

Question 1. (25 pts)

Consider the following relations for a database that keeps track of student enrollement in courses. Also consider that the relations are bags that allow duplicates:

- student<id, name, age, deptid> storing 100000 tuples
- department<id, name> storing 600 tuples
- course<id, name, semester> storing 3000 tuples
- enrollment<student_id, course_id, year, grade> storing 4000000 tuples

(a) For the following query:

List id and name of all students who didn't take courses in the first semester of 2021

give two RA expressions for this query (for example, by applying equivalence rules), one of which is optimized in terms of performance, and briefly discuss the estimated performance of the two expressions

$$\text{Select } s.id, s.name$$

$$\text{From student } s, \text{ enrollment } e$$

$$\text{Where } s.id = e.studentid \text{ AND NOT EXISTS}$$

$$\left(\begin{array}{l} \text{SELECT } * \\ \text{FROM course } c \\ \text{where } c.course_id = e.courseid \\ \text{AND } c.semester = 1 \text{ AND } e.year = 2021 \end{array} \right)$$

⇒ Not optimised as there is a big subquery

① $\pi_{s.id, s.name} (\sigma_{c.semester \neq 1 \wedge e.year \neq 2021} (p_c(course) \bowtie p_e(enrollment) \bowtie p_s(student)))$
 ↳ Not optimised as join THEN filtering, it is better to do the inverse for memory

② $\pi_{s.id, s.name} (p_s(student) \bowtie (\sigma_{e.year \neq 2021} (p_e(enrollment)) \bowtie \sigma_{c.semester \neq 1} (p_c(course))))$
 ↳ optimised as the filtering is made BEFORE the joins

① Filter 100.000 × 3000 × 4000.000 tuples

② Filter on 100.000 + 3000 + 4000.000 tuples ⇒ this is non negligible!!

(b) Write an RA expression to express this query:

List all department id and name of the students who attend the Database course ⇒ set-based query

Make sure not to show duplicate results

$$\text{SELECT } s.name, s.deptid \text{ DISTINCT FROM student } s, \text{ course } c, \text{ enrollment } e$$

$$\text{WHERE } c.name = 'Database' \text{ AND } c.id = e.course_id \text{ AND } e.studentid = s.id$$

$$\pi_{s.name, s.deptid} (p_c(\sigma_{name = 'Database'}(course)) \bowtie_{c.id = e.course_id} p_e(enrollment) \bowtie_{e.student_id = s.id} p_s(student))$$

$\pi_{s.name, s.dept_id} (\rho_c (\sigma_{name = 'Database'} (course)) \bowtie_{c.id = e.course_id} \rho_e (enrollment) \bowtie_{e.student_id = s.id} \rho_s (student))$

↳ There will be no duplicates as π is a projection and eradicates them! :P

(c) Normalize/flatten the following query and give its equivalent RA expression

```
SELECT course.name
FROM course
WHERE course.id in (
  SELECT enrollment.course_id
  FROM student, enrollment
  WHERE student.id = enrollment.student_id AND
  student.name LIKE '%son');
```

Normalizing → transforming to a EXIST subquery

```
SELECT course.name
FROM course
WHERE EXISTS ( SELECT 1 FROM enrollment, student
  WHERE enrollment.student_id = student.id AND student.name like '%son' )
```

Flattening is done in 6 steps

1. Translate the subquery

$\rho_s (\sigma_{name \text{ like } '%son'} (student)) \bowtie_{e.student_id = s.id} \rho_e (enrollment)$

2. Add context relation and parameters

$\pi_{course_id} (\rho_s (\sigma_{name \text{ like } '%son'} (student)) \bowtie_{e.student_id = s.id} \rho_e (enrollment))$

3. Translate the FROM of the outer query

$\rho_c (course)$

4. Syntactize both expression

$\rho_c (course) \bowtie \pi_{course_id} (\rho_s (\sigma_{name \text{ like } '%son'} (student)) \bowtie_{e.student_id = s.id} \rho_e (enrollment))$

5. simplify

6. Complete the expression

$\pi_{name} (\rho_c (course) \bowtie \pi_{course_id} (\rho_s (\sigma_{name \text{ like } '%son'} (student)) \bowtie_{e.student_id = s.id} \rho_e (enrollment)))$

(d) Suppose relation $R \langle A, B, C, D \rangle$ has the tuples

(1,2,3,4)

(1,2,3,5)

(3,2,1,0)

(A) Using **bag** projection and theta-join, how many tuples appear in the result of:

$\pi_{A,B}(R) \bowtie_{R.B < S.B} \rho_S(\pi_{B,C}(R))$

(B) Using **set** projection and theta-join, how many tuples appear in the result?

A	B	C	D
1	2	3	4

(A) $\pi_{A,B}(R)$

A	B
1	2

$\rho_S(\pi_{B,C}(R))$

B	C
2	3

A	B	C	D
1	2	3	4
1	2	3	5
3	2	1	0

(A) $\pi_{A,B}(R)$

A	B
1	2
1	2
3	2

$\rho_S(\pi_{B,C}(R))$

B	C
2	3
2	3
2	1

$\Rightarrow \pi_{A,B}(R) \bowtie_{R.B < S.B} \rho_S(\pi_{B,C}(R))$ outputs 0 results as $2 \nless 2$

(B) If it is set based

$\pi_{A,B}(R)$

$$R$$

A	B
1	2
3	2

$\rho_S(\pi_{B,C}(R))$

$$S$$

B	C
2	3
2	1

$\Rightarrow 0$ output