

# Midterm 1

Saturday, May 18, 2019 4:06 PM



dbms\_quiz

66 / 76 =

87%

Name: \_\_\_\_\_

Class Login: cs186-\_\_\_\_\_

Midterm 1: CS186, Spring 2016

You should receive 1 double-sided answer sheet and an 11-page exam. Mark your name and login on both sides of the answer sheet, and in the blanks above. For each question, place **only your final answer** on the answer sheet—do not show work or formulas there. You may use the blank spaces for scratch paper, but do not tear off any pages. You will turn in both question and answer sheets.

## I. Storage: Disk, Files, Buffers [11 points]

1. [3 points] Write down the letters of true statements *in alphabetical order*. (If none are true, write Ø.) B C D E F

- A. When querying for a 16 byte record, exactly 16 bytes of data is read from disk. F. Disks deal with pages  
B. Writing to an SSD drive is more costly than reading from an SSD drive. T. In fact, it's faster.  
C. In a heap file, all pages must be filled to capacity except the last page. F. Some like a record  
D. If the file size is smaller than the number of buffer frames, a sequential scan of the file using either MRU or LRU (starting with an empty buffer pool) will have the same hit rate.  
E. Assuming integers take 4 bytes and pointers take 4 bytes, a slot directory that is 256 bytes can address 64 records in a page.  
F. In a page containing fixed-length records with no nullable fields, the size of the bitmap never changes.

2. [2 points] Write down the true benefits of using a record header for **variable** length records *in alphabetical order* (or if none are true benefits, write Ø.) A

- A. Does not need delimiter character to separate fields in the records.  
B. Always matches or beats space cost when compared to fixed-length record format.  
C. Can access any field without scanning the entire record.  
D. Has compact representation of null values

Nah you can do w/e u want I guess  
Like supporting deletes and stuff

32. You have 128 bytes for pointers  
33. 128 bytes for records on that page  
34. The rest points

3. [6 points] Assume we have 4 empty buffer frames and the following access pattern, in which pages are immediately unpinned.

T A M E T E A M M A T E M E A T L I D

Use the replacement policy listed, and list the four pages in the buffer pool at the end, *in alphabetical order*. Hint: you don't need to draw a big chart for every access — look for patterns.

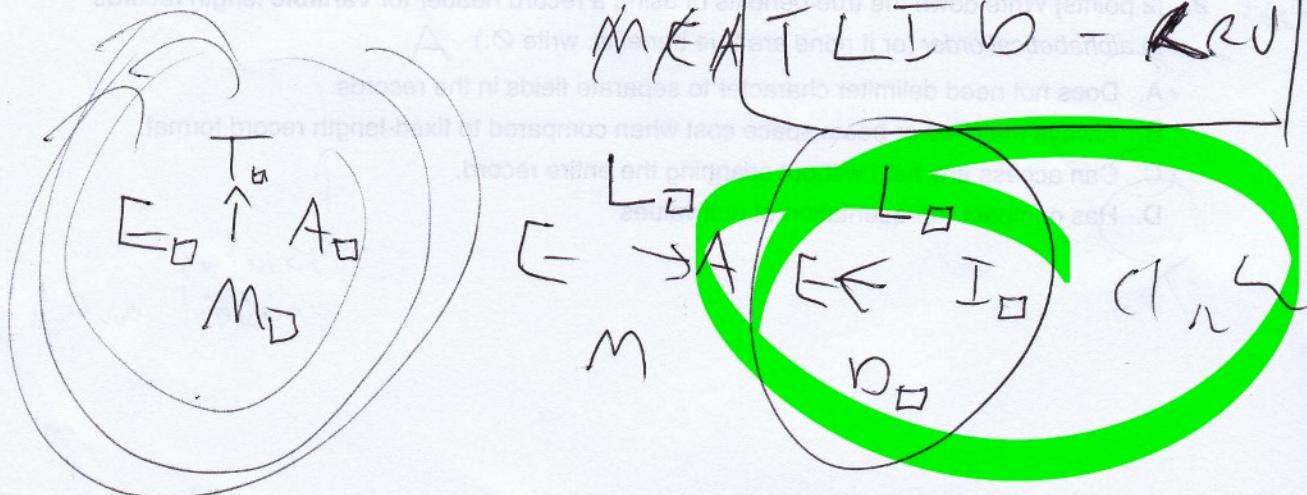
- A. MRU
- B. LRU
- C. Clock. (*Assume the clock hand starts on the first buffer and does not move unless a page needs to be replaced.*)

This space left intentionally blank for scratch work.

XAMTETEAM AMATE MEATLID

MRU → MEAT

MEATLID - LRU



## II. Joins [12 points]

A C D

1. [4 points] Alphabetically, write down the letters of statements that apply (or write Ø.)
- A. Sometimes, adding more memory to our system **will not** reduce I/O costs for a sort-merge join.
  - B. A Grace hash join will always perform better than a naïve hash join.
  - C. Sometimes, replacing an Alternative 2 index with an Alternative 1 index on the same key will speed up an index-nested-loops join.
  - D. A Grace hash join can often complete in 2 passes if the size of the smaller relation is less than roughly the square of the number of buffers available for the join.

For the following questions in this section (Joins), assume that we are streaming our join output to a terminal. Consider the cost of the initial table scan, but **do not consider the cost of writing the final output**.

We have the following schema:

```
CREATE TABLE Cheesemakers (
cm_id INTEGER PRIMARY KEY,
name VARCHAR(33),
ranking INTEGER,
skill_level INTEGER)
```

```
CREATE TABLE Products (
cheese_id INTEGER PRIMARY KEY,
cm_id INTEGER REFERENCES Cheesemakers,
ctype VARCHAR(16),
smelliness INTEGER,
cheesiness INTEGER)
```

Until instructed otherwise, assume that:

- Cheesemakers has **[C] = 500 pages**
- Products has **[P] = 2000 pages**
- we fix **B = 102 pages** of memory for computing joins.

Consider the following query:

```
SELECT C.name, C.ranking, P.ctype, P.smelliness, P.cheesiness
FROM Cheesemakers C, Products P
WHERE C.cm_id = P.cm_id ← Equi join
```

2. [4 points] Using the smaller relation as the “outer” one, what is the I/O cost of using a block nested loops join to evaluate the query above? Please provide the final number.

$$500 * 2000 + 500 = 10500$$

Page 3 of 11

(100+2000) \* 5 = 10500

Load 100 pages of C  
Load each page of P 1 at a time  
Equijoin to a buffer and pump out the results

3. [4 points] What is the I/O cost of using a sort-merge join to evaluate the query above? Please provide the final number. (Remember to take advantage of the "important refinement" discussed in lecture for merge-joining partitions during the last pass of sort!)

*This space left intentionally blank for scratch work.*

Sort in two runs

Q1 2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

2000.2

$$\text{passes} = \log_{39} \left( \frac{2000}{40} \right) + 1$$

### III. Sort/Hash [16 points]

$$60840 \Rightarrow 1521 \rightarrow 40 \text{ runs}$$

For this question, consider our table of Products from the previous Joins section, but assume:

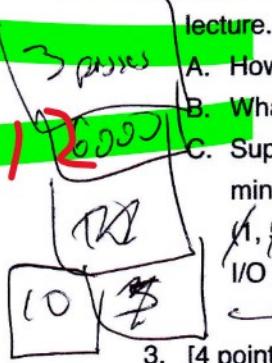
- $[P] = 2000$  pages
- $p_c = 100$  tuples/page
- $B = 40$  pages of buffer
- In all parts, we'll use **QuickSort** for our internal sort algorithm.
- **Include the cost of the initial scan and cost of writing output in your I/O calculations.**

1. [3 points] Alphabetically, write down the letters of statements that apply (or write  $\emptyset$ .)

- A. Given a buffer of size B, the largest file you can sort in a single pass is B.
- B.** All files can be externally sorted, whereas not all files can be externally hashed.
- C. For sort-merge joins, quicksort is always a better choice than heapsort for the internal sort algorithm.

50

2. [6 points] First, let's sort our table of Products (P) using the external algorithm we learned in lecture.



1 pass to make  $\lceil \frac{2000}{40} \rceil$  runs of size 40

- A. How many passes are needed to sort this file?
- B. What is the I/O cost (in pages) of sorting this file?
- C. Suppose we want to decrease the I/O cost of sorting this file, but we want to add the minimum number of buffer pages possible. Among the numbers on the answer sheet (1, 5, 10, 50) circle the *smallest* number of additional buffers to add that decreases the I/O cost.

