

Zhimin Wang zhiminw2

CS 411: Database Systems

Fall 2021

Homework 4 (Due by 23:59 CT on November 12th)

Logistics

1. This homework is due on November 12th at 23:59 CT. **We DO NOT accept late homework submissions.**
2. You will be using Gradescope to submit your solutions. Answer each sub-question (e.g. "a.") on a **new page** and submit your solution as a single PDF file on Gradescope.
IMPORTANT: Please make sure to link PDF pages with the corresponding question outline on GradeScope.
3. Please write down any **intermediate steps** to receive full credit.
4. Keep your solutions brief and clear.
5. Please use Campuswire if you have questions about the homework but do not post answers. Feel free to come to office hours.

Q1. Relational Algebra Queries (20 points)

Consider a relational database of a Medical College/Hospital system. The corresponding relational schema is described below:

Doctor (DoctorID, Name, DepartmentID, Phone, City, salary)

Patient (PatientID, Patient_Name, Gender, Phone, City, dob)

Appointments(Date, PatientID, DoctorID)

Courses(CourseID, Course_Name, DoctorID, DepartmentID)

Department (DepartmentID, Dept_Name, Building)

Please answer the following queries using **relational algebra expressions**: (do not use expression tree).

You can solve it in steps and combine your answers. For example,

R1: $A \bowtie_{A.abc=B.abc} B$

R2: $R1 \bowtie C$

R3: $\Pi_x(\sigma_{x=y} R2)$

Note: You can use Relax: Relational Algebra Calculator to try your queries:

<https://dbis-uibk.github.io/relax/calc/local/uibk/local/Q>. But please note that Relax is not meant for grading your answers. It is a tool that can help you write RA queries and visualize your answer.

We have provided the schema and instance of Question 1 here:

https://docs.google.com/document/d/1LfRoVz373b_irQQzaFwB71igvModsW9ulcpXik4ecl0

You can paste the above script into 'Group Editor' and execute your queries in the 'Relational algebra' section to test your RA queries. Please watch [this](#) short tutorial (made by Abdu) on how to load the data and use Relax to write RA queries.

1. Find the names of patients who have not yet made an appointment. (5 points)
2. Find the names of all the departments that had zero male patients on 10/28/2021.
Please note that the patient belongs to the same department as that of the doctor assigned to them. Ignore departments that did not have any patients on that day. (5 points)
3. Find all the dates when there was at least one appointment having either a patient from Nashville or a doctor who has never taught a course. (10 points)

Q2. Relational Algebra Equivalence (20 points)

Consider the following supermarket database schema, with primary key(s) underlined:

Doctor (DoctorID, Name, DepartmentID, Phone, City, salary)

Patient (PatientID, Patient_Name, Gender, Phone, City, dob)

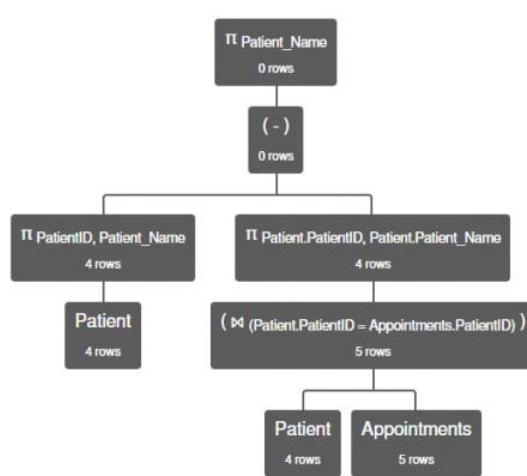
Appointments(Date, PatientID, DoctorID)

Courses(CourseID, Course_Name, DoctorID, DepartmentID)

Department (DepartmentID, Dept_Name, Building)

Q1:

1. $\pi \text{Patient_Name} ((\pi \text{PatientID}, \text{Patient_Name} \text{ Patient}) - (\pi \text{Patient.PatientID}, \text{Patient.Patient_Name} (\text{Patient} \bowtie (\text{Patient.PatientID} = \text{Appointments.PatientID}) \text{ Appointments}))$



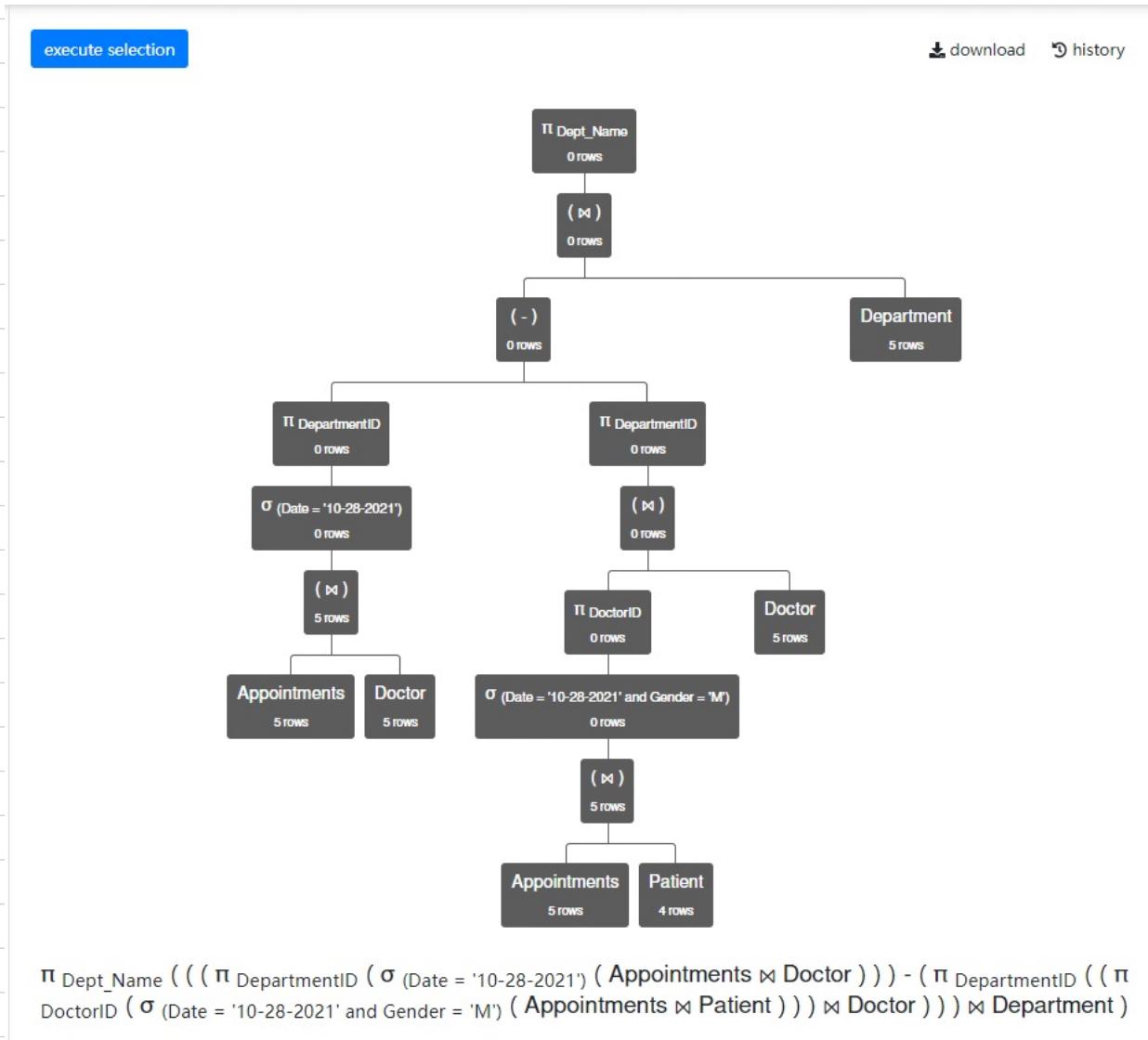
$\pi \text{Patient_Name} ((\pi \text{PatientID}, \text{Patient_Name} \text{ Patient}) - (\pi \text{Patient.PatientID}, \text{Patient.Patient_Name} (\text{Patient} \bowtie (\text{Patient.PatientID} = \text{Appointments.PatientID}) \text{ Appointments}))$

Patient.Patient_Name

1

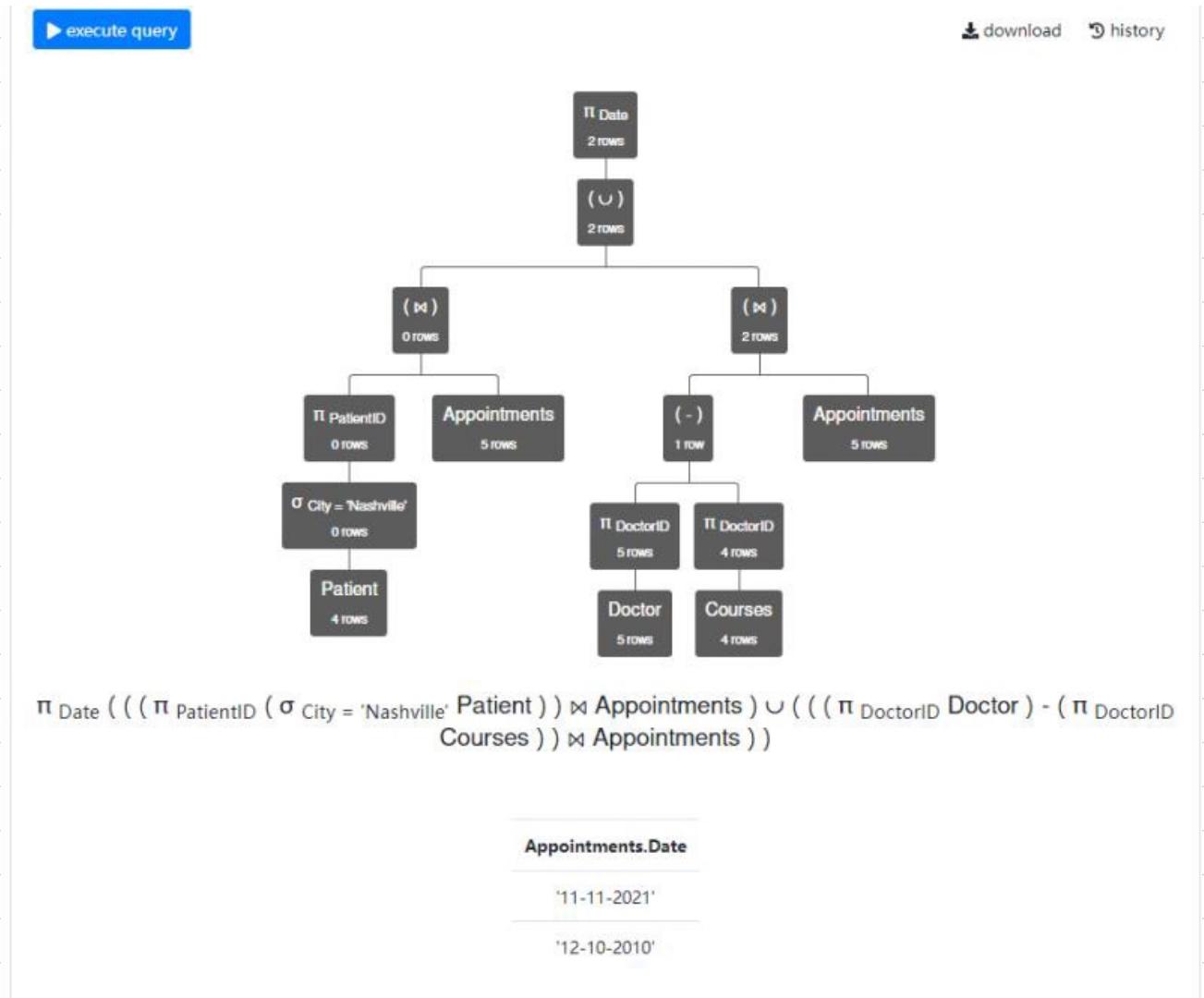
Q1:

2. $\pi_{Dept_Name} (((\pi_{DepartmentID} (\sigma_{(Date = '10-28-2021')} (Appointments \bowtie Doctor))) - (\pi_{DepartmentID} ((\pi_{DoctorID} (\sigma_{(Date = '10-28-2021' \wedge Gender = 'M')} (Appointments \bowtie Patient))) \bowtie Doctor))) \bowtie Department)$



Q 1

3: $\pi \text{ Date } (((\pi \text{ PatientID } (\sigma \text{ City} = \text{'Nashville'} \text{ Patient})) \bowtie \text{ Appointments}) \cup (((\pi \text{ DoctorID } \text{ Doctor}) - (\pi \text{ DoctorID } \text{ Courses})) \bowtie \text{ Appointments}))$



For each pair of relational algebra expressions (E1 and E2) shown below, state whether or not the two expressions are equivalent in general. If your answer is "true", justify the equivalency by explaining which RA rule(s) can be used to transform one query expression into the other. If your answer is "false", give a sample instance of the database where the two expressions are not equal. Make no assumptions about the existence of foreign keys.

1. E1: $\Pi_{\text{Doctor.Name}, \text{Doctor.DepartmentID}}(\text{Doctor} \bowtie_{\text{Doctor.DepartmentID} = \text{Courses.DepartmentID}} \text{Courses})$

E2: $\sigma_{\text{Doctor.DepartmentID} == \text{Courses.DepartmentID}}(\Pi_{\text{Doctor.Name}, \text{Doctor.DepartmentID}}(\text{Doctor}) \times \Pi_{\text{Courses.DepartmentID}}(\text{Courses}))$ (5 points)

2. E1: $(\Pi_{\text{PatientID}}(\sigma_{\text{Patient.City} == 'Chicago'} \text{Patient}) \bowtie \text{Appointments}) \cup (\text{Appointment} \bowtie \Pi_{\text{DoctorID}}(\sigma_{\text{Doctor.City} == 'Nashville'} \text{Doctor}))$

E2: $\Pi_{\text{Appointments.PatientID}, \text{Appointments.DoctorID}, \text{Appointments.Date}}((\sigma_{\text{Patient.City} == 'Chicago'} \text{Patient} \bowtie \text{Appointments}) \bowtie \sigma_{\text{Doctor.City} == 'Nashville'} \text{Doctor})$ (5 points)

3. E1: $\Pi_{\text{DoctorID}}(\sigma_{\text{Doctor.City} == 'Nashville'} \text{Doctor}) - \Pi_{\text{Doctor.DoctorID}}(\text{Doctor} \bowtie_{\text{Doctor.DoctorID} = \text{Courses.DoctorID}} \text{Courses})$

E2: $\Pi_{\text{DoctorID}} \text{Doctor} - \Pi_{\text{DoctorID}}(\sigma_{\text{Doctor.City} != 'Nashville'} \text{Doctor}) - \Pi_{\text{DoctorID}}((\sigma_{\text{Doctor.City} == 'Nashville'} \text{Doctor}) \bowtie_{\text{Doctor.DoctorID} = \text{Courses.DoctorID}} \text{Courses})$ (10 points)

Q3. Physical Operators

In this problem, you may choose between two options, both of which will provide you with full credit.

If you **choose to do the MP option** and implement it successfully, you will **receive 10 extra credit** points.

You should not submit solutions for both options, just do one.

Q₂

1: False.

That's how join works.

$$\sigma_c(A \times B) = A \bowtie_{(c)} B$$

But the output here is slightly different. For E₁, it will return a table with attributes $\begin{matrix} \text{Name} \\ \text{Doctor.} \end{matrix}$ and $\begin{matrix} \text{DepartmentID} \\ \text{Doctor.} \end{matrix}$. For E₂, it will return a table with three attributes: Doctor.Name, Doctor.DepartmentID, Course.DepartmentID

Q2.

2. False.

whose

For E_1 , it is like select rows \checkmark Patient . City = Chicago OR Doctor . City = Nashville

But for E_2 , its Patient . City = Chicago AND Doctor . City = Nashville

Example :

```
Doctor = { DoctorID, Name, DepartmentID, Phone, City, salary
01, John, D10, 777-777-7771, Chicago, 60000
02, Amelia, D11, 747-777-5772, Nashville, 60000
03, Joe, D12, 767-787-8773, Champaign, 50000
04, Triv, D14, 787-797-0792, Nashville, 80000
05, Rose, D13, 587-897-0991, SanJose, 90000
}
```

```
Department = { DepartmentID, Dept_Name, Building
D10, Cardiology, 2009
D11, ER, 2010
D12, Orthopedic, 2020
D13, Gynecology, 2021
D14, General, 2015
}
```

```
Patient = { PatientID, Patient_Name, Gender, Phone, City, dob
11111, Tori, F, 557-897-9995, Chicago, 09-04
11112, Thomas, M, 555-555-9999, Champaign, 08-04
14413, Cary, F, 555-555-8888, Champaign, 06-04
55554, Michelle, M, 444-444-4444, Champaign, 09-07
}
```

```
Appointments =
AptID, PatientID, DoctorID Date
624, 11111, 02, 01-11-2012 E2's result E1's result
625, 14413, 04, 04-13-2011
145, 11112, 01, 11-11-2021 E1's result
644, 55554, 03, 11-09-2020
345, 14413, 01, 12-10-2010
}
```

```
Courses = { CourseID, Course_Name, DoctorID, DepartmentID
C1, Ortho101, 02, D12
C2, Anatomy101, 04, D10
C3, Pharmacology, 03, D14
C4, Gynae, 05, D13
}
```

}, True.

These two parts are equivalent:

$$\Pi_{\text{DoctorID}}(\sigma_{\text{Doctor.City} = \text{'Nashville'}} \text{Doctor}) \subseteq \Pi_{\text{DoctorID}} \text{Doctor} - \Pi_{\text{DoctorID}}(\sigma_{\text{Doctor.City} \neq \text{'Nashville'}} \text{Doctor})$$

$$\Pi_a(\sigma_c A) = \Pi_a A - \Pi_a(G_{1c} A)$$

I don't know what RA rules can it apply.
It's just like $A \cup B = C$ then $A = C - B$

The rest part, in E_1 , it removes doctor who has taught course.

In E_2 , it removes doctor ^{who} has taught courses and live in Nashville.

But since the previous step has already remove all doctors who live in other cities,

it's like $\sigma_c(R-S) = G_c(R) - S = G_c(R) - G_c(S)$ in such case.

Option 1: Implement a Physical Operator (10 points + 10 EC = 20 points)

Implement the one-pass, sort-based INTERSECTION physical operator using Python. A skeleton code file, along with test cases, can be found [here](#).

The **intersection** operator should combine tables with the same columns (order of columns and names of columns must be the same), as described in the lecture notes and textbook. I/O operations are handled by our code skeleton, you are only responsible for the Set INTERSECTION operator.

For each test case provided, we have three input files and one expected output file which can be used to check the correctness of your output. You can use **diff** to compare your output and the expected output:

- **config.txt**: define memory size and block size.
- **table1.csv**: input table 1
- **table2.csv**: input table 2
- **expectedOutput.csv**: the results of **table1 INTERSECTION table2**.

Output Format:

1. If the intersection cannot be carried out under the constraints outlined in **config.txt**, please write "INVALID MEMORY" or "INVALID SCHEMA/INPUT" (as shown in the skeleton) to the first cell in the output file.
2. Please include the column names and the data of the resulting table. The rows in the output table should maintain the same relative orders as the bigger input table.

Submission:

Please submit your code as a submission for "Homework 4 Q3 Option 1" (**TODO: Attach link**) on Gradescope, as a script called **mp.py**. We will run your code on each of the provided test cases, along with a number of held-out test inputs in order to verify correctness.

Option 2: Two-Pass "Multi-Way" Merge Sort (10 points)

Consider the following relation R with 12 blocks of data stored on the disk ($B(R) = 12$):

[63, 33, 89] [98, 21, 10] [78, 89, 51] [80, 100, 98] [66, 18, 36] [17, 80, 67]

[42, 5, 19] [34, 22, 52] [7, 96, 98] [10, 76, 88] [31, 35, 7] [8, 13, 88]

Each block holds 3 values and there are 4 blocks of memory available for the operation ($M = 4$). Use Two-pass "multi-way" merge sort algorithm to sort the blocks above.

Solution Format:

- Please write one memory update per line and don't leave empty lines.

Q } ;

option 2 ;

[63,33,89][98,21,10][78,89,51][80,100,98]=>[10,21,33][51,63,78][80,89,89][98,98,100]
[66,18,36][17,80,67][42,5,19][34,22,52]=>[5,17,18][19,22,34][36,42,52][66,67,80]
[7,96,98][10,76,88][31,35,7][8,13,88]=>[7,7,8][10,13,31][35,76,88][88,96,98]
[10,21,33][5,17,18][7,7,8]=>[5,7,7]=>out
[10,21,33][5,17,18][7,7,8]=>[8,_,_]
[10,21,33][5,17,18][10,13,31]=>[8,10,10]=>out
[10,21,33][5,17,18][10,13,31]=>[13,17,18]=>out
[10,21,33][19,22,34][10,13,31]=>[19,21,22]=>out
[10,21,33][19,22,34][10,13,31]=>[31,_,_]
[10,21,33][19,22,34][35,76,88]=>[31,33,_]
[51,63,78][19,22,34][35,76,88]=>[31,33,34]=>out
[51,63,78][36,42,52][35,76,88]=>[35,36,42]=>out
[51,63,78][36,42,52][35,76,88]=>[51,52,_]
[51,63,78][66,67,80][35,76,88]=>[51,52,63]=>out
[51,63,78][66,67,80][35,76,88]=>[66,67,76]=>out
[51,63,78][66,67,80][35,76,88]=>[78,_,_]
[80,89,89][66,67,80][35,76,88]=>[78,80,80]=>out
[80,89,89][66,67,80][35,76,88]=>[88,_,_]
[80,89,89][66,67,80][88,96,98]=>[88,88,89]=>out
[80,89,89][66,67,80][88,96,98]=>[89,_,_]
[98,98,100][66,67,80][88,96,98]=>[89,96,98]=>out
[98,98,100][66,67,80][88,96,98]=>[98,98,100]=>out

- First, write sorted runs and then merge sorted runs.
- For sorted runs, write as: [1,3][5,4][7,9]=>[1,3][4,5][7,9]
- For every merge run, mark in bold the position of pointers in each block at the beginning of the run. Eg: [1,3]**[2,4]**=>[1,2]=>out
- If two slots have the same number, choose the leftmost first by default.
For example: [2,3]**[2,4]**=>**[2,2]**=>out
- If the output block cannot be outputted and needs a memory reload, write as
[1,2]**[2,4]**=>[2,_]

Q4. Block-based Nested Loop Join (20 points)

The following three tables from a recruiting database of a Company are stored on the disk:

Candidates (CandidateID, FirstName, LastName, PhoneNumber)

WorkExperience (CandidateID, CompanyID, Role, Years)

Companies (CompanyID, CompanyName, StateName, CountryName)

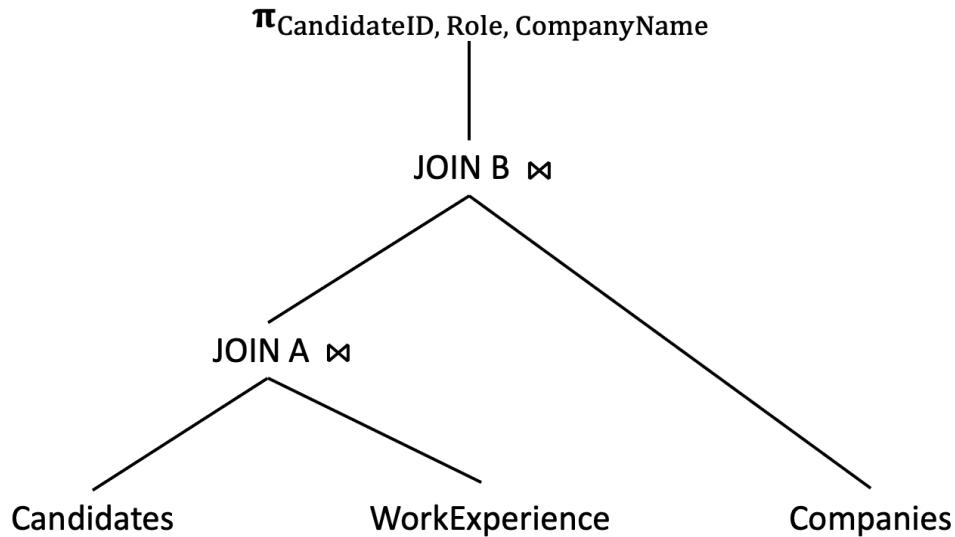
The statistics of these tables are shown below:

Table	Total Tuples	Total Occupied Blocks
Candidates	3000	40
WorkExperience	5000	70
Companies	200	30

Consider the following query:

```
SELECT CandidateID, Role, CompanyName
FROM Candidates, WorkExperience, Companies
WHERE Candidates.CandidateID=WorkExperience.CandidateID
    AND WorkExperience.CompanyID=Companies.CompanyID
```

The Joining strategy is as follows:



Assume the number of blocks in memory is **M = 20**. Answer the following questions:

- (i) Calculate the cost to perform block-based nested loop join in "Join A". You can choose either **Candidates** or **WorkExperience** as the outer relation for your computation. Please clearly state your choice in the solution.

- (ii) Assume that after "Join A", the joining result would be stored back to the disk first and the resulting table has 8000 tuples. Each block can contain up to 50 tuples in the resulting joined relation. Assume that the tuples are densely stored. Calculate the estimated cost to perform a block-based nested loop join in "Join B". You can choose either **Universities** or the '**JOIN A**' relation as the outer relation for **JOIN B**. Please clearly state the order in which you joined **Universities** and the '**JOIN A**' relation in your solution. [Hint: You need to compute the total number of blocks for "Join A" to be able to estimate "Join B"] (10 points)

Q5. Cost-based Optimization (15 points)

Note: On your exam, the auto-grader only rounds (not ceiling/floor) the final result and not the intermediary results. For this question, follow the same process, and do not round intermediate results.

