COP 4710 Database Systems

Fall 2010

Final Exam

December 9, 2010

Time: 120 minutes

Your Name: _____

USF ID:  __U_____

Remember: always write down your intermediate results for partial credits.

## Problem I. Basic concepts (10pts)

True or false, mark your choice clearly. Ambiguous answers will get zero points.

1. SQL is an implementation of domain relational calculus.  T (   ) F (   )

2. In a relation, a superkey must be a candidate key.  T (   ) F (   )

3. In performing natural joins between table A and table B, if we have an index on the joining attribute(s) of table A, we should use A as the OUTER table in the nested-loop join algorithms to gain better performance.  T (   ) F (   )

4. A valid relational algebraic expression always returns a relation.  T (   ) F (   )

5. The resulting table of a relational algebraic expression is duplicate free.  
T (   ) F (   )

6. Relational calculus is a non-procedural language.  T (   ) F (   )

7. If a relation is in Boyce-Codd Normal Form, it must also be in $3^{rd}$ Normal Form.  
T (   ) F (   )

8. The schema of a large table A is decomposed into two smaller tables A' and A'' by **lossless** decomposition. If we join A' with A'', we can actually get more tuples than those in the original table A.  
T (   ) F (   )

9. A DBMS is a software system that enables users to create and maintain a database.  T (   ) F (   )

10. Given two relations R and S with x and y attributes in their schemas, respectively, the Cartesian product of R and S must have exactly x + y attributes.  
T (   ) F (   )


## Problem II. Query languages (27pts, 3 pts each)

Given the following relational database

> Movies (<u>title</u>, year, length, genre, studioName, ProducerName)  
> StarsIn (<u>movieTitle, starName</u>, movieYear)  
> MovieStar (<u>name</u>, address, gender, birthdate)

where the underlined attributes are the primary keys, and the following foreign keys:

> StarsIn.starName → MovieStar.name  
> StarsIn.movieTitle → Movies.title

Movies.ProducerName → MovieStar.name

You can make reasonable assumptions about the database, please state the assumptions clearly if you do.

**II-A. Write the following queries in <u>SQL</u>:**

1. Print the year of the movie(s) titled "2012";

2. Print the title and year of the movie(s) produced by MGM in 2003;

3. Print the names of actors who have starred in a Spielberg (who acted as a producer) movie in 2003;

4. How many movies did Disney produce in each year since 2000? (Hint: first solve "how many movies did Disney produce in each year")

5. Print the title and year of all movies Tom Hanks got involved (as a producer or an actor). (Hint: find the movies Tom acted in, then find the movies he was a producer of, …)

6. Print the names of movie stars who have acted in ALL categories (genre) of movies.

**II-B. Write the following queries in <u>relational algebra</u>:**

7. Print the title and year of the movie(s) produced by MGM in 2003;

8. Print the names of actors who starred in the movie "Casablanca" but NOT in the movie "Spellbound".

9. Print the names of movie stars who have acted in ALL categories (genre) of movies.

**Problem III. (20 pts, 4 pts each) Database Design**

Consider a relation R(A, B, C, D, E) and the following set of functional dependencies:

$$AB \rightarrow C, \ CD \rightarrow E, \ DE \rightarrow B.$$

Finish the following tasks.

1. Compute (AB)+;

2. Compute (AD)+;

3. Find all the candidate keys of relation R;

4. Does R follow Boyce-Codd Normal Form? Explain how you reach your conclusion.

5. What is the highest normal form R follows? Explain how you reach your conclusion.

**Problem IV. Query processing (18 pts)**

Consider a natural join between relation $R$ and relation $S$, and the following information about the relations to be joined. The cost metric is the number of page I/Os unless otherwise noted, and the cost of writing out the result should be uniformly ignored.

>    Relation R contains 10,000 data pages and each page contains 10 tuples.
>    Relation S contains 5,000 data pages and each page has 20 tuples in it.
>    Both relations are stored as simple heap files.
>    102 buffer pages are available.

Answer the following question with your justifications.

1. (4pts) What is the cost of joining R and S using a page-based simple nested loop join?

2. (4 pts) What is the cost of joining R and S using a block-based nested loops join?

Let us further assume that we have built B+-tree indexes on the join attributes for both R and S.  We obviously can perform index-based loop join.

3.  (5 pts) To get the best join performance, which index would you use? Briefly explain your answer.

4. (5 pts) What is the cost of performing the join with the index you choose?
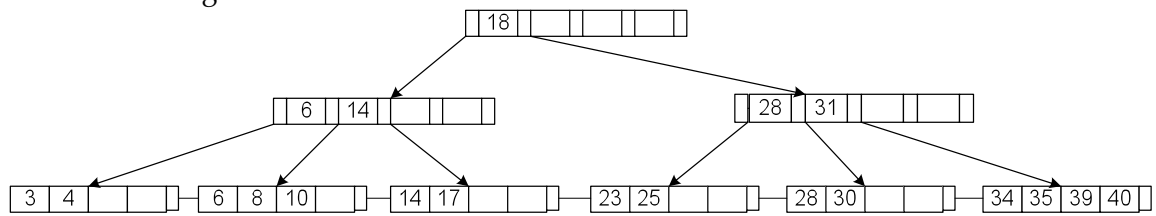
**Problem V. Indexing (25pts)**

1.  (10pts) Construct a B$^+$-tree for the following set of key values:

    27, 20, 43, 40, 67, 17, 55, 33, 11, 75, 48, 9, 95, 60, 65, 107, 98

Assume the tree is initially empty and the values are added in the order specified by the above list.  The maximum number of data entries that will fit in one node is FOUR (i.e., d = 2 as we discussed in class). Draw the intermediate trees for partial credits.

Given the following B$^+$-tree:



2. (2pts) Draw the tree after deleting value 10 from the above tree;

3. (4pts) Draw the tree after deleting key value 3 from the **original** tree (not the one you got from question 2) shown in the figure.

4.  (5pts) Draw the tree after deleting key value 23 from the **original** tree (not the one you got from question 3) shown in the figure.

5.  (4 pts) Starting from the original tree, list two values such that the height of the tree will decrease by one when these two values are deleted from the tree.