$16000 / 40$ , not in / out

$2000 / 40$

60840

3. [4 points] Given the resources at the top of Question III, answer the following two questions.

Do not bother to simplify arithmetic expressions over constants, like  $247^*(36^3 + \log(4))$ .

- A. What is the largest file size (in pages) that we can sort in 3 passes?
- B. What is the smallest file size (in pages) that will require 3 passes to sort?
- C. Suppose I want to eliminate duplicates from our (unsorted) table of Products, using external hashing. Write down the letters of true statements. (If none are true, write  $\emptyset$ .)

why?

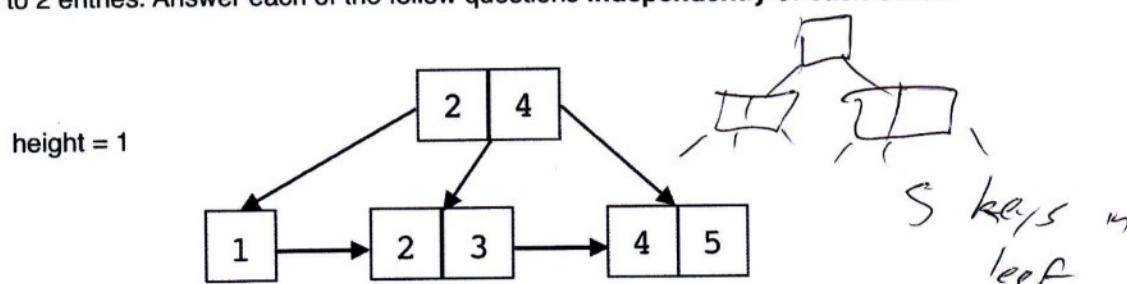
A B

#### IV. Indexes and B+ Trees [12 points]

Note: for B+ tree page splits with an odd number of items, assume that the majority of the items is placed on the right-hand page after the split.

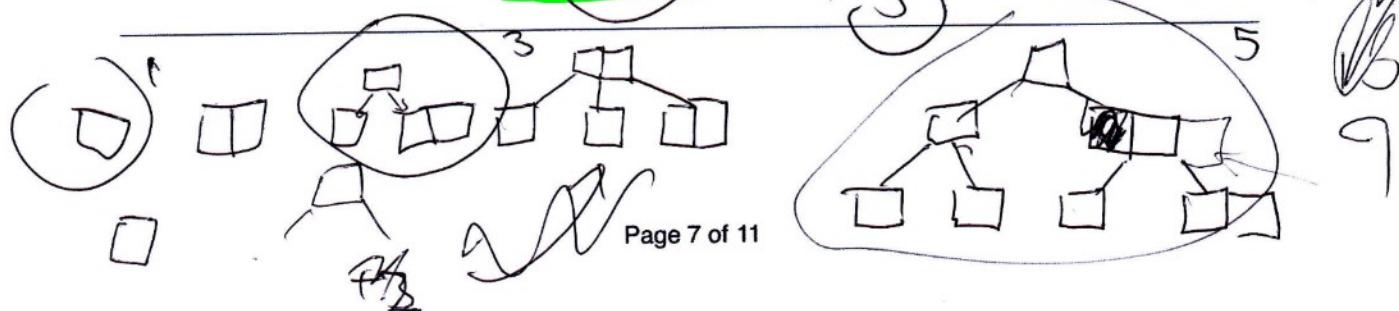
- [4 points] Alphabetically, write down the letters of statements that apply (or write Ø).
  - A. All internal keys in a B+ tree also appear in its leaf nodes. *not if branching de (left)*
  - B. The height of a B+ tree increases whenever any node splits. *No, not necessarily*
  - C. An ISAM index is similar to a B+ tree, but does not allow for insertion of new values. *No, it can't*
  - D. The column(s) we select for our index key must have a unique value for every row in the table. *not necessarily?*
  - E. An Alternative 1 index may be either clustered or unclustered. *? go for yes or no*

- [5 points] The following B+ Tree has order 1 (max fanout of 3) and each leaf node can hold up to 2 entries. Answer each of the follow questions **independently of each other**.



- What value(s) would be in the root node if we were to insert 0? *[2|4]*
- What value(s) would be in the root node if we were to insert 6? *[4]*
- Starting with the height 1 tree in the picture above, suppose we start inserting keys 6, 7, 8, ... and so on. After inserting what key will the height of the tree become 3?

- How many leaf nodes will there be? *6 leaf nodes*
- How many internal (non-leaf) nodes will there be? *1 root + 4 internal*
- How many internal (non-leaf) nodes do we traverse to do an equality search? *2 + 1 = 3*



## V. SQL [ 17 points]

select lname, count(\*)  
 from People, Calls  
 where calls.to = people.id AND people.areaCode = location.areaCode  
 group by areaCode

Consider the following schema:

```
CREATE TABLE Location (
    lname TEXT,
    areacode INTEGER PRIMARY KEY)
```

```
CREATE TABLE People(
    id INTEGER PRIMARY KEY,
    pname TEXT,
    areacode INTEGER REFERENCES Location)
```

```
CREATE TABLE Calls(
    cid INTEGER PRIMARY KEY,
    date DATETIME,
    from INTEGER REFERENCES People,
    to INTEGER REFERENCES People,
    duration INTEGER)
```

Select count(\*), lname  
 from People, Calls  
 where Calls.to = people.id  
 group by lname

You should assume that referential integrity is being enforced, and no NULL values appear.

For parts 1 and 2, fill in the blanks in the SQL queries.

1. [7 points] Names of pairs of people who called each other on 2014-12-25

2. [7 points] Find the location to which the most phone calls have been made, and return its location name as well as the number of calls made to it.

① with caller\_Pair as (from, to) AS  
 (Select p1.pname, p2.pname  
 From Calls as C1, Calls as C2, People as P1, People as P2  
 Where C1.from = C2.to AND C1.to = C2.from  
 and C1.from < C2.from AND  
 C1.from = P1.id AND C1.to = P2.id)  
 Select p1.pname, p2.pname  
 From caller\_Pair, People as P1, People as P2  
 Where P1.id = from AND P2.id = to

3. [3 points] On the answer sheet, alphabetically write down letters of the statement(s) that find the name of the location of the person who made the longest call. (If none match, write Ø.)

A. WITH (SELECT P.areacode AS areacode, C.duration AS duration  
FROM calls C, People P  
WHERE C.from = P.id) AS DC,  
SELECT L.lname  
FROM DC, Location L  
WHERE L.areacode = DC.areacode  
ORDER BY DC.duration DESC  
LIMIT 1;

B. SELECT L.lname  
FROM Location L, (SELECT P.areacode AS areacode,  
C.duration AS duration  
FROM Calls C, People P  
WHERE C.from = P.id) AS DC  
WHERE L.areacode = DC.areacode  
ORDER BY DC.duration ASC;

C. WITH (SELECT L.lname AS lname, L.areacode AS LAC, P.id AS pid  
FROM Location L FULL OUTER JOIN People P  
ON L.areacode = P.areacode) AS Names,  
SELECT Names.lname  
FROM Names, Calls C  
WHERE C.cid = Names.pid  
ORDER BY C.duration DESC  
LIMIT 1;

## VI. Relational Algebra [8 points]

Consider the schema from the SQL question, but with the relation names shortened:

- Location  $\rightarrow L$
- People  $\rightarrow P$
- Calls  $\rightarrow C$

$$\text{PI\_areacode}(L) - \text{PI\_areacode}(P \text{ JOIN } P.\text{to} = C.\text{id } C)$$

1. To answer the following two questions, put one relational algebra operator or relation name in each blank:
  - a. [3 points] Areacodes to which no call has ever been made.
  - b. [3 points] Names of people who live in the location "City of Joe"
2. [2 points] The equivalence relation  $\equiv$  means that two expressions produce the same results on all possible databases. Write down the letters of the equivalences below that are true. (If none are true, write  $\emptyset$ .)

A.  $\rho(\text{Temp1}(1 \rightarrow id), (\pi_{id}(P) \cup \pi_{from}(C))$

$$\rho(\text{Temp2}(1 \rightarrow id), ((\sigma_{id=5}(P) \bowtie_{id=from} C))$$

$$\sigma_{id=5}(\text{Temp1}) \equiv \pi_{id}(\text{Temp2})$$

B.  $\pi_{areacode}(P) = \pi_{areacode}(\sigma_{areacode=415}(P \bowtie L))$

$$= \pi_{areacode}(P) - \pi_{areacode}(\sigma_{areacode=415}(\pi_{areacode}(P) \cap \pi_{areacode}(L)))$$

*This space left intentionally blank.*