

National University of Computer and Emerging Sciences, Lahore Campus

Course:	Advanced Database Concepts	Course Code:	CS451
Program:	BS(Computer Science)	Semester:	Spring 2017
Duration:	3 hours	Total Marks:	50
Paper Date:	Tue 23-May-2017	Weight	40%
Section:	CS	Page(s):	
Exam:	Final		

Instruction/Notes: Scratch sheet can be used for rough work however, all the questions and steps are to be shown on this question paper. No extra/rough sheets should be submitted with question paper. You will not get any credit if you do not show proper working, reasoning and steps as asked in question statements. Calculators are allowed.

Q1. (10 points) Consider a disk with block size $B=512$ bytes. A block pointer is $P=6$ bytes long, and a record pointer is $P_R=7$ bytes long. A file has $r=100,000$ EMPLOYEE records of fixed-length. Record length R is 115 bytes long and DEPTCODE field is 15 bytes long.

Suppose the file is ordered by the non-key field DEPTCODE and we want to construct a clustering index on DEPTCODE that uses block anchors (every new value of DEPTCODE starts at the beginning of a new block). Assume there are 500 distinct values of DEPTCODE, and that the EMPLOYEE records are evenly distributed among these values. Calculate:

- The index blocking factor (bfr_i).
- The number of first-level index entries (r_1) and the number of first-level index blocks (b_1).
- The number of levels needed (x) if we make it a multi-level index.
- The total number of blocks required by the multi-level index (b_i).
- The number of block accesses needed to search for and retrieve all records in the file having a specific DEPTCODE value using the clustering index (assume that multiple blocks in a cluster are either contiguous or linked by pointers).

Ans:

a) the index blocking factor (bfr_i).

Index record size $R_i = (V \text{ DeptCode} + P) = (15 + 6) = 21$ bytes

$bfr_i = fo = \text{floor}(B/R_i) = \text{floor}(512/21) = 24$

b) the number of first-level index entries (r_1) and the number of first-level index blocks (b_1).

$r_1 = \text{number of distinct Department_code values} = 500$ entries

$b_1 = \text{ceiling}(r_1 / bfr_i) = \text{ceiling}(500/24) = 21$ blocks

c) the number of levels needed (x) if we make it a multi-level index.

We can calculate the number of levels as follows:

$r_2 = \text{number of 1st-level index blocks } b_1 = 21$ entries

$b_2 = \text{ceiling}(r_2 / bfr_i) = \text{ceiling}(21/24) = 1$ block;

Hence, the index has $x = 2$ levels

d) the total number of blocks required by the multi-level index (b_i).

$b_i = b_1 + b_2 = 21 + 1 = 22$ blocks

e) the number of block accesses needed to search for and retrieve all records in the file having a specific DEPTCODE value using the clustering index (assume that multiple blocks in a cluster are either contiguous or linked by pointers).

Number of block accesses to search for the first block in the cluster of blocks $= x + 1 = 2 + 1 = 3$

The 200 records are clustered in $\text{ceiling}(200/bfr) = \text{ceiling}(200/24) = 9$ blocks.

RollNo: _____

Name: _____

Hence, total block accesses needed on average to retrieve all the records with a given DeptCode=
 $x+50=2+50=52$ block accesses

RollNo: _____

Name: _____

Q2. (4 points) Assume a relation R (A, B, C) is given; R is stored as an ordered file (un-spanned) on non-key field C and contains 500,000 records. Attributes A, B and C need 5 byte of storage each, and blocks have a size of 2048 Bytes. Each A value occurs at an average 5 times in the database, each B value occurs 50 times in the database, and each C value occurs 50,000 times in the database. Assume there is no index structure exists.

Estimate the number of block fetches needed to compute the following queries (where C_a and C_c are integer constants):

a) SELECT B, C FROM R WHERE A = C_a ;

b) SELECT B, C FROM R WHERE C = C_c ;

Ans: bfr= $2048/15=136$; b= $500,000/136= 3677$

a) $O(b) = 3677$

b) $O(\log(b) + s)$ i.e. $O(12 + 368 - 1) = 379$

RollNo: _____

Name: _____

Q3. (3 points) Consider the student table:

<u>RollNo</u>	Name	Address	Gender	Age	Grade
1001	Khadija	Faisal	F	16	B
1002	Tahree	Town	F	16	C
1003	m	Model	F	18	A
1004	Isbah	Town	M	18	B
1005	Izaan	DHA	F	20	A
1006	Alia	Model	F	17	B
1007	Tahree	Town	M	19	A
1008	m	Faisal	M	17	D
	Ismail	Town			
	Izaan	DHA			
		Johar Town			
		DHA			

Find the selectivity (*s*) of the condition to retrieve:

a) RollNo=1004

b) Gender='F'

c) Age=16

Ans:

a) RollNo=1004 : $1/8 = 0.125$ (12.5 %)

b) Gender='F' : $5/8 = 0.625$ (62.5 %)

c) Age=16 : $2/8 = 0.25$ (25 %)

RollNo: _____

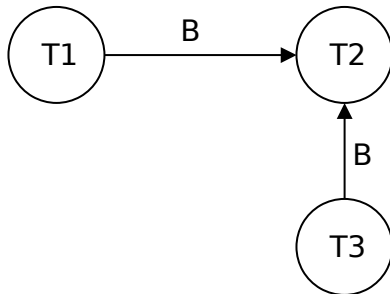
Name: _____

Q4. (4+6= 10 points) Consider the schedule:

Op #	T1	T2	T3
1	Read(A)		
2		Read(C)	
3			Read(A)
4		Write(C)	
5			Read(B)
6	Read(B)		
7		Write(B)	

a) Draw the precedence graph of the schedule given above. In case it is conflict-serializable then list down all **the equivalent serial schedules**.

Ans: a) it is conflict-serializable; Equivalent serial Schedules are $T1 \rightarrow T3 \rightarrow T2$ & $T3 \rightarrow T1 \rightarrow T2$.



RollNo: _____

Name: _____

b) Show that the schedule will be accepted/rejected by the below protocols. Provide proper reason and show your working.

i) The basic two-phase locking protocol (add locks to the transactions)

ii) The timestamp-ordering protocol (you have $T1 < T2 < T3$)

Ans: b) i) Accepted:

Op #	T1	T2	T3
1	RI-A Read(A)		
2		RI-C Read(C)	
3			RI-A Read(A)
4		WI-C; (<u>upgrade lock</u>) Write(C)	
5			RI-B Read(B) C3; ul-A, ul-B
6	RI-B Read(B) C1; ul-A, ul-B		
7		WI-B Write(B) C2; ul-C, ul-B	

ii) Rejected

Op #	T1	T2	T3
1	Read(A)		
2		Read(C)	
3			Read(A)
4		Write(C)	
5			Read(B)
6	Read(B)		
7		Write(B); abort T2 R-TS(B) > TS(T2)	

RollNo: _____

Name: _____

Q5. (3 points) Suppose that the most often used query on the Student database is:

```
SELECT StudentName, CourseCode, LetterGrade
FROM student S JOIN grade G ON S.RollNo=G.RollNo WHERE S.BatchId='2014';
```

On which column(s) would you create an index? Write down the column name(s) and one sentence why you choose the column(s).

Ans: S.BatchId (filter column) , S.RollNo (joining col), G.Rollno (joining col)

Q6. (3+4+3= 10 points) Figure below shows the log corresponding to a particular schedule at the point of a system crash for five transactions. Suppose that we use the immediate update (undo/redo) protocol with *checkpointing*. Describe the recovery process from the system crash.

Assume that the initial values of items are $A=100$, $B=200$, and $C=300$. Isolation level of all transactions is *READ COMMITTED*.

- a) Identify which transactions need undo/ redo operation(s)?
- b) Specify which operations in the log are redone (in correct order) and which are undone.
- c) Write down the values of items A, B, and C after system recovery.

```
[start_transaction, T1]
[read_item, T1, B, 200]
[start_transaction, T2]
[read_item, T2, A, 100]
[write_item, T2, A, 100, 50]
[read_item, T2, B, 200]
[write_item, T2, B, 200, 120]
[commit, T2]
[start_transaction, T3]
[read_item, T1, A, 50]
[write_item, T1, A, 50, 20]
[read_item, T3, C, 300]
[write_item, T3, C, 300, 250]
[checkpoint]
[commit, T3]
[start_transaction, T4]
[read_item, T4, C, 250]
[start_transaction, T5]
[read_item, T5, C, 250]
[write_item, T5, C, 250, 210]
[commit, T5]
```

System crash

RollNo: _____

Name: _____

Ans:

a)- T2 was committed before the checkpoint and hence is not involved in the recovery.

- The list of committed transactions T since the last checkpoint contains transactions T3 and T5. Hence T3 and T5 need redo operations.

- The list of active transactions T' contains transactions T1 and T4. Hence they are cancelled and must be resubmitted. Hence T1 and T4 need undo operations.

b)- Only the WRITE operations of the committed transactions (i.e. T3 and T5) are to be redone. Hence, REDO is applied to:

[write_item, T3, C, 250]

[write_item, T5, C, 210]

Only the WRITE operations of the cancelled transactions (i.e. T1 and T4) are to be undone. Hence, UNDO is applied to:

[write_item, T1, A, 50]

- The transactions that are active and did not commit i.e., transactions T1 and T4 are canceled and must be resubmitted. Their operations have to be undone.

c) The values of items are A=50, B=120, C=210

Q7. (10 points) Consider the bank database, and the following SQL query:

Customer (custID, custName, cnic, birthDate, address, ...)

Account (accNo, custID, accTitle, accType, openingDate, ...)

Transaction (tID, accNo, transType, amount, transDate, ...)

```
SELECT C.cnid, A.accNo, A.Title, T.noOfTrans
FROM customer C JOIN account A ON C.custID=A.custID JOIN (SELECT accNo, COUNT(*) AS noOfTrans
FROM transaction GROUP BY accNo) T ON A.accNo=T.accNo
```

Write an efficient relational-algebra expression that is equivalent to this query and draw the optimal query plan for this query.

Ans: