

Name

CWID

Quiz

1

September 22th, 2015
Due September 29th, 11:59pm

Quiz 1: CS525 - Advanced Database Organization

Please leave this empty!

1.1

1.2

1.3

1.4

Sum

Instructions

- **You have to hand in the assignment using your bitbucket account**
- **This is an individual and not a group assignment**
- Multiple choice questions are graded in the following way: You get points for correct answers and points subtracted for wrong answers. The minimum points for each questions is **0**. For example, assume there is a multiple choice question with 6 answers - each may be correct or incorrect - and each answer gives 1 point. If you answer 3 questions correct and 3 incorrect you get 0 points. If you answer 4 questions correct and 2 incorrect you get 2 points. ...
- For your convenience the number of points for each part and questions are shown in parenthesis.
- There are 4 parts in this quiz
 1. SQL
 2. Relational Algebra
 3. Index Structures
 4. Result Size Estimation

Part 1.1 SQL (Total: 31 + 10 bonus points Points)

Consider the following disaster event database schema and example instance. **The example data should not be used to formulate queries. SQL statements that you write should return the correct result for every possible instance of the schema!**

city

name	state	population
Chicago	IL	6,500,200
Honolulu	HI	300,000
Madison	WI	400,000

event

city	year	eventType	deaths
Honolulu	1982	volcano	10
Chicago	1871	fire	300
Chicago	1915	ship wreck	844

predictions

city	eventType	deaths
Chicago	ice age	5,000,000
Honolulu	volcano	150
Madison	riot	1000

measures

type	cost	effectPerc	worksFor
police	\$5,000	95	riot
firemen	\$3,000	80	fire
firemen	\$3,000	10	volcano
atomic heater	\$100,000,000	70	ice age

Hints:

- Attributes with black background are the primary key attributes of a relation
- The attribute *city* of relation *event* and *predictions* are both foreign keys to attribute *name* of relation *city*.

Question 1.1.1 (2 Points)

Write a query that returns cities with predicted disasters (relation `predictions`) which have not been the site of disasters in the past (relation `event`). Return each such city only once.

Question 1.1.2 (3 Points)

Write an SQL query that returns a rolling sum of the total number of deaths per city by year, i.e., in the example database there were two disasters in Chicago in 1871 and 1915. Thus, the output for Chicago would be (1871,300) and (1915,1144).

Question 1.1.3 (2 Points)

Return a list of cities ordered by safety. The safety of a city is the number of deaths in the past and predicted disasters divided by the number of residents, e.g., if a city has 100 total deaths and 1000 inhabitants then its safety is 1/10. Note that lower safety values are better.

Question 1.1.4 (4 Points)

Write an SQL query that returns events (all attributes from the event table) such that there where no other events in the same city for at least 100 years.

Question 1.1.5 (3 Points)

Write an SQL query that returns the names of cities where all types of disasters existing in the database have taken place (e.g., in the example instance the disaster types are `volcano`, `fire`, `ship wreck`, `ice age`, and `riot`).

Question 1.1.6 (3 Points)

Write an SQL query that returns the population of cities in the future assuming that all predicted disasters will occur and that the city population will not change otherwise.

Question 1.1.7 (4 Points)

Return the name of the city/cities with the most disasters.

Question 1.1.8 (4 Points)

Write an SQL query that returns the names of cities without any past (events) and future (predictions) disasters.

Question 1.1.9 (6 Points)

Write a query that returns the most effective combination of up to 3 measures for future predicted disasters for which the combined cost is less than \$1,000,000. Measures can only be applied to the disaster types indicated in the `measures` table (attribute `worksFor`). The effectiveness of a set of measures is the sum of the prevented deaths of the measure. The number of prevented deaths is computed as the number of deaths of a disaster multiplied by the `effectPerc` of the measure.

Hint: This is a relatively complex query. Recall that you can use `WITH` in SQL to define temporary views.

Question 1.1.10 Optional Bonus Question (10 bonus points Points)

Write an interpreter of stack operations as a **recursive** SQL query. The initial stack state is stored in a relation **stack(pos,element)** that records which stack position stores which element (starting from position 0). The operations to be executed by your interpreter are stored in a relation **ops(seq,op,element)** where **seq** stores the order of operations, **op** is the type of operation, and **element** is the element which is used by the operation. You have to support the following operations: 1) **pop** removes the top element (at position 0) from the stack and 2) **push** pushes the element stored in attribute **element** onto the stack. The result of your query should be the new content of the **stack** relation. An example instance of the **stack** and **ops** relations are shown below.

stack

pos	element
0	5
1	19
2	23

event

seq	op	element
0	pop	
1	push	13
2	push	35
3	push	37

Part 1.2 Relational Algebra (Total: 29 Points)

Question 1.2.1 Relational Algebra (3 Points)

Write a relational algebra expression over the schema from the SQL part (part 1) that returns the number of deaths for all volcano and riot disasters.

Question 1.2.2 Relational Algebra (4 Points)

Write a relational algebra expression over the schema from the SQL part (part 1) that returns the cities with more than 2 predicted disasters.

Question 1.2.3 Relational Algebra (4 Points)

Write a relational algebra expression over the schema from the SQL part (part 1) that returns cities with volcano eruptions but without riots.

Question 1.2.4 SQL \rightarrow Relational Algebra (3 Points)

Translate the SQL query from Question 1.1.3 into relational algebra (bag semantics).

Question 1.2.5 SQL \rightarrow Relational Algebra (5 Points)

Translate the SQL query from question 1.1.4 into relational algebra (bag semantics).

Question 1.2.6 SQL \rightarrow Relational Algebra (6 Points)

Translate the SQL query from question 1.1.5 into relational algebra (bag semantics).

Question 1.2.7 Equivalences (4 Points)

Consider the following relation schemas:

$R(A, B)$, $S(B, C)$, $T(C, D)$.

Check equivalences that are correct under **set semantics**. For example $R \bowtie R \equiv R$ should be checked, whereas $R \equiv S$ should not be checked.

- $\sigma_{A=5}(R \bowtie_{B=C} T) \equiv \sigma_{A=5}(R) \bowtie_{B=C} T$
- $R \triangleright S \equiv S \bowtie R$
- $R \bowtie (S \cup T) \equiv (R \bowtie S) \cup (R \bowtie \rho_{B \leftarrow C}(T))$
- $\pi_A(R \bowtie S) \equiv \pi_A(S \bowtie R)$
- $R - (S - T) \equiv (R \cup T) - S$
- $(R \cap S) \cup (R \cap T) \equiv R \cap (S \cup T)$
- $\sigma_{A < 3}(\sigma_{B < 4}(R)) \equiv \sigma_{A < 3}(R) \bowtie \sigma_{B < 4}(R)$
- $\delta(R) \equiv \alpha_{A, B}(R)$

Part 1.3 Index Structures (Total: 30 Points)

Assume that you have the following table:

Item		
SSN	name	age
1	Pete	13
2	Bob	23
45	Alice	77
44	John	49
43	Joe	45
42	Lily	3
88	Gertrud	29
89	Heinz	14

Question 1.3.1 Construction (12 Points)

Create a B+-tree for table *Item* on key *SSN* with $n = 2$ (up to two keys per node). You should start with an empty B+-tree and insert the keys in the order shown in the table above. Write down the resulting B+-tree after each step.

When splitting or merging nodes follow these conventions:

