

Name	NetID

CS411: Database Systems

Spring 2017

Midterm 2, May 1

- READ THESE INSTRUCTIONS CAREFULLY BEFORE YOU START. DO NOT turn this page UNTIL the proctor instructs you to.
- Check that your exam booklet consists of 13 printed sides, including this one.
- Write your NetID at the top of all the sheets.
- The exam lasts for 90 minutes, i.e., from 8–9.30am.
- We will not answer any questions during the exam. If you need to make any assumptions for any of the questions, please feel free to do so and clarify the assumption in your answer.
- All questions are compulsory.
- If you need more space, feel free to use the back side of each page. We will not provide extra sheets, so please use your space wisely. Show all necessary steps as part of your calculation to get partial credit.
- The maximum score you can obtain is $25 + 15 + 15 + 30 + 15 = 100$.
- You must stop writing when time is called by the proctors.
- **Cheating: No.**

Question	1	2	3	4	5	Total
Points						

True or False - 25 points

Please **circle** the right answer, and provide a short 1–2 line description justifying your choice.

1. [1] You know the answer to this question.

True False

2. [3] Consider the following three join mechanisms taught in class: 1) block-based nested loop join, 2) one-pass hash join, and 3) two-pass hash join. The memory requirements of these mechanisms are in the following order: $\text{requirement}(1) \leq \text{requirement}(3) \leq \text{requirement}(2)$.

True False

3. [3] For undo logging, if `<ABORT T>` is seen in the log on disk, then transaction T has no dirty data on disk, i.e., all of the changes made by transaction T have been written to disk.

True False

4. [3] Consider a relation $R(A, B, C)$, with two functional dependencies $A \rightarrow BC$, $B \rightarrow C$. Suppose we decompose R into $R_1(A, B)$, and $R_2(B, C)$, we can no longer determine if the functional dependency $A \rightarrow C$ holds in R , using R_1 and R_2 , and therefore the dependency $A \rightarrow C$ is not preserved.

True False

5. [3] Consider the following algebraic law which we know is valid for set semantics:
 $\sigma_C(R \cup S) = \sigma_C(R) \cup \sigma_C(S)$. This law remains valid for bag semantics as well.

True False

6. [3] We always prefer bushy join trees over left/right-deep join trees. This is because in left/right-deep join trees, we incur more cost because we have more intermediate join operations to compute.

True False

7. [2] While estimating the size of a join operation between two relations R and S , with the join attribute as A , we implicitly assume that all of the values of A in one of the two relations, is present in the other.

True False

8. [2] Unclustered indices are always dense, in the sense that every unique search key has an entry in the index.

True False

9. [3] B-trees provide better response time for range queries than hashing-based indexing schemes.

True False

10. [2] For a given relation that is not in BCNF, there is always a decomposition of the given relation into two relations which are in BCNF—this decomposition may not be unique.

True False

SQL - 15 points

You are the owner of a pet hotel. Since you are so smart and hard working, the size of your hotel has grown so large that it is impossible to keep track of pets' information with handwritten logs anymore and therefore you decide to use a relational database to help you manage your hotel data.

Joey, the almighty SQL madman, has designed the following relational schema for your database:

Pet (*pid*, *name*, *type*, *owner*)

Contract (*cid*, *pid*, *startDate*, *duration*)

PricePerDay (*type*, *price*)

You can make the following assumptions:

- The attributes *startDate* and *duration* are both expressed in days, an integer. For example, *startDate* = 400 means the 400th day, while *duration* = 10 means the duration is 10 days.
- The attribute *type* in the Pet table is a foreign key to the attribute *type* in the PricePerDay table.
- The attribute *pid* in the Contract table is a foreign key to the attribute *pid* in the Pet table.

Specify the following queries using SQL:

1. [4] How much time does each type of pet spend in the hotel? Returned tuples should be formatted as (*type*, *totalDuration*)

3. [6] Who is the most frequent pet owner that comes to your hotel, i.e., who is the owner who provides the most contracts?

B-Trees and Indexing - 15 points

1. Consider a partially filled B-tree of order 2 as shown in Figure 1.
 - (a) [4] In Figure 1, fill in missing values in the root node of the tree. Use the figure itself to save time on re-drawing the B-tree.

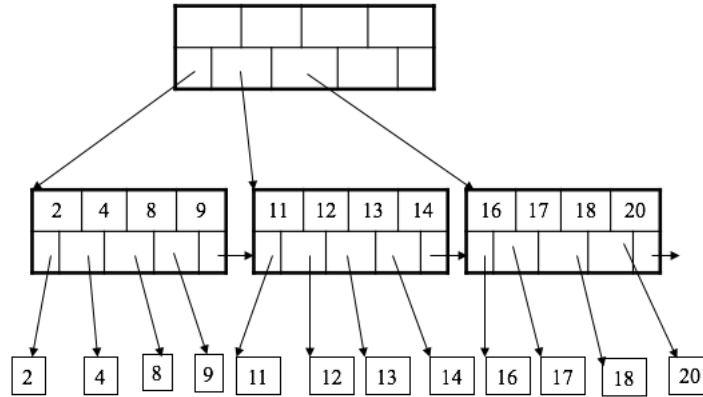


Figure 1: B-tree with root values missing

- (b) [6] There are many sequences of exactly 3 insertions that will cause the root node in the B-tree shown in Figure 1 to split. List one such sequence of three insertions, and justify your choice (1–2 lines are sufficient)

2. [5] Draw an extensible hash table with a block capacity of 2, i.e., each block holds two records, after the insertion of the following search keys (in the same order as provided):
1000, 1101, 0111, 1011, 1111

You only need to draw the final table. Indicate the bit counter for each bucket as well as the array of pointers.

