

Data Management Assignment 5

Cuihuan Zhang

May 11, 2025

1. (a)
 - i. Scenarios when insertion from a heap would take 1 operation when the new record is directly insert in front (such like the head of a linked list) or null value of the block, without access the other data.
 - ii. Scenarios when insertion from a heap would take N operations when the heap is compaction or cannot easily access the back of the heap. For example, insert the value in between, then we need to move the rest of the data.
 - (b) Scenarios when insertion from a sorted sequence would take $\log N + N$ operations: a sorted sequence requires $\log N$ using binary search to find the position, then move the rest of the data cost N when the sequence is compaction.
 - (c) Scenarios when insertion from a 2-3 tree would take $\log N$ operations when insertion a integer without insert a new node. Just find the position (node) and insert the value when the node has null value of the block.
2. (a)
 - i. A dense index contains a dedicated pointer (index entry) for every unique key value in the data file (number of indexes = number of records). A sparse index contains pointers only for selected key values (e.g., the first key of each data block) (number of indexes ; number of records).
For example:
Block1: [ID=1, Alice] \rightarrow [ID=2, Bob]
Block2: [ID=3, Mary] \rightarrow [ID=4, David]
Dense index file:
1 \rightarrow Block1-Offset0 (Alice)
2 \rightarrow Block1-Offset1 (Bob)
3 \rightarrow Block2-Offset0 (Mary)
4 \rightarrow Block2-Offset1 (David)
Sparse index file:
1 \rightarrow Block1
3 \rightarrow Block2

- ii. Sparse index requires less storage of the index file, which also leads to less index maintenance operations (no need for index maintenance operations when every time we need to insert or delete a value).
- (b) A clustered index ensures that the physical order of records in the data file matches the logical order of the index keys. An unclustered index is independent of the physical order of data records.
For example:

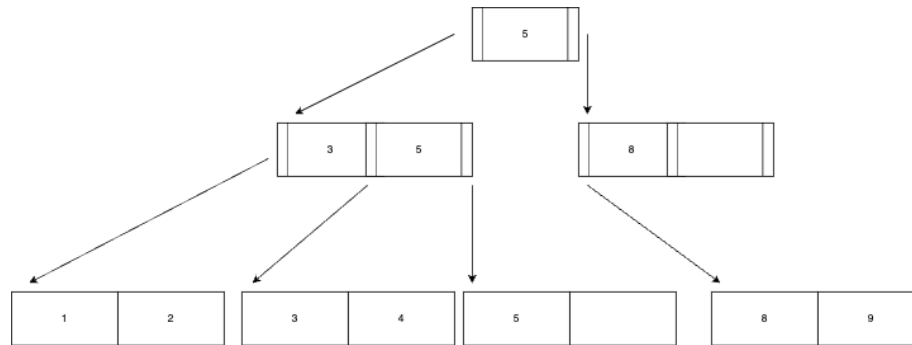


Figure 1: clustered index

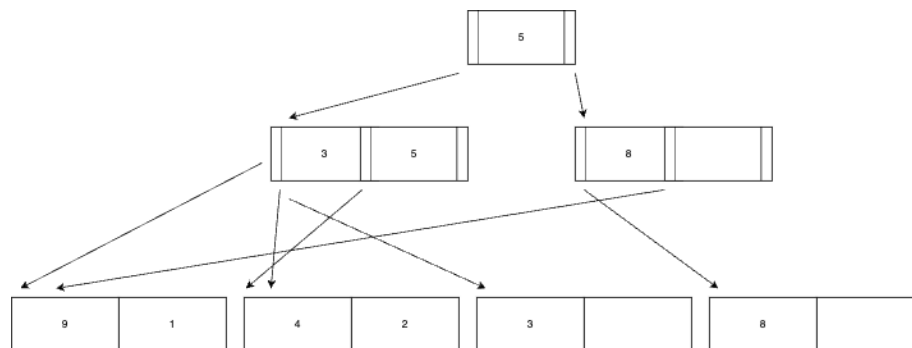


Figure 2: unclustered index

- (c) Best: sparse index and clustered file. High storage efficiency since only record the head of blocks. Less index maintenance operation since when inserting/deleting data, the update frequency of sparse indexes is lower than the dense indexes. High query speed since after locates the block by sparse index, the target record must be in the block.
- Worst: sparse Index and unclustered file. Low query efficiency since after locates the certain block by sparse index, the target record may not include in the block due to the unclustered file. Require more storage since the localization of sparse indexes fails under the unclustered file, and additional redundant index entries may required to

cover the scattered data. For example, require an additional dense and clustered index.

3. (a)
 - **Disk size:** 5 GiB = 5×2^{30} bytes
 - **Block size:** 256 bytes
 - **Key size:** 8 bytes
 - **Pointer size:** ?

Step-by-Step Calculation

1. Total Number of Blocks:

$$\text{Number of blocks} = \frac{\text{Disk size}}{\text{Block size}} = \frac{5 \times 2^{30}}{256} = 5 \times 2^{22} = 20,971,520 \text{ blocks}$$

2. Block Address Size:

$$\log_2(20,971,520) \approx 24.32 \text{ bits}$$

- Round up to 25 bits (requires **4 bytes** ($4 \times 8 = 32$ bits, 1 byte = 8 bits) for alignment)

3. The largest possible number of children: Each node contains m pointers and $m - 1$ keys:

$$m \times (\text{Pointer size}) + (m - 1) \times (\text{Key size}) \leq \text{Block size}$$

Substitute values (4 bytes pointers, 8 bytes keys):

$$4m + 8(m - 1) \leq 256$$

$$12m \leq 264 \implies m \leq 22$$

$$\boxed{m = 22}$$

The smallest number of children are 11, and the largest number of children are 22.

(b)

Level	Nodes in a narrow tree ($m = 11$)	Nodes in a wide tree ($m = 22$)
1	1	1
2	2	22
3	22	484
4	242	10,648
5	2,662	234,256
6	29,282	—
7	322,102	—
8	3,543,122	—

- (c) I will use B+ trees because the quantity involves range queries. Quickly locate all records within the interval through the traversal of the B+ tree structure.

4. Order of joins

- Table R : 2000 blocks.
- Table S : 4000 blocks.
- Available RAM: