

COMP207 Assignment 2 – Query Processing

Issue Date: Monday, 18 November 2019

Submission Deadline: Tuesday, 03 December 2019, 17:00

About This Assignment

This is the second of two assignments for COMP207. It is worth 10% of the total marks for this module. It consists of four questions, which you can find at the end of this document.

Submit your solutions to these questions in PDF format by the given submission deadline. Your solutions must be submitted on Vital (see the detailed submission instructions below). Accuracy and relevance are more important in your answers, so don't write large volumes in your submission, but do ensure that what you write covers what is asked for and keeps to the problem statement.

Submission Details

Please submit **one PDF file with your solutions**. Name your file as follows:

<your student ID>-Assignment-2.pdf

If your student ID is 12345678, then your file should be named:

12345678-Assignment-2.pdf.

Please submit only this file (no archives).

To act as your ‘signature’ for the assignment, at the top of your PDF document put your Student ID number.

ID: 201447569

Your solutions must be submitted on Vital (see Vital for submission instructions).

The submission deadline for this assignment is **Tuesday, 03 December 2019, 17:00**. Earlier submission is possible, but any submission after the deadline attracts the standard lateness penalties. Plagiarism and collusion guidelines will apply throughout the assignment submission. For details on late submissions, how to claim extenuating circumstances, etc., please see the undergraduate student handbook, which can be found at <http://intranet.csc.liv.ac.uk/student/ug-handbook.pdf>, or in Section 6 of the Code of Practice on Assessment.¹

Assessment information at a glance

Assignment Number	2 (of 2)
Weighting	10% of the final module mark
Assignment Circulated	Monday, 18 November 2019
Deadline	Tuesday, 03 December 2019, 17:00
Submission Mode	Electronically on Vital
Learning Outcome Assessed	LO2: Demonstrate an understanding of advanced SQL topics
Purpose of Assessment	Assessment of knowledge of SQL query processing
Marking Criteria	See description of this assignment
Submission necessary in order to satisfy module requirements?	N/A
Late Submission Penalty	Standard UoL Policy

¹ https://www.liverpool.ac.uk/media/livacuk/tqsd/code-of-practice-on-assessment/code_of_practice_on_assessment.pdf

Question 1 (10 marks)

The following tables form part of a hotel booking database held in a relational DBMS (primary keys are underlined):

Hotel (hotelNo, hotelName, city)
Room (roomNo, hotelNo, type, price)
Booking (hotelNo, guestNo, dateFrom, dateTo, roomNo)
Guest (guestNo, guestName, guestAddress)

- **Hotel** contains hotel details and hotelNo is the primary key.
- **Room** contains room details for each hotel and (roomNo, hotelNo) forms the primary key.
- **Booking** contains details of the bookings and (hotelNo, guestNo, dateFrom) forms the primary key.
- **Guest** contains guest details and guestNo is the primary key.

Give the relational algebra expressions to return the results for the following two queries:

- (a) List the names and cities of those hotels who charge more than £85 for a room. **(5 marks)**
- (b) List the names and addresses of guests who have made a booking to stay Christmas Day 2019. **(5 marks)**

(a)

$$\pi_{\text{hotelName}, \text{city}}(\sigma_{\text{price} > 85}(\text{Hotel} \bowtie \text{Room}))$$

(b)

$$\pi_{\text{guestName}, \text{guestAddress}}(\sigma_{\text{dateFrom} \leq \text{date('2019-12-25')} \wedge \text{dateTo} \geq \text{date('2019-12-25')}}(\text{Guest} \bowtie \text{Booking}))$$

Question 2 (10 marks)

Consider the following database schema and example instance for a property management system:

property

pId	price	owner	sqrFeet	location
1	100,000	Alice	560	Lake View
2	3,400,000	Bob	2,000	Hyde Park
3	1,200,000	Bob	1,200	Hyde Park
4	5,000,000	Martha	800	Evanston

repairs

rId	pId	company	date	type
1001	1	M.M. Plumbing Ltd.	2013-12-12	Bathroom
1002	2	M.M. Plumbing Ltd.	2013-12-13	Kitchen
1003	4	Rob's Double Glazing	2012-01-01	Windows

Hints:

- Attributes with a grey background form the primary key of a relation (e.g, *pId* for relation *property*).
- The attribute *pId* of relation *repairs* is a foreign key to relation *property*.

Give the relational algebra expressions to return the results for the following two queries:

- (a) Get the pId, owner and location details of all properties that are larger than 900 square feet (sqrFeet). **(5 marks)**
- (b) Get the names of repair companies (company) that did a repair on a property in Hyde Park. **(5 marks)**

(a) $\pi_{pId, owner, location} (\sigma_{sqrFeet > 900} (Property))$

(b).

$\pi_{company} (\sigma_{location = 'HydePark'} (Property \bowtie repairs))$

Question 3 (20 marks)

- (a)** Consider the following relation:

studentCourses(StudentID, CourseNo, Quarter, Year, Units, Grade)

The relation contains the grades for the courses completed by students. Assume that in studentCourses there are 200,000 different students, each identified by their StudentID. On average, a student took 40 different courses.

If the file blocks hold 2000 bytes and each studentCourses tuple requires 50 bytes, how many blocks will then be needed to store the relation studentCourses? **(5 marks)**

- (b)** A database includes two relations Student (S) and Program (P).

S			
Student_No	F_Name	L_Name	Prog_Code
04009991	Alicia	Smith	0001
04009992	Alan	Smith	0002
04009995	Alicia	Bush	0001
04009996	John	Smith	0001

P	
Prog_Code	P_Name
0001	Computing
0002	Software Engineering

Give a relational expression that could possibly return the following result:

F_Name	L_Name	P_Name
Alicia	Smith	Computing
John	Smith	Computing

(5 marks)

- (c)** Translate the following relational algebra into SQL:

$\pi_{\text{studId}, \text{lName}}(\sigma_{\text{course}='BSc'}(\text{STUDENT}))$ **(5 marks)**

- (d)** Given these relations, write the SQL statement that produced the equivalent queries below:

Course (courseNo, courseDept, courseLeader)

Student (studNo, name, type, tutorId, courseNo)

Two sample equivalent corresponding queries have been produced:

$\pi_{\text{studno}, \text{name}}(\sigma_{\text{type}='undergrad'} \wedge (\text{courseDept}='CompSci'))(\text{Student} \bowtie \text{s}. \text{courseNo} = \text{c}. \text{courseNo})$ Course
and
 $\pi_{\text{studno}, \text{name}}(\sigma_{\text{type}='undergrad'}(\text{Student})) \bowtie \text{s}. \text{courseNo} = \text{c}. \text{courseNo} (\sigma_{\text{courseDept}='Comp Sci'}(\text{Course}))$

(5 marks)

Q3

(a)

$$200000 \cdot 40 = 8 \cdot 10^6 \text{ tuples}$$

$$\text{size} = 8 \cdot 10^6 \cdot 50 \text{ bytes} = 4 \cdot 10^8 \text{ bytes}$$

$$\text{block} = \frac{4 \cdot 10^8 \text{ bytes}}{2 \cdot 10^3 \text{ bytes}} = 2 \cdot 10^5 \text{ blocks.}$$

(b)

TL F-Name, L-Name, P-Name (6s. prog_code = 0001 AND S.L-Name = 'Smith' (SXP))

(c) SELECT StudId, LName

FROM Student

WHERE course = 'BSc'

(d)

SELECT Student.studNo, Student.name

FROM Student, Course

WHERE Student.courseNo = Course.courseNo AND

Student.type = 'Undergrad' AND

Course.courseDept = 'Comp Sci'

Question 4 (60 marks)

Consider a database with relations $R(A, B, C)$, $S(D, E)$, and $T(F, G)$.

- (a) Give the initial query plan (constructed as in Lecture 13) for the SQL query

```
SELECT B, E, G
FROM R, S, T
WHERE A = 10 AND C = D AND E = F AND A > G;
```

Then use the heuristics from Lecture 16 to transform the initial query plan into an optimised (logical) query plan. Perform the transformation step-wise, pushing a single operator over a single operator in each step, and indicate the heuristics you apply. **(20 marks)**

- (b) Suppose that

- $|R| = 1000$, $|\pi_A(R)| = 1000$, $|\pi_B(R)| = 100$, $|\pi_C(R)| = 500$;
- $|S| = 5000$, $|\pi_D(S)| = 300$, $|\pi_E(S)| = 10$;
- $|T| = 4000$, $|\pi_F(T)| = 4000$, $|\pi_G(T)| = 1500$.

Estimate the number of tuples returned by the following queries. Explain your calculations.

- i) $\sigma_{A=10}(R)$ **(6 marks)**
- ii) $\sigma_{A=10 \text{ OR } B=b}(R)$ **(6 marks)**
- iii) $R \bowtie_{C=D} S$ **(6 marks)**

- (c) Suppose that in addition to the assumptions on R , S , and T from part (ii), we also have the following:

- Each disk block can hold up to 10 tuples.
- All relations are stored in consecutive blocks on disk.
- No indexes are available.

What is the best physical query plan (in terms of the number of disk access operations) you can find for $\sigma_{B=b \text{ AND } E=100}(R \bowtie_{C=D} S)$? Describe your plan and show the calculation of the number of disk access operations. **(22 marks)**

Q4
(a)

$$\pi_{B,E,G}$$

①

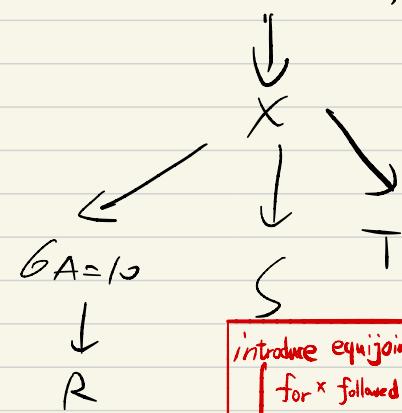
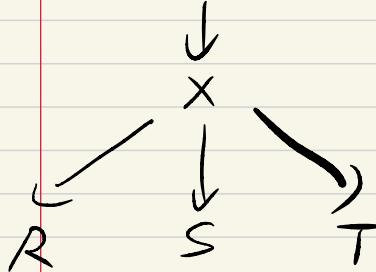
push selection as far down the tree as possible

②

$$\pi_{B,E,G}$$

$$6 A=10 \text{ AND } c=D \text{ AND } E=F \text{ AND } A>G$$

$$6 c=D \text{ AND } E=F \text{ AND } A>G$$



④

$$\pi_{B,E,G}$$

$$6 A>G$$

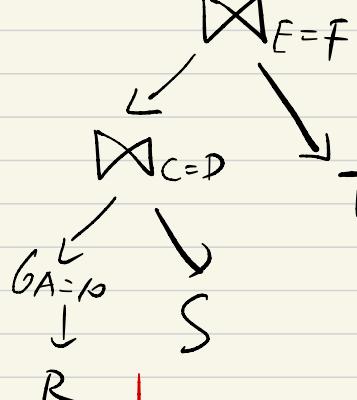
S

introduce equijoin for X followed by 6

③

$$\pi_{B,E,G}$$

$$6 A>G \text{ AND } E=F$$



←

|

$$6 A=10$$

R

$$6 A=10$$

R

$$6 A=10$$

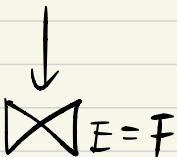
R

push selections as far down the tree as possible

(3)

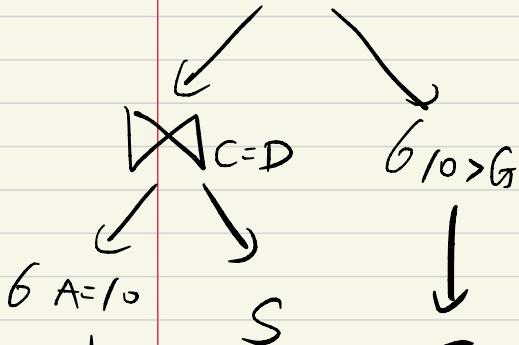
(6)

$\pi_{B,E,G}$

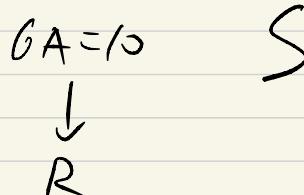
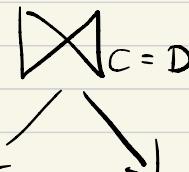


insert projections where appropriate

$\pi_{B,E,G}$



$\pi_{B,E}$



insert projections
where appropriate

7

$\pi_{B,E,G}$



$\bigcap_{E=F}$



$\pi_{B,E}$



$G_A > G$

$\bigcap_{C=D}$



$\pi_{B,C}$

S



$G_A = 10$



R

(b) Assume that each value is used equally often.

i) $\frac{|R|}{\text{No. distinct A or R}} = \frac{1000}{1000} = 1$

ii)

$$\frac{|R|}{\text{No. distinct A or R}} + \frac{|R|}{\text{No. distinct B or R}} = 1 + \frac{1000}{100} = 11$$

minus the interaction part

So the value can be $11 - 1 = 10$ or 11
(with interaction) (without interaction)

iii)

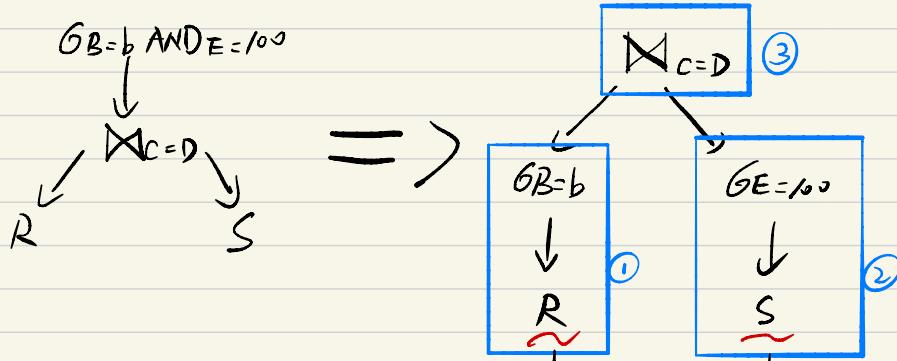
$$\frac{|R| \times |S|}{\text{max. number of distinct values for A in R or S}} = \frac{1000 \cdot 5000}{500(\text{Max})} = 10000$$

Q4

(c) Assume that the M which is number of disk block fit into RAM.

Consider that $GB = b$ AND $GE = 100$ ($R \bowtie_{C=D} S$)

Using the tree representation:



Using External merge sort:

Consider ①: Total disk operations:

$$O\left(\frac{N_1}{B} \log_m\left(\frac{N_1}{MB}\right)\right) = O\left(\frac{1000}{10} \log_m\left(\frac{1000}{10M}\right)\right) = O\left(100 \log_m\left(\frac{100}{M}\right)\right)$$

($B = 10$ because Each disk block can hold up to 10 tuples)

Consider ②: Total disk operations:

$$O\left(500 \log_m\left(\frac{500}{M}\right)\right)$$

②

Consider ②: No. of elementary operations :

$$O(|R| \log_2 |R| + |S| \log_2 |S|)$$

So No. of disk operations :

$$O\left(\frac{|R|}{B} \log_m \frac{|R|}{B} + \frac{|S|}{B} \log_m \frac{|S|}{B}\right)$$

$$= O\left(\frac{10}{10} \log_m 1 + \frac{500}{10} \log_m 50\right)$$

$$= O(1 \log_m 1 + 50 \log_m 50)$$

Consider ①, ② and ③.

Total disk operation is :

$$\underline{\underline{O(1 \log_m 1 + 50 \log_m 50 + 500 \log_m (500) + 100 \log_m (100))}}$$

③ ② ①