Query Execution and Optimization - 30 points

1. [10] Consider four relations with the following statistics:

A(w,x)	B(x,y)	C(y,z)	D(w,z)
T(A) = 500	T(B) = 200	T(C) = 400	T(D) = 100

- R(a,b) means relation R has attributes a and b. T(R) refers to the number of tuples in relation R.
- Use formula $T(R1 \bowtie R2) = T(R1) \times T(R2) \times 0.01$ for size estimation, where R1 and R2 are any two relations.
- Join the relations only where a natural join is feasible.
- You may assume that the join operation is symmetric, i.e., the plan (R1R2) is the same as the plan (R2R1).

Use the Dynamic Programming algorithm to find the optimal plan to join the above relations. Complete the following table:

Subset	Size	Lowest cost	Lowest cost plan
AB	1000	0	AB
BC	800	0	BC
CD	400	0	CD
AD	500	0	AD
ABC	4000	800	A(BC)
BCD			
ABD			
ACD			
ABCD			

2. Use Figure 2 to answer the following questions. Use $B(R)$, $B(S)$, $B(T)$, $B(U)$ to denote the number of blocks in each of the relations. For convenience, use $K1 = B(R \bowtie S)$ and $K2 = B(R \bowtie S \bowtie T)$. Let M denote the number of blocks in the buffer. Consider the following physical execution plans; **for each plan, estimate both the cost and memory requirements for that plan**. You can skip the cost of the final write of $(R \bowtie S \bowtie T \bowtie U)$. You should include cost of reading and writing the initial relations and any intermediates. Use the formulae discussed in class; if you wish to use other formulae, then clarify the assumptions you are making.

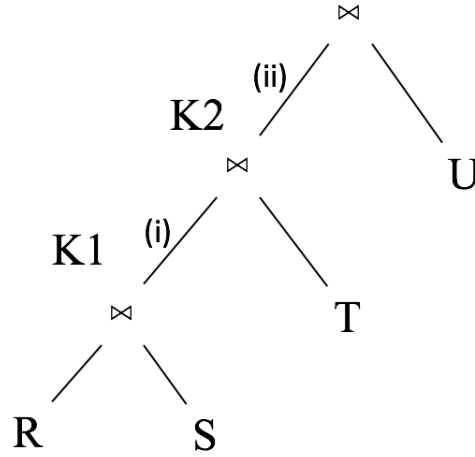


Figure 2: Query plan for Question 2

- (a) [6] Construct hash tables of S , T , U , all of which fit in the memory simultaneously. Use a **one-pass hash join** for $R \bowtie S$, $(R \bowtie S) \bowtie T$ and $(R \bowtie S \bowtie T) \bowtie U$, while pipelining the output of $R \bowtie S$ and $(R \bowtie S) \bowtie T$ without materialization i.e., you don't materialize at (i) and (ii) (see Figure 2).

- (b) [6] Construct hash tables for S and T , both of which fit in memory simultaneously. Use a **one-pass hash join** for $R \bowtie S$ and $(R \bowtie S) \bowtie T$ while pipelining the output of $R \bowtie S$ at (i). Then, materialize $(R \bowtie S) \bowtie T$ at (ii), construct a hash table for U , and use a one-pass hash join for $(R \bowtie S) \bowtie T$ with U .
- (c) [8] Use a block nested loop join for $R \bowtie S$ (with R as the outer relation), pipeline $R \bowtie S$ at (i), while maintaining a hash table for T , join $R \bowtie S$ with T using a **one-pass hash join**. Then, materialize $(R \bowtie S) \bowtie T$ at (ii), construct a hash table for U , and use a one-pass hash join for joining $(R \bowtie S) \bowtie T$ with U .

Transaction Management - 15 points

Consider the following **UNDO** log, and use it to answer the following questions.

<u>LogID</u>	<u>Log</u>
1	<START T1>
2	<T1, A, 20>
3	<START T2>
4	<START T3>
5	<T2, C, 7>
6	<T1, B, 15>
7	<T2, D, 9>
8	<COMMIT T1>
9	<START T4>
10	<T3, E, 2>
11	<T4, A, 6>
12	<ABORT T2>
13	<T3, B, 9>
14	<COMMIT T3>
15	<T4, B, 21>
16	<START T5>
17	<START T6>
18	<T6, D, 8>
19	<COMMIT T6>
20	<T5, E, 6>

1. [4] At the end of the log, what is the value of Variable C? Explain why (1–2 lines are sufficient).

3. [6] Continue from (2) after starting checkpointing. Suppose the system crashes right after LogID 13. Show which actions (e.g.: $\langle T1, A, 15 \rangle$) need to be undone and in what order.