Consider the following relations:

Q 4 :

1. Work Experience as outer loops.

$$M-2 = 20-2 = 18 \text{ blocks/each fine.}$$

For Work Experience : $B(W) = 70$

$$\text{We need to loop } \frac{B(N)}{(M-2)} = \frac{70}{18} \approx 4$$

∴ For Candidates : $4 \times B(C) = 4 \times 40 = 160$
Cost = $160 + 70 = 230 \text{ blocks.}$

$\alpha_4,$

2. * Regard University as Companies.

$$800 / 50 = 160 \text{ blocks}$$

Companies as outer loop.

$$B(C) = 30$$

$$\frac{B(C)}{M-2} \approx 2$$

$$\therefore 2 \times B(A) = 2 \times 160 = 320$$

$$\therefore \text{Cost} = 320 + 30 = 350 \text{ blocks}.$$

A(x,y,z)	B(w,x)	C(w,x,z)	D(y,z)
T(A) = 2000	T(B) = 4000	T(C) = 6000	T(D) = 5000
	V(B, w) = 100	V(C, w) = 30	
V(A, x) = 50	V(B, x) = 20	V(C, x) = 50	
V(A, y) = 20			V(D, y) = 50
V(A, z) = 40		V(C, z) = 20	V(D, z) = 20

Determine the most efficient way to join the 4 relations, i.e., $A \bowtie B \bowtie C \bowtie D$. Note that $T(R)$ is the number of tuples in relation R and $V(R, a)$ is the number of distinct values of attribute a in relation R.

Assume the distribution of values of each attribute in all relations are uniform and both Containment Of Values and Preservation Of Value Sets hold in the relations above. Show your work by completing the following table. Each step in the dynamic programming algorithm should be one row.

Subquery	Size	Cost	Plan
AB	$2000 \times 4000 / 50 = 160000 = 160k$	0	AB
AC	$2000 \times 6000 / 50 = 240000 = 240k$	0	AC
AD	$2000 \times 5000 / 50 = 200000 = 200k$	0	AD
BC	$4000 \times 6000 / 100 = 240000 = 240k$	0	BC
BD	$4000 \times 5000 / 100 = 200000 = 200k$	0	BD
CD	$6000 \times 5000 / 20 = 1500000 = 1.5M$	0	CD
ABC	48000	$= \text{size}(BC) = 4.8k$	$(BC)A$
ABD	400000	$= \text{size}(AB) = 400k$	$(AB)B$
ACD	30000	$= \text{size}(AD) = 30k$	$(AD)C$
BCD	120000	$= \text{size}(DC) = 12k$	$(DC)B$
$ABCD$	240000	$= \text{size}(BCA) + \text{cost}(BCA)$ $= 48k + 4.8k$ $= 52.8k$	$(BCA)D$

Q6. Physical Query Plan (15 points)

Consider the following schema for a bibliographic library (such as DBLP). Assume we have the following statistics for the tables and columns with the primary key(s) underlined:

Researcher (ORCID, name, affiliation, email)

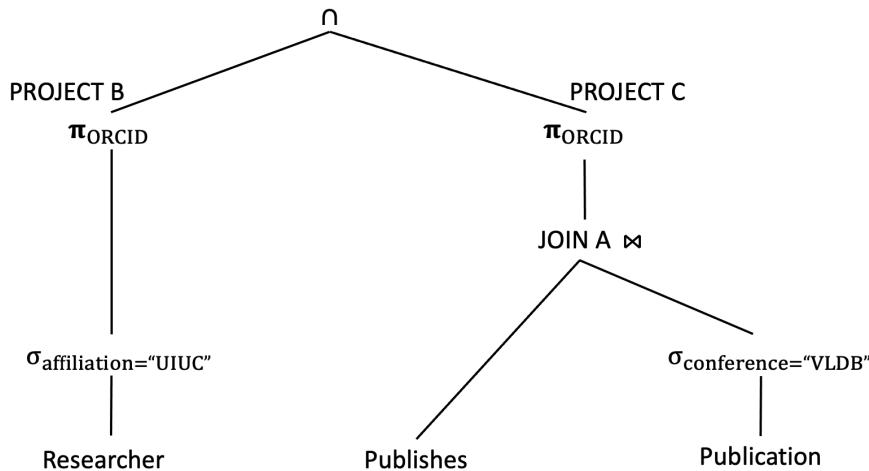
Publishes (ORCID, DOI, year)

Publication (DOI, title, conference)

Select the ORCIDs of researchers who are affiliated with 'UIUC' and have published at least one publication in conference 'VLDB'.

