

CS 377: Database Systems

Homework #4

Due: Friday, April 21, 2017 at 11:59 PM on Gradescope

1. Extensible Hashing ($8 \times 3 + 3 \times 2 = 30$ points)

Assume we have the following records where we indicate the hashed key (binary representation) in brackets. We want to use an Extensible Hashing structure where each bucket can hold up to 3 records and the structure initially starts out empty (one empty bucket).

- I. 'a' [00001]
- II. 'b' [01000]
- III. 'c' [00000]
- IV. 'd' [10100]
- V. 'e' [10101]
- VI. 'f' [00010]
- VII. 'g' [11000]
- VIII. 'h' [10001]

Consider the result after the records above have been inserted in the order shown, using the highest-bits for the hash function. In other words, records in a bucket of local depth d , agree on their *leftmost* d bits.

- (a) What does the hash structure look like after the insert of each tuple. You should have one diagram for each additional tuple.
- (b) What is the global depth of the resulting directory?
- (c) After inserting these records, suppose that we insert 'i' [10000]. How many buckets will we have now?
- (d) What other keys are in the same bucket with record 'i'?

2. B⁺-Tree ($3 \times 6 + 8 = 26$ points)

Consider the B⁺-Tree of order $d = 2$, height $h = 3$ levels, and number of values per node $n = 4$ shown in Figure 1. For each subpart, draw the tree after the specified operation starting from the one (Figure 1) provided below (i.e., the subparts are independent and do not affect one another). Please follow the guidelines listed below:

- The left pointer denotes $<$ while the right pointer denotes \geq .

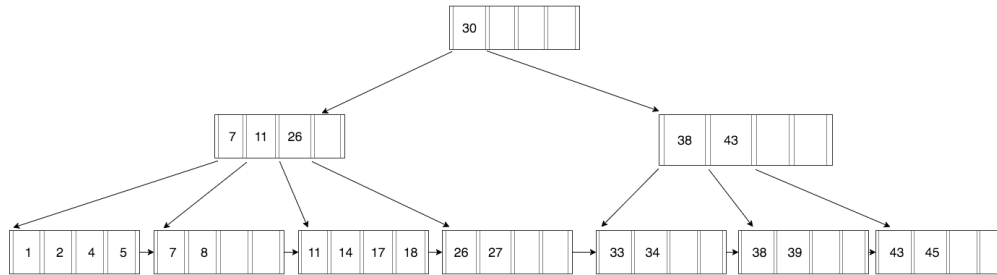


Figure 1: B⁺-Tree for problem 2

- In the case of underflow, if you can borrow from both siblings, choose the one on the *left*.
- In the case of overflow, promote the key in the middle from the left sibling.

- (a) Insert tuple with search key value 10.
- (b) Insert tuple with search key value 20.
- (c) Delete tuple with search key value 7.
- (d) Delete tuples with search key values 17, 18, 26.

3. Using Indexes on Yelp Reviews (4 × 3 = 12 points)

Consider the following Yelp review database:

Business(bid, name, city, state)

BusinessCategory(bid, category)

User(uid, name, fans)

Review(bid, uid, stars, date, votes_useful, votes_funny, votes_cool)

In addition to the primary keys that have been specified, we've also created a B⁺ tree index on the attribute city on the relation business using the following command:

```
CREATE INDEX bcity_idx ON business(city);
```

For the following queries, answer if the index we created was used or not and your logic behind it.

- (a) `SELECT * FROM business
WHERE city > 'Atlanta';`
- (b) `SELECT * FROM business
WHERE city > 'B' AND city < 'G' AND name = 'A';`
- (c) `SELECT * FROM business
WHERE city > 'B' AND city < 'G' OR state = 'GA';`
- (d) `SELECT * FROM business
WHERE city != 'Atlanta';`

4. **Selection Search via B⁺-Tree** ($5 \times 3 = 15$ points)

Consider the following bank database (it was used as an example during the midterm review):

```
Branch(bname, bcity, assets)
Customer(cname, cstreet, ccity)
Loan(lID, bname, amount)
Borrower(cname, lID)
Account(aID, bname, balance)
Depositor(cname, aID)
```

Suppose a B⁺-Tree index was created on bcity for relation Branch, and there are no other indexes in the database. Describe the optimal selection algorithm, specifying whether the index on bcity was used or not:

- (a) $\sigma_{\text{NOT}(\text{bcity} < \text{'Brooklyn'})}(\text{Branch})$
- (b) $\sigma_{\text{NOT}(\text{bcity} = \text{'Brooklyn'})}(\text{Branch})$
- (c) $\sigma_{\text{NOT}(\text{bcity} < \text{'Brooklyn'} \text{ OR } \text{assets} < 5000)}(\text{Branch})$

5. **Estimating Cost of Joins** ($5 + 5 + 10 = 20$ points)

Let relations $R_1(A, B, C)$ and $R_2(C, D, E)$ have the following properties: R_1 has 20,000 tuples, R_2 has 45,000 tuples. 25 tuples of R_1 fit on one block and 30 tuples of R_2 fit on one block. Assume that the join can not be all done in main memory, so it requires disk access. Estimate the cost for each of the following join strategies for $R_1 \bowtie R_2$, assuming that there are M blocks of memory:

- (a) Nested-loop join
- (b) Merge join (assume neither are sorted on the join key but the tuples with the same bucket fit in memory)
- (c) Hash join (assume no overflow occurs and calculate for both recursive partitioning and without recursive partitioning)