of MongoDB. You have
# in HW3 and when using Mongo Compass.
#
connect_url = "mongodb://localhost:27017/"
mongodb_check.set_connect_url(connect_url)
```

```
In [11]: #
# Run the following function. This will load information into MongoDB
#
df = mongodb_check.load_and_test_mongo()
df
```

Out[11]:

	_id	customerNumber	customerName	country
0	61c29b5c86f61fdf6b6bb381	103	Atelier graphique	France
1	61c29b5c86f61fdf6b6bb384	119	La Rochelle Gifts	France
2	61c29b5c86f61fdf6b6bb38e	146	Saveley & Henriot, Co.	France
3	61c29b5c86f61fdf6b6bb397	171	Daedalus Designs Imports	France
4	61c29b5c86f61fdf6b6bb398	172	La Corne D'abondance, Co.	France
5	61c29b5c86f61fdf6b6bb3a6	209	Mini Caravy	France
6	61c29b5c86f61fdf6b6bb3b0	242	Alpha Cognac	France
7	61c29b5c86f61fdf6b6bb3b3	250	Lyon Souvenirs	France
8	61c29b5c86f61fdf6b6bb3b4	256	Auto Associés & Cie.	France
9	61c29b5c86f61fdf6b6bb3d0	350	Marseille Mini Autos	France
10	61c29b5c86f61fdf6b6bb3d1	353	Reims Collectables	France
11	61c29b5c86f61fdf6b6bb3df	406	Auto Canal+ Petit	France

1. Database Core Concepts (10 points)

- There is a [Google Doc \(https://docs.google.com/document/d/1b0VVAS_LC25iMjIBx9zP9UhDg5eqsVoQ6h6Wdb68KwQ/edit?usp=sharing\)](https://docs.google.com/document/d/1b0VVAS_LC25iMjIBx9zP9UhDg5eqsVoQ6h6Wdb68KwQ/edit?usp=sharing).
- Make a copy of the Google Doc. Answer the questions in the document. You will submit a PDF of the document and your answers in the zip file you submit. The file must be in the folder and name **question1.pdf**.

2. Relational Algebra

Instructions

You will use the [online relational \(.\)](#) calculator to answer some of the subquestions. For these questions, your answer must contain:

- The text of the relational statement. The TAs may cut, paste and run the statement and it must work.
- An image showing the results of your execution.
- There is an example below.

Example

Question

- Use the "Silberschatz - UniversityDB" for this question.
- Professor Wu taught only one section. Produce the following information for the section.

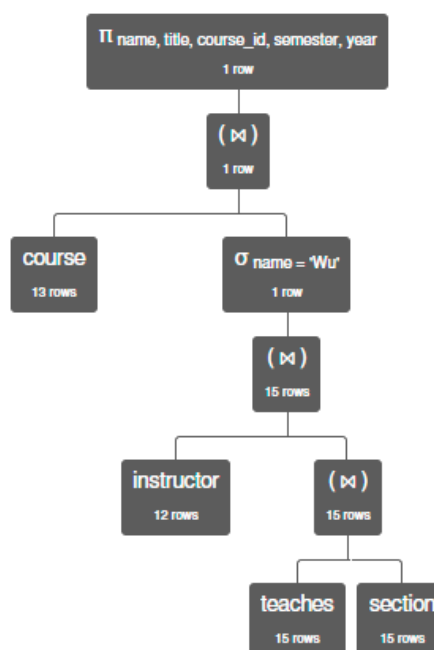
instructor.name	course.title	course.course_id	teaches.semester	teaches.year
'Wu'	'Investment Banking'	'FIN-201'	'Spring'	2010

Answer

```
 $\pi$  name, title, course_id, semester, year
  (course  $\bowtie$  (  $\sigma$  name='Wu' (instructor  $\bowtie$  (teaches  $\bowtie$  section)
  )))
```

execute selection

download history



```
 $\pi$  name, title, course_id, semester, year ( course  $\bowtie$  (  $\sigma$  name = 'Wu' ( instructor  $\bowtie$  ( teaches  $\bowtie$  section )
  )))
```

instructor.name	course.title	course.course_id	teaches.semester	teaches.year
'Wu'	'Investment Banking'	'FIN-201'	'Spring'	2010

2.1 Relation Model Schema (2 points)

Question

- The following is a simple MySQL table definition.

```
CREATE TABLE `new_table` (
  `product_category` INT NOT NULL,
  `produce_code` VARCHAR(45) NOT NULL,
  `product_name` VARCHAR(45) NULL,
  `product_description` VARCHAR(45) NULL,
  PRIMARY KEY (`product_category`, `produce_code`));
```

- Using the notation from chapter 2 slides for defining a relational schema, provide the corresponding relation schema definition.
 - Ignore the column types, NOT NULL, etc.
 - Two under-bar text, you can use $\underline{\text{cat}}$ to produce cat.

Answer (In Markdown cell below)

`new_table(product_category, produce_code, product_name, product_description)`

2.2 Relational Algebra (4 points)

Question

- Provide your answer following the examples' format.
- Use the Relax calculator the "Silberschatz - UniversityDB" for this question.
- A section r overlaps with another section l if and only if: They occurred at the same time (year, semester, time_slot_id).
- Produce the following table that shows the overlapping sections.

r.title	r.course_id	r.sec_id	l.title	l.course_id	l.sec_id	l.year	l.semester	l.time_slot_id
'Image Processing'	'CS-319'	2	'World History'	'HIS-351'	1	2010	'Spring'	'1'

Answer

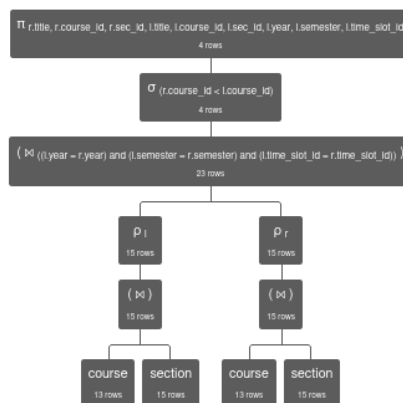
```


$$\pi_{r.title, r.course\_id, r.sec\_id, l.title, l.course\_id, l.sec\_id, l.year, l.semester, l.time\_slot\_id} (\sigma_{(r.course\_id < l.course\_id)} ((\rho_l (course \bowtie section)) \bowtie ((l.year=r.year) \wedge (l.semester=r.semester) \wedge (l.time\_slot\_id=r.time\_slot\_id)) (\rho_r (course \bowtie section))))$$


```

execute selection

download history



```


$$\pi_{r.title, r.course\_id, r.sec\_id, l.title, l.course\_id, l.sec\_id, l.year, l.semester, l.time\_slot\_id} (\sigma_{(r.course\_id < l.course\_id)} ((\rho_l (course \bowtie section)) \bowtie ((l.year = r.year) \wedge (l.semester = r.semester) \wedge (l.time\_slot\_id = r.time\_slot\_id)) (\rho_r (course \bowtie section))))$$


```

r.title	r.course_id	r.sec_id	l.title	l.course_id	l.sec_id	l.year	l.semester	l.time_slot_id
'Image Processing'	'CS-319'	1	'Investment Banking'	'FIN-201'	1	2010	'Spring'	'B'

2.3 Relational Algebra (4 points)

Question

- The relation algebra has additional operators for ordering, aggregation/group by, etc.
- A simple analysis of the data in the data in "Silberschatz - UniversityDB" shows that the takes table must be incomplete. Produce the following table, where sum_of_credits is the sum of a student's credits based on the information in takes.

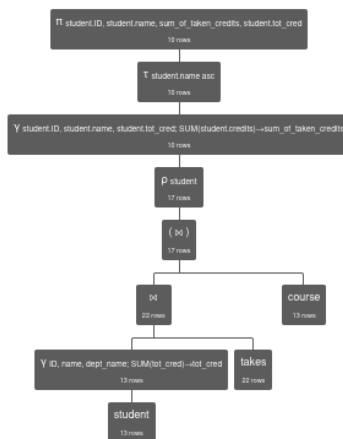
student_ID	student_name	sum_of_taken_credits	student.tot_cred
76653	'Aoi'	3	60
19991	'Brandt'	3	80
76543	'Brown'	7	58
23121	'Chavez'	3	110
44553	'Peltier'	4	56
55739	'Sanchez'	3	38
12345	'Shankar'	14	32
98988	'Tanaka'	8	120
54321	'Williams'	8	54
128	'Zhang'	7	102

Answer

```

π student.ID, student.name, sum_of_taken_credits, student.tot_cred (
  τ student.name asc (
    γ student.ID, student.name, student.tot_cred;
    sum(student.credits)→sum_of_taken_credits
  ) ρ student
  ((γ ID, name, dept_name; sum(tot_cred)→tot_cred student) ⋈ takes ⋈ course)))

```

[execute selection](#)
[download](#) [history](#)


```

π student.ID, student.name, sum_of_taken_credits, student.tot_cred (
  τ student.name asc (
    γ student.ID, student.name, student.tot_cred; SUM(student.credits)→sum_of_taken_credits
  ) ρ student
  ((γ ID, name, dept_name; SUM(tot_cred)→tot_cred student) ⋈ takes) ⋈ course
)

```

student.ID	student.name	sum_of_taken_credits	student.tot_cred
76653	'Aoi'	3	60
19991	'Brandt'	3	80
76543	'Brown'	7	58
23121	'Chavez'	3	110
44553	'Peltier'	4	56
55739	'Sanchez'	3	38
12345	'Shankar'	14	32
98988	'Tanaka'	8	120
54321	'Williams'	8	54
128	'Zhang'	7	102

3. SQL Query

Instructions and Example

You must follow and comply with the instructions for completing the questions in this section. Any deviation from the format is a score of 0.

1. The zip file you downloaded contains a file `question_2_sql.py`. The file contains:

- An example of the format and approach to answers.
- An empty function for each answer. You answer the question by completing the function's implementation.

2. The sample returns a Pandas data frame with `customerNumber`, `customerName` and `Country`. The country is a parameter to the function call.

```
In [12]: import question_3_sql
```

```
In [13]: #
# Call the function with the parameter France and display the results
#
result = question_3_sql.question_3_example_get_customers('France')
result
```

Out[13]:

	customerNumber	customerName	country
0	103	Atelier graphique	France
1	119	La Rochelle Gifts	France
2	146	Saveley & Henriot, Co.	France
3	171	Daedalus Designs Imports	France
4	172	La Corne D'abondance, Co.	France
5	209	Mini Caravy	France
6	242	Alpha Cognac	France
7	250	Lyon Souvenirs	France
8	256	Auto Associés & Cie.	France
9	350	Marseille Mini Autos	France
10	353	Reims Collectables	France
11	406	Auto Canal+ Petit	France

3.1 Revenue by Country (2 points)

Question

1. An `order` is a set of `orderdetails`.
 2. The value/revenue for an `orderdetails` is `priceEach*quantityOrdered`
 3. The value/revenue for an `order` is the sum of the value/revenue of the `orderdetails`.
- Implement the function `revenue_by_country`. We provide an example for the output. The company can only claim revenue if the order has `shipped`.

Answer

```
In [14]: result = question_3_sql.question_3_revenue_by_country()
result
```

Out[14]:

```
country SUM(revenue)
```

	country	SUM(revenue)
0	USA	3032204.26
1	Germany	196470.99
2	Norway	270846.30
3	Spain	947470.01
4	Denmark	176791.44
5	Italy	360616.81
6	Philippines	87468.30
7	UK	391503.90
8	Sweden	120457.09
9	France	965750.58
10	Belgium	91471.03
11	Singapore	263997.78
12	Austria	161418.16
13	Australia	509385.82
14	New Zealand	416114.03
15	Finland	295149.35
16	Canada	205911.86
17	Hong Kong	45480.79
18	Japan	167909.95
19	Ireland	49898.27

3.2 Customer Payments and Customer Purchases (2 points)

Question

1. `classicmodels.payments` records customer payments.
2. You can use the the formula above for computing the cost of an order.
3. The total owed by a customer is the total value/revenue for all orders. For the purposes of this problem, you should include all orders and not just the ones that shipped.
4. Implement the functions `purchases_and_payments`. The function returns a data frame with the following columns.
 - `customerNumber`
 - `customerName`
 - `total_spent` is the total value/cost over all orders by the customer.
 - `total_payments` is the total paid by the customer over all payments.
 - `total_unpaid` is the difference between `total_spent` and `total_payments`.

5. Order the result by `customerName`.
6. You must use at least one sub-query in your answer.

Answer

```
In [15]: #
# Execute this cell to display your answer.
#
result = question_3_sql.question_3_purchases_and_payments()
result
```

Out[15]:

	customerNumber	customerName	total_spent	total_payment	total_unpaid
0	242	Alpha Cognac	60483.36	60483.36	0.00
1	249	Amica Models & Co.	82223.23	82223.23	0.00
2	276	Anna's Decorations, Ltd	137034.22	137034.22	0.00
3	103	Atelier graphique	22314.36	22314.36	0.00
4	471	Australian Collectables, Ltd	55866.02	44920.76	10945.26
...
93	201	UK Collectables, Ltd.	106610.72	61167.18	45443.54
94	298	Vida Sport, Ltd	108777.92	108777.92	0.00
95	181	Vitachrome Inc.	72497.64	72497.64	0.00
96	144	Volvo Model Replicas, Co	66694.82	43680.65	23014.17
97	475	West Coast Collectables Co.	43748.72	43748.72	0.00

98 rows × 5 columns

3.3 What Customers Buy What? (1 point)

Question

1. Products are in `productLines`.
2. Produce a table that contains the `customerNumber` and `customerName` for all customers that have not orders a product from line `Planes` and not ordered a product from line `Trucks` and `Buses`.

Answer

```
In [16]: #  
# Run the cell below.  
#  
result = question_3_sql.question_3_customers_and_lines()  
result
```

Out[16]:

	customerNumber	customerName
0	103	Atelier graphique
1	406	Auto Canal+ Petit
2	187	AV Stores, Co.
3	219	Boards & Toys Co.
4	344	CAF Imports
5	171	Daedalus Designs Imports
6	362	Gifts4AllAges.com
7	357	GiftsForHim.com
8	350	Marseille Mini Autos
9	456	Microscale Inc.
10	209	Mini Caravy
11	486	Motor Mint Distributors Inc.
12	204	Online Mini Collectables
13	314	Petit Auto
14	112	Signal Gift Stores
15	259	Toms Spezialitäten, Ltd
16	201	UK Collectables, Ltd.
17	298	Vida Sport, Ltd

4 MongoDB

Instructions and Example

You must follow and comply with the instructions for completing the questions in this section. Any deviation from the format is a score of 0.

1. The final exam folder has a subdirectory `MongoDB` that contains MongoDB collections dumped in JSON format.
 - `actors_imdb.json`
 - `got_characters.json`
 - `got_episodes.json`
 - `imdb_titles.json`
 - `title_ratings.json`

2. Use MongoDB Compass:

- Create a MongoDB database F21_Final.
- Import the data from the files into collections. You can do this by using MongoDB Compass to create a collection, and then selecting the import data function.

3. You will implement your answers in functions in the file ``

2. The sample returns a data frame of the form (seasonNum, episodeNum, sceneNum, characterName) for the characters that appeared in season one, episode one.

```
In [18]: import question_4_mongo
```

```
In [19]: result = question_4_mongo.question_4_example()
result
```

Out[19]:

	seasonNum	episodeNum	sceneNum	characterName
0	1	1	1	Gared
1	1	1	1	Waymar Royce
2	1	1	1	Will
3	1	1	2	Gared
4	1	1	2	Waymar Royce
...
148	1	1	35	Summer
149	1	1	36	Bran Stark
150	1	1	36	Summer
151	1	1	36	Jaime Lannister
152	1	1	36	Cersei Lannister

153 rows × 4 columns

4.1 Implementing a JOIN (2 points)

Question

1. You will need to implement an aggregation for this problem. You can use MongoDB Compass to produce and test the aggregation, and then copy into the implementation template.
2. The aggregation operator `$lookup` implements a join-like function for MongoDB.
3. The aggregation operator (in a project) for getting substrings is `$substr`.
4. Write a query that joins episodes and ratings and produces a list of documents of the form:
 - seasonNum, episodeNum, episodeTitle, episodeDescription,

- episodeDate from got_episodes.
- tconst, averageRating, numVotes from title ratings.

Answer

```
In [20]: #
# Run your test here.
#
result = question_4_mongo.question_4_ratings()
result
```

Out[20]:

	tconst	averageRating	numVotes	seasonNum	episodeNum	episodeTitle	episodeDesc
0	tt1480055	9.1	44596	1	1	Winter Is Coming	Jon Arryn, th of the King, i
1	tt1668746	8.8	33936	1	2	The Kingsroad	While Bran r from his f ta
2	tt1829962	8.7	32139	1	3	Lord Snow	Lord Stark daughters i k
3	tt1829963	8.8	30562	1	4	Cripples, Bastards, and Broken Things	Eddard inve Jon Arryn's J
4	tt1829964	9.1	31777	1	5	The Wolf and the Lion	Cate captured Tyr plans to
...
68	tt6027908	7.8	126920	8	2	A Knight of the Seven Kingdoms	The I Win apprc
69	tt6027912	7.4	210187	8	3	The Long Night	The Night K his arr arrived
70	tt6027914	5.4	160721	8	4	The Last of the Starks	In the w costly victi and L
71	tt6027916	5.9	187221	8	5	The Bells	Daene Cersei wei options i
72	tt6027920	4	238095	8	6	The Iron Throne	In the after the dev atta

73 rows × 8 columns

4.2 Just Kidding

- We did not spend a lot of time on MongoDB and that previous query was not fun.

- So, 4.1 is actually with 5 points and you are done with MongoDB. For now.

5 Neo4j

Instructions and Example

You must follow and comply with the instructions for completing the questions in this section. Any deviation from the format is a score of 0.

1. You will use the Movie Graph for this question.
2. Implement the answers in functions in the Python file

The example function returns a table with information about which people directed Tom hanks in which movies.