Table	# tuples	# blocks
Publishes	200000	2500
Publication	15000	1500
Researcher	5000	100

Column	# distinct
Researcher.Affiliation	10
Publication.Conference	30
Publishes.DOI	10000



Assume that the Researcher, Publishes, and Publication tables have a clustered index on their respective primary keys. Assume the distribution of Affiliation and Conference values are uniform, to compute estimated sizes after selection.

- (a) What is the cost and size of the selection conference = "VLDB"? Assume we have an unclustered index on "conference". (2 points)
- (b) What is the cost and size of executing selection affiliation = "UIUC"? (2 points)
- (c) What is the cost of executing Join A if we use an index-based nested loop join? Follow other assumptions from (a). Use the smaller relation as the inner relation. (3 points)
- (d) What is the cost of executing Join A if we use a two-pass hash-based join? Follow other assumptions from (a). (3 points)
- (e) Assuming the block size for the result table, after projection B/C (Use part d for Join A) is 100 tuples and memory capacity is 50 blocks, can we perform a one-pass algorithm for the intersection operation? If yes, estimate the cost of intersection. If not, explain why and estimate the total cost of the intersection with the appropriate algorithm. (Hint: First, estimate the size of result tables after projections B and C). (5 points)

Q6. Publishes : Psh Publications : Pca Researcher : R

a: $T(P_{ca}) = 15000$

$$V(P_{ca}, \text{Conference}) = 30$$

For unclustered attribute.

$$\therefore \text{cost} = \frac{T}{V} = \frac{15000}{30} = 500$$

$$\text{size} = T(\text{Select}) = \frac{T}{V} = 500 \text{ tuples.}$$

Q6 Publishers: P_{sh} Publications: P_{ca} Researcher: R

b. Cost: T(R) = 500, V(R, filtration) = 10, B(R) = 500

$$\therefore T(\text{Select}) = 8(R) = 100$$

$$\text{Size: } T(S) = \frac{T(R)}{\sqrt{V}} = 500 \text{ tuples}$$

Q6. Publishers: P_{sh} Publications: P_{ca} Researcher: R

C. according to a's result:

$$\text{For: } \sigma_{\text{conference}} = \cap_{P_{ca}} (P_{ca})$$

size: 500 tuples. = 50 blocks

D₂ is primary key.

$$V(P_{ca}, D_2) = T = 500$$

$$\text{For } P_{sh}: T(P_{sh}) = 200000$$

$$B(P_{sh}) = 2500$$

So P_{ca} will be inner relation. D₂, as primary key

∴ Since P_{sh} is clustered, index on S is clustered:

$$B(P_{sh}) + T(P_{sh}) B(P_{ca}) / V(P_{ca}, D_2)$$

$$= 2500 + 200000 \times 50 / 500$$

$$= 2500 + 20000$$

$$= 22500$$

Q6.

$$d: \beta(P_{sh}) = 25^\circ, \beta(P_{ca}) = 50^\circ$$

$$\therefore \text{Cost} \{ (25^\circ + 5^\circ) = 765^\circ \}$$

Q 6:

R : Before intersection:

$$T(P_{sh} \times (\sigma_{\text{conference} = "VLDB"} P_{ch})) = T(P_{sh}) \cdot T(R_1) / V(P_{sh}, D_{IT})$$

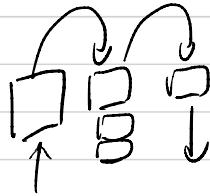
\downarrow bigger than R_1 's
 R_1
 $= 1000 \text{ tuples}$

For ($\sigma_{\text{affiliation} = "VLUC"} \text{ Reservation}$):

$$T(R_2) = 500 \text{ according to (b)}$$

$$\text{? block_size} = 100$$

$$\therefore 10000/100 = 100 \text{ blocks for (Form A)}$$
$$500/100 = 5 \text{ blocks for } R_2$$



Then it's similar to Union i

$$M \geq \min(B(R), B(S)) + 2$$
$$= 5 + 2 = 7$$

\therefore It meets the requirements
Cost = $B(R) + B(S) = 105$