

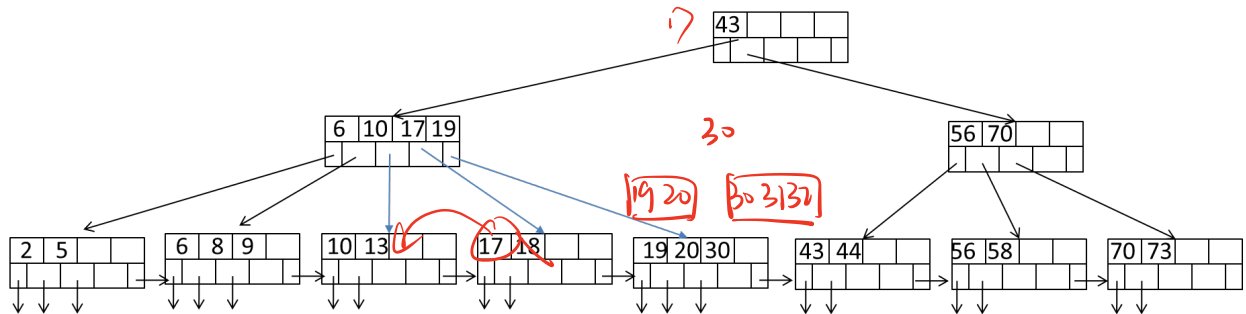
DSCI 551 – HW4

(Indexing and Query Execution)

(Spring 2022)

100 points, Due 4/11, Monday

1. [40 points] Consider the following B+tree for the search key "age". Suppose the degree d of the tree = 2, that is, each node (except for root) must have at least two keys and at most 4 keys. Note that sibling nodes are nodes with the same parent.



- [10 points] Describe the process of finding keys for the query condition "age ≥ 10 and age ≤ 50 ". How many blocks I/O's are needed for the process?
 - [15 points] Draw the B+tree after inserting 31 and 32 into the tree. Only need to show the final tree after the insertions.
 - [15 points] Draw the tree after deleting 18 from the original tree.
2. [60 points] Consider natural-joining tables $R(a, b)$ and $S(a, c)$. Suppose we have the following scenario.
- R is a clustered relation with 1000 blocks.
 - S is a clustered relation with 500 blocks.
 - 102 pages available in main memory for the join.
 - Assume the output of join is given to the next operator in the query execution plan (instead of writing to the disk) and thus the cost of writing the output is ignored.

$$B(R) = 1000 \quad B(S) = 500 \\ M = 102$$

Describe the steps (including input, output, and their sizes at each step, e.g., sizes of runs or buckets) for each of the following join algorithms. What is the total number of block I/O's needed for each algorithm? Which algorithm is most efficient?

- [10 points] (Block-based) nested-loop join with R as the outer relation.
- [10 points] (Block-based) nested-loop join with S as the outer relation.
- [20 points] Sort-merge join (assume only 100 pages are used for sorting and 101 pages for merging). Note that if join can not be done by using only a single merging pass, runs from one or both relations need to be further merged, in order to reduce the number of runs.

Step 1: Sort R into 10 runs, Sort S into 5 runs

$$\text{cost } (1000 + 500) \times 2 = 3000 \text{ blocks}$$

Step 2: merge 15 runs. cost = 1500 blocks \Rightarrow total cost: 4500

Select the relation with a larger number of runs for further merging first if both have too many runs.

- d. [20 points] Partitioned-hash join (assume $M=101$ pages used in partitioning of relations and no hash table is used to lookup in joining tuples).

step 1: hash R into 100 buckets, 10 / bucket
 hash S 100 buckets, 5 / bucket
 $cost = 2B(R) + 2B(S) = 3000$

step 2: join R_i with S_i , $cost = B(R) + B(S) = 150$