```
In [21]: import question_5_neo4j
```

```
In [22]: result = question_5_neo4j.directed_tom_hanks()
result
```

Out[22]:

	0	1	2
0	Tom Hanks	You've Got Mail	Nora Ephron
1	Tom Hanks	Sleepless in Seattle	Nora Ephron
2	Tom Hanks	Joe Versus the Volcano	John Patrick Stanley
3	Tom Hanks	That Thing You Do	Tom Hanks
4	Tom Hanks	Cloud Atlas	Tom Tykwer
5	Tom Hanks	Cloud Atlas	Andy Wachowski
6	Tom Hanks	Cloud Atlas	Lana Wachowski
7	Tom Hanks	The Da Vinci Code	Ron Howard
8	Tom Hanks	The Green Mile	Frank Darabont
9	Tom Hanks	Apollo 13	Ron Howard
10	Tom Hanks	Cast Away	Robert Zemeckis
11	Tom Hanks	Charlie Wilson's War	Mike Nichols
12	Tom Hanks	The Polar Express	Robert Zemeckis
13	Tom Hanks	A League of Their Own	Penny Marshall
14	Tom Hanks	You've Got Mail	Nora Ephron
15	Tom Hanks	Sleepless in Seattle	Nora Ephron
16	Tom Hanks	Joe Versus the Volcano	John Patrick Stanley
17	Tom Hanks	That Thing You Do	Tom Hanks
18	Tom Hanks	Cloud Atlas	Tom Tykwer

	0	1	2
19	Tom Hanks	Cloud Atlas	Lilly Wachowski
20	Tom Hanks	Cloud Atlas	Lana Wachowski
21	Tom Hanks	The Da Vinci Code	Ron Howard
22	Tom Hanks	The Green Mile	Frank Darabont
23	Tom Hanks	Apollo 13	Ron Howard
24	Tom Hanks	Cast Away	Robert Zemeckis
25	Tom Hanks	Charlie Wilson's War	Mike Nichols
26	Tom Hanks	The Polar Express	Robert Zemeckis

5.1 People Who Directed Themseves (2 points)

Question

- Implement the function `people_who_directed_themselves`.
- The format of the answer is a data frame of the form `(name, title, name)` where the person `ACTED_IN` and `DIRECTED` the movie.

Answer

```
In [24]: #
# Test you answer
#
result = question_5_neo4j.directed_themselves()
result
```

Out[24]:

	0	1	2
0	Tom Hanks	That Thing You Do	Tom Hanks
1	Clint Eastwood	Unforgiven	Clint Eastwood
2	Danny DeVito	Hoffa	Danny DeVito

5.2 People Who Reviewed the same Movie (3 points)

Question

- Implement the function `both_reviewed(person_1_name, person_2_name)`
- The function returns a data frame of the form `person_1_name, movie_title, person_2_name` if the two people with the names reviewed the movie.
- Test you answer with the names below. You cannot hard code names in your query.

Answer

```
In [25]: #
result = question_5_neo4j.both_reviewed('James Thompson', 'Jessica Th
result
```

Out[25]:

	0	1	2
0	James Thompson	The Replacements	Jessica Thompson
1	James Thompson	The Da Vinci Code	Jessica Thompson

6 Data Modeling – RACI

Question

- [RACI](https://www.softwareadvice.com/resources/what-is-a-raci-chart/) (<https://www.softwareadvice.com/resources/what-is-a-raci-chart/>) is an acronym for an approach to defining the relationships between people/stakeholders and a project.
- For this question, you will:
 - Do a Crow's Foot ER diagram defining a data model for representing RACI.
 - Create a SQL schema to represent the tables, constraints, etc. that you determine are necessary.
- The core entity types are:
 - Project(project_id, project_name, start_date, end_date)
 - Person(UNI, last_name, first_name, email)
- Implementing RACI is about understanding relationships between people and projects. The table below explains the concept.

Role	Description
Responsible	Who is responsible for doing the actual work for the project task.
Accountable	Who is accountable for the success of the task and is the decision-maker. Typically the project manager.*
Consulted	Who needs to be consulted for details and additional info on requirements. Typically the person (or team) to be consulted will be the subject matter expert.
Informed	Who needs to be kept informed of major updates. Typically senior leadership.

- There are two constraints:
 - There is exactly one person who is Accountable for a project.
 - A specific person can have at most one relationship to a project, for example "Bob" cannot be both Consulted and Informed for the same project.
- To answer this question, you must:
 - Draw the Crow's Foot ER diagram using LucidChart.
 - Create a database schema implementing the data model you define.
- You do not need to populate the data model with data or query the data, but YOU MUST execute your DDL statements.

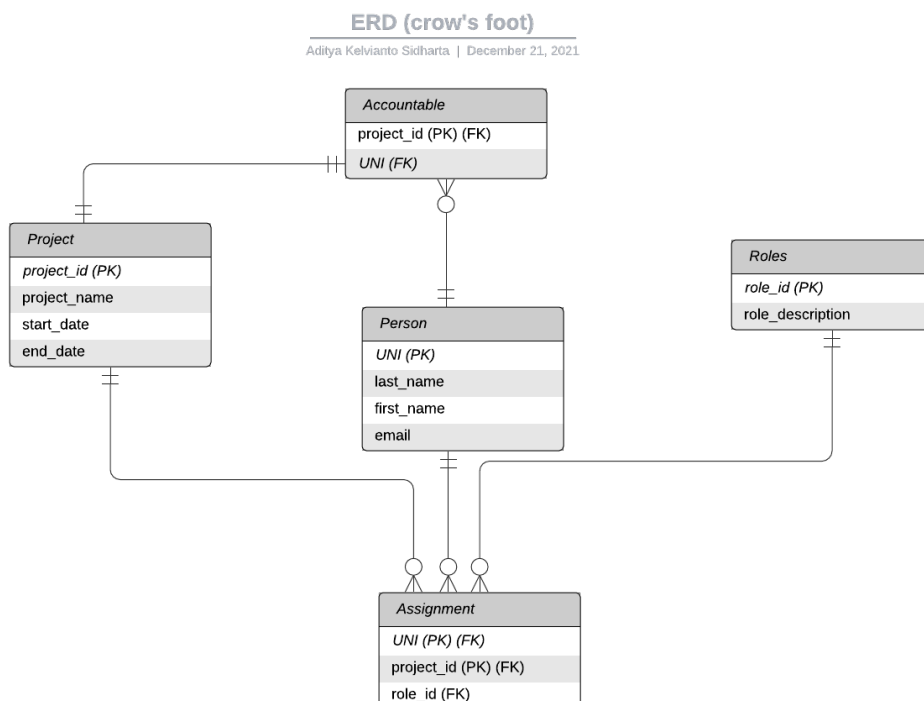
- You execute the DDL statements implementing functions in the file `question_6_schema`.
- You may use DataGrip or other tools to design the schema and test your statements, but for the final answer you must have one function in `question_6_schema` for each DDL statement and you must execute each function in a cell below.
- Name your database schema `RACI`.
- There is no single, correct answer. Document any assumptions or design decisions that you make.

Design Decisions and Assumptions

Document any design decisions or assumptions that you make.

- Each of the Project ID must appear at least once in the Accountable Table (Because there must be at least one person that is Accountable for each project. Therefore, we are using the `project_id` as our primary key.
- A person might or might not be assigned to a project at all. If he/she is assigned to a project, he or she is only able to take on one role within a single project (Responsible, Accountable, Consulted, or Informed). However, a person is able to take on one role in zero, one or more projects. This is done by setting both `UNI` and `project_id` as composite key for assignment table
- The only requirement is that there is only one person who is Accountable for a project. In other words, in a project, its possible for zero, one, or more than one people become responsible / consulted or informed.

ER Diagram



Schema Creation