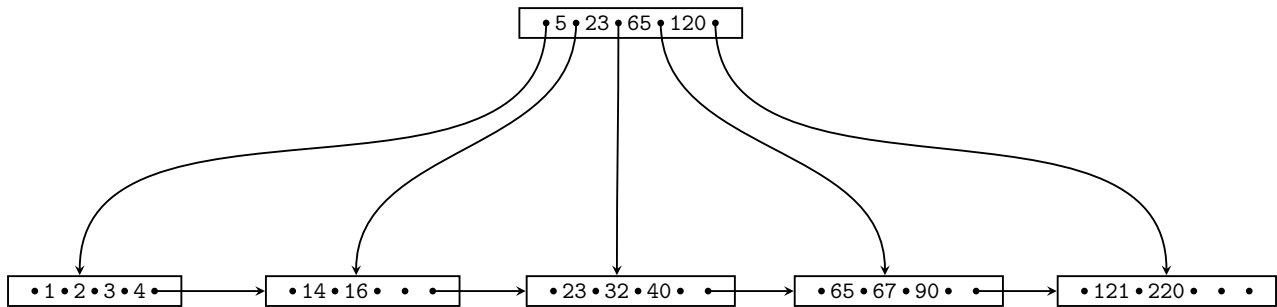
- **Leaf Split:** In case a leaf node needs to be split during insertion and n is even, the left node should get the extra key. E.g, if $n = 2$ and we insert a key 4 into a node [1,5], then the resulting nodes should be [1,4] and [5]. For odd values of n we can always evenly split the keys between the two nodes. In both cases the value inserted into the parent is the smallest value of the right node.
- **Non-Leaf Split:** In case a non-leaf node needs to be split and n is odd, we cannot split the node evenly (one of the new nodes will have one more key). In this case the “middle” value inserted into the parent should be taken from the right node. E.g., if $n = 3$ and we have to split a non-leaf node [1,3,4,5], the resulting nodes would be [1,3] and [5]. The value inserted into the parent would be 4.
- **Node Underflow:** In case of a node underflow you should first try to redistribute values from a sibling and only if this fails merge the node with one of its siblings. Both approaches should prefer the left sibling. E.g., if we can borrow values from both the left and right sibling, you should borrow from the left one.

Question 1.3.2 Operations (10 Points)

Given is the B+-tree shown below ($n = 4$). Execute the following operations and write down the resulting B+-tree after each operation:

insert(44), insert(59), insert(5), delete(121), delete(200), inser(300)

Use the conventions for splitting and merging introduced in the previous question.



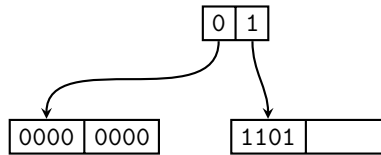
Question 1.3.3 Extensible Hashing (8 Points)

Consider the extensible Hash index shown below that is the result of inserting values 2, 7, and 4. Each page holds two keys. Execute the following operations

`insert(1), insert(5), insert(6), insert(8), delete(2)`

and write down the resulting index after each operation. Assume the hash function is defined as:

x	h(x)
0	1100
1	0001
2	0000
3	1010
4	1101
5	0111
6	1110
7	0000
8	1010



Part 1.4 Result Size Estimations (Total: 10 Points)

Consider a table *church* with attributes *name*, *city*, *confession*, *capacity*, a table *person* with *name*, *confession*, *age*, and a table *attendsService* with attributes *person* and *church*. *attendsService.person* is a foreign key to *person.name*. Attribute *church* of relation *attendsService* is a foreign key to attribute *name* of relation *church*. Given are the following statistics:

$$\begin{array}{lll} T(\textit{church}) = 50 & T(\textit{person}) = 300,000 & T(\textit{attendsService}) = 250,000 \\ V(\textit{church}, \textit{name}) = 50 & V(\textit{person}, \textit{name}) = 300,000 & V(\textit{attendsService}, \textit{person}) = 250,000 \\ V(\textit{church}, \textit{city}) = 30 & V(\textit{person}, \textit{confession}) = 6 & V(\textit{attendsService}, \textit{church}) = 45 \\ V(\textit{church}, \textit{confession}) = 5 & V(\textit{person}, \textit{age}) = 100 & \\ V(\textit{church}, \textit{capacity}) = 45 & & \end{array}$$

Question 1.4.1 Estimate Result Size (3 Points)

Estimate the number of result tuples for the query $q = \sigma_{\textit{confession}=\textit{catholic}}(\textit{church})$ using the first assumption presented in class (values used in queries are uniformly distributed within the active domain).

Question 1.4.2 Estimate Result Size (3 Points)

Estimate the number of result tuples for the query $q = \sigma_{\textit{age}>30}(\textit{person})$ using the first assumption presented in class. The minimum and maximum values of attribute *age* are 1 and 100.

Question 1.4.3 Estimate Result Size (4 Points)

Estimate the number of result tuples for the query $q = (person \bowtie_{name=person} attendsService \bowtie_{church=church.name} church)$ using the first assumption presented in class.