

Name: _____

CSE 544, Winter 2009, Final Examination
11 March 2009

Rules:

- Open books and open notes.
- No laptops or other mobile devices.
- Calculators allowed.
- Please write clearly.
- Relax! You are here to learn.

Question	Max	Grade
1	15	
2	15	
3	10	
4	14	
5	14	
6	12	
7	8	
8	12	
Total	100	

Name: _____

1. (15 points) **Relational Model**

(a) (3 points) Briefly explain what is a “record-at-a-time” data manipulation language.

(b) (3 points) Briefly explain what is a “set-at-a-time” data manipulation language.

(c) (4 points) What are the benefits of a “set-at-a-time” over a “record-at-a-time” data manipulation language?

(d) (3 points) What is an integrity constraint? Why are integrity constraints important?

(e) (2 points) Give examples of two types of integrity constraints.

Name: _____

2. (15 points) **Storage Manager**

(a) (4 points) Why do DBMSs implement their own buffer managers?

(b) (3 points) To store its data, a DBMS can either use OS files or it can talk directly to a disk device. Give one benefit of each approach.

(c) (4 points) What is GiST and when is it helpful?

(d) (4 points) Give one limitation of the GiST approach.

Name: _____

3. (10 points) **Operator Algorithms**

Relation R has 90 pages. Relation S has 80 pages. Explain how a DBMS could efficiently join these two relations given that only 11 pages can fit in main memory at a time. Your explanation should be **detailed**: specify how many pages are allocated in memory and what they are used for; specify what exactly is written to disk and when.

(a) (5 points) Present a solution that uses a hash-based algorithm.

(b) (5 points) Present a solution that uses a sort-based algorithm.

Name: _____

4. (14 points) **Query Optimization**

Consider the following SQL query that finds all applicants who want to major in CSE, live in Seattle, and go to a school ranked better than 10 (i.e., rank < 10).

Relation	Cardinality	Number of pages	Primary key
Applicants (<u>id</u> , name, city, sid)	2,000	100	id
Schools (<u>sid</u> , sname, srank)	100	10	sid
Major (<u>id</u> , major)	3,000	200	(id,major)

```

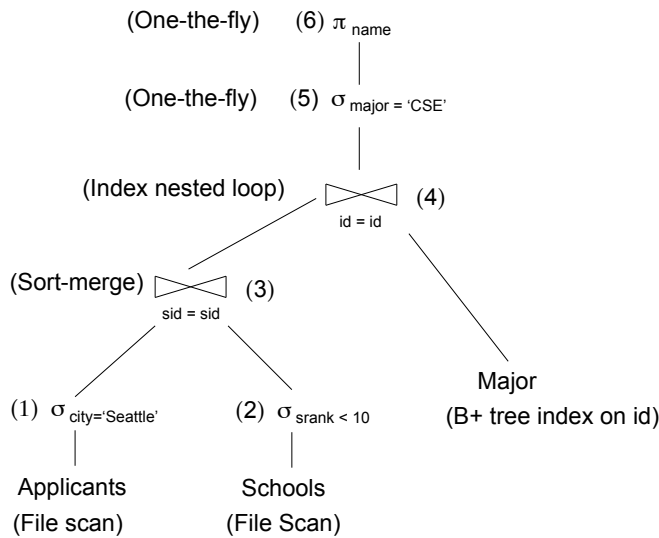
SELECT A.name
FROM Applicants A, Schools S, Major M
WHERE A.sid = S.sid AND A.id = M.id
AND A.city = 'Seattle' AND S.rank < 10 AND M.major = 'CSE'

```

And assuming:

- Each school has a *unique* rank number (srank value) between 1 and 100.
- There are 20 different cities.
- Applicants.sid is a foreign key that references Schools.sid.
- Major.id is a foreign key that references Applicants.id.
- There is an unclustered, secondary B+ tree index on Major.id and all index pages are in memory.

(a) (10 points) What is the cost of the query plan below? Count only the number of page I/Os.



Name: _____

- (b) (4 points) The Selinger optimizer uses a dynamic programming algorithm coupled with a set of heuristics to enumerate query plans and limit its search space. Draw two query plans for the above query that the Selinger optimizer (as described in the paper) would NOT consider. For each query plan, indicate why it would not be considered.

Name: _____

5. (14 points) **Transactions**

For each statement below, indicate if it is true or false.

- (a) (2 points) Serializability is the property that a (possibly interleaved) execution of a group of transactions has the same effect on the database and produces the same output as some serial execution of those transactions.

TRUE FALSE

- (b) (2 points) The following schedule is serializable:

$r_0[A] \rightarrow w_0[A] \rightarrow r_1[B] \rightarrow w_1[B] \rightarrow r_1[A] \rightarrow w_1[A] \rightarrow r_0[C] \rightarrow w_0[C] \rightarrow c_0 \rightarrow c_1$

TRUE FALSE

- (c) (2 points) A NO-STEAL buffer manager policy means that all pages modified by a transaction are forced to disk before the transaction commits.

TRUE FALSE

- (d) (2 points) Strict two-phase locking (2PL) ensures that transactions never deadlock.

TRUE FALSE

- (e) (2 points) Strict two-phase locking (2PL) ensures serializability.

TRUE FALSE

- (f) (2 points) In the ARIES protocol, at the end of the analysis phase, the Dirty Page Table contains the exact list of all pages dirty at the moment of the crash.

TRUE FALSE

- (g) (2 points) The ARIES protocol uses the “repeating history” paradigm, which means that updates for all transactions (committed or otherwise) are redone during the REDO phase.

TRUE FALSE

Name: _____

6. (12 points) **Parallel Data Processing**

(a) (4 points) In a parallel DBMS, why is it difficult to achieve linear speedup and linear scaleup?

(b) (4 points) List two features common to a traditional DBMS and MapReduce.

Clarification: Here, “traditional DBMS” means a traditional *parallel DBMS*.

(c) (4 points) List two features that are different between the two types of systems: i.e., features that are present in one system but not in the other. For example, you can give one feature present in MapReduce but absent in a parallel DBMS and one feature present in a parallel DBMS but missing from MapReduce (or any other combination).

Name: _____

7. (8 points) **Database as a Service**

Frank and Betty own a small Internet based company that sells collector pens over the web. Recently, they decided to get rid of all their infrastructure and move to using Amazon Web Services. As part of this transformation, they plan to get rid of their DBMS and run their application on top of Amazon SimpleDB.

(a) (4 points) List three potential benefits of this move.

(b) (4 points) List three potential challenges that they will face.

Name: _____

8. (12 points) **Column-Stores**

(a) (4 points) Explain why column-oriented DBMSs are advantageous for OLAP workloads compared with row-oriented DBMSs?

(b) (4 points) Describe one approach to simulate a column-oriented DBMS in a row-oriented DBMS.

(c) (4 points) Explain why such simulation yields worse performance compared with using a column-store directly?