```
In [32]: #  
#  
import question_6_schema
```

```
In [30]: #  
# Execute each function in a single cell.  
#  
res = question_6_schema.schema_operation_1()  
res
```

Out[30]: 0

```
In [33]: #  
# Execute each function in a single cell.  
#  
res = question_6_schema.schema_operation_2()  
res
```

Out[33]: 0

```
In [34]: #  
# Execute each function in a single cell.  
#  
res = question_6_schema.schema_operation_3()  
res
```

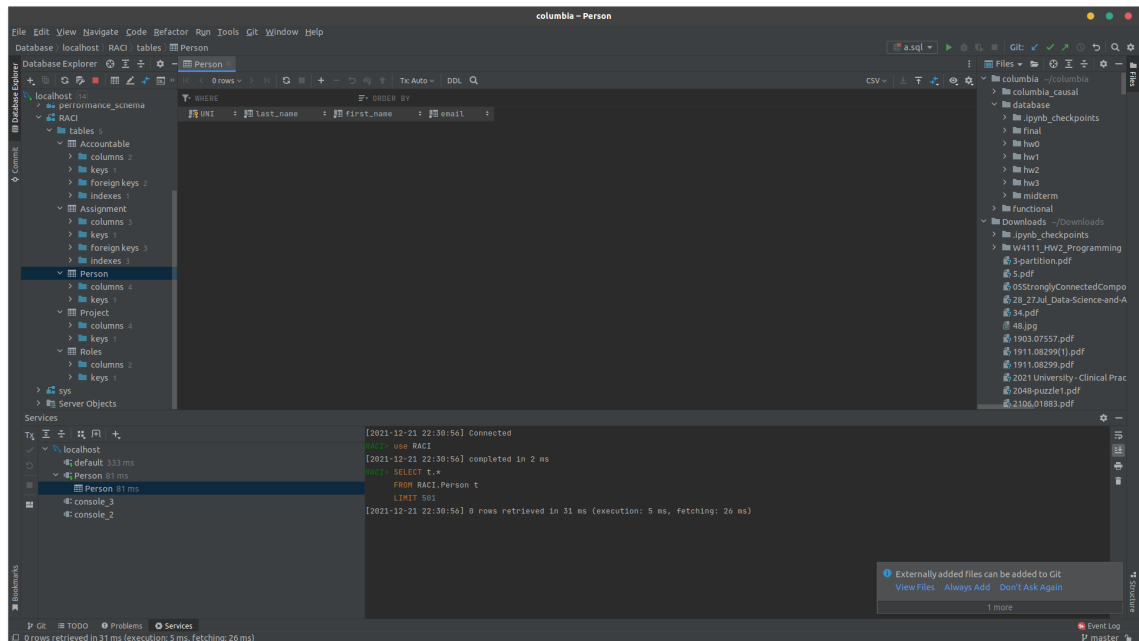
Out[34]: 0

```
In [35]: #  
# Execute each function in a single cell.  
#  
res = question_6_schema.schema_operation_4()  
res
```

Out[35]: 0

```
In [36]: #
# Execute each function in a single cell.
#
res = question_6_schema.schema_operation_5()
res
```

Out[36]: 0



7. Data Transformation

Question

- In this question, you will produce a star schema and populate with data from `classicmodels`.
- A star schema has a fact table and dimensions. The core fact is:
 - A customer (`customerNumber`)
 - Some quantity of a product (`quantityOrdered`) at a price (`priceEach`)
 - On a given date (`orderDate`)
- We will consider three dimensions:
 - date is (month, quarter, year)
 - location is the dimension representing where the customer is and is of the form (city, country, region). Region is one of (EMEA, NA, AP).
 - USA and Canada are in NA.
 - Philippines, Hong Kong, Singapore, Japan, Australia and New Zealand are in AP
 - All other countries are in EMEA.
 - product_type is (scale, product line).
- You will follow the same approach for implementation as for question 6.

- There is an implementation template `question_7_sql`. You will implement three sets of SQL operations.
 - The functions of the form `schema_operation_n()` implement creating the star schema, tables, constraints, etc. There is one function for each statement. Name your schema `classicmodels_star`
 - The functions `data_transformation_n()` contain SQL statements for loading the `classicmodels_star` schema. You can have at most 3 SQL statement per function.
 - There are three queries you must implement:
 - `sales_by_year_region()` returns the total value of orders broken down by region and year.
 - `sales_by_quarter_year_county_region()` drills down to show the same information expanded to include quarter and year.
 - `sales_by_product_line_scale_year()` shows sales by product line, product scale and year.

Answer

In the following cells, execute your various functions that invoke SQL.

```
In [44]: #  
#  
import question_7_sql
```

```
In [51]: question_7_sql.schema_operation_1()
```

```
Out[51]: 0
```

```
In [52]: question_7_sql.schema_operation_2()
```

```
Out[52]: 0
```

```
In [53]: question_7_sql.schema_operation_3()
```

```
Out[53]: 0
```

```
In [54]: question_7_sql.schema_operation_4()
```

```
Out[54]: 0
```

```
In [ ]: question_7_sql.data_transformation_1()
```

```
In [ ]: question_7_sql.data_transformation_2()
```

```
In [ ]: question_7_sql.data_transformation_3()
```

```
In [ ]: question_7_sql.data_transformation_4()
```

```
In [59]: question_7_sql.sales_by_year_region()
```

```
Out[59]:
```

	year	region	SUM(sales)
0	2003	NA	1.228194e+07
1	2003	EMEA	1.665505e+07
2	2003	AP	6.088170e+06
3	2003	None	1.011002e+05
4	2004	EMEA	2.523000e+07
5	2004	NA	1.961989e+07
6	2004	AP	7.960323e+06
7	2004	None	4.154632e+05
8	2005	NA	6.380509e+06
9	2005	EMEA	9.036608e+06
10	2005	AP	3.221823e+06
11	2005	None	3.279811e+05

```
In [60]: question_7_sql.sales_by_quarter_year_county_region()
```

```
Out[60]:
```

	quarter	year	country	region	SUM(sales)
0	1	2003	USA	NA	4.659626e+05
1	1	2003	Germany	EMEA	5.274505e+04
2	1	2003	Norway	EMEA	2.510947e+05
3	1	2003	Spain	EMEA	2.010310e+05
4	1	2003	Denmark	EMEA	1.618776e+05
...
121	2	2005	Spain	EMEA	1.257368e+06
122	2	2005	Australia	AP	6.383612e+05
123	2	2005	Italy	EMEA	4.243471e+05
124	2	2005	Austria	EMEA	6.290408e+05
125	2	2005	Belgium	EMEA	1.031728e+05

126 rows × 5 columns

```
In [61]: question_7_sql.sales_by_product_line_scale_year()
```

```
Out[61]:
```

	year	scale	SUM(sales)
0	2003	1:18	1.550907e+07
1	2003	1:24	6.651887e+06
2	2003	1:10	2.645018e+06
3	2003	1:12	4.186672e+06

	year	scale	SUM(sales)
4	2003	1:32	1.736587e+06
5	2003	1:50	9.639826e+05
6	2003	1:700	2.620276e+06
7	2003	1:72	8.127648e+05
8	2004	1:12	6.462946e+06
9	2004	1:18	2.216434e+07
10	2004	1:24	9.873237e+06
11	2004	1:50	1.461928e+06
12	2004	1:700	4.538730e+06
13	2004	1:72	1.336395e+06
14	2004	1:10	4.568298e+06
15	2004	1:32	2.819803e+06
16	2005	1:10	1.716769e+06
17	2005	1:12	2.205585e+06
18	2005	1:18	7.931748e+06
19	2005	1:24	3.649010e+06
20	2005	1:32	1.049892e+06
21	2005	1:50	5.086841e+05
22	2005	1:700	1.436068e+06

The screenshot shows a Jupyter Notebook interface with a SQL query executed in a database. The query is a SELECT statement with a JOIN and a WHERE clause. The results are displayed in a table with columns: date_id, month, quarter, and year. The status bar at the bottom indicates 265 rows retrieved.

```
USE classicmodels_star;
SELECT t.*
FROM classicmodels_star.date t
WHERE t.date_id = 1;
```

date_id	month	quarter	year
1	1	1	2003
2	2	1	2003
3	3	1	2003
4	4	1	2003
5	5	1	2003
6	6	2	2003
7	7	2	2003
8	8	2	2003
9	9	3	2003
10	10	3	2003
11	11	3	2003
12	12	3	2003
13	1	1	2004
14	2	1	2004
15	3	1	2004
16	4	2	2004
17	5	2	2004
18	6	2	2004
19	7	3	2004
20	8	3	2004
21	9	3	2004
22	10	4	2004

The screenshot shows a Jupyter Notebook interface with a SQL query executed in a database. The query is a SELECT statement with a JOIN and a WHERE clause. The results are displayed in a table with columns: order_id, customer_number, quantity_ordered, price_each, product_type_id, location_id, date_id, order_date, and order_number. The status bar at the bottom indicates 1500 rows retrieved.

```
USE classicmodels_star;
SELECT t.*
FROM classicmodels_star.orders t
WHERE t.order_id = 1;
```

order_id	customer_number	quantity_ordered	price_each	product_type_id	location_id	date_id	order_date	order_number
1	363	30	136	12	62	5	2003-01-06	10100
2	363	30	136	12	62	4	2003-01-06	10100
3	363	30	136	12	62	3	2003-01-06	10100
4	363	30	136	12	62	2	2003-01-06	10100
5	363	30	136	12	62	1	2003-01-06	10100
6	363	50	55.09	12	62	5	2003-01-06	10100
7	363	50	55.09	12	62	4	2003-01-06	10100
8	363	50	55.09	12	62	3	2003-01-06	10100
9	363	50	55.09	12	62	2	2003-01-06	10100

The screenshot shows a Jupyter Notebook with a SQL query and its results. The query is:

```
use classicmodels.star
SELECT t.*
FROM classicmodels.star.fact t
LIMIT 501
```

The results are displayed in a table with 10 columns: `location_id`, `city`, `country`, `region`, `product_type_id`, `scale`, `product_line`, `date`, `fact`, and `fact`. The table contains 101 rows of data, with the first row being:

location_id	city	country	region	product_type_id	scale	product_line	date	fact	fact
1	Aachen	Germany	EMEA	1	1:18	Classic Cars	2003-01-06	10180	10180

The Services panel shows the execution of the query, with a message: "265 rows retrieved starting from 1 in 62 ms (execution: 3 ms, fetching: 59 ms)".

The screenshot shows a Jupyter Notebook with a SQL query and its results. The query is:

```
use classicmodels.star
SELECT t.*
FROM classicmodels.star.location t
LIMIT 501
```

The results are displayed in a table with 5 columns: `location_id`, `city`, `country`, `region`, and `product_type_id`. The table contains 101 rows of data, with the first row being:

location_id	city	country	region	product_type_id
1	Aachen	Germany	EMEA	1

The Services panel shows the execution of the query, with a message: "95 rows retrieved starting from 1 in 29 ms (execution: 5 ms, fetching: 24 ms)".

The screenshot shows a Jupyter Notebook with a SQL query and its results. The query is:

```
use classicmodels.star
SELECT t.*
FROM classicmodels.star.product_type t
LIMIT 501
```

The results are displayed in a table with 7 columns: `product_type_id`, `scale`, `product_line`, `date`, `fact`, `fact`, and `fact`. The table contains 101 rows of data, with the first row being:

product_type_id	scale	product_line	date	fact	fact	fact
1	1:18	Classic Cars	2003-01-06	10180	10180	10180

The Services panel shows the execution of the query, with a message: "30 rows retrieved starting from 1 in 37 ms (execution: 4 ms, fetching: 33 ms)".

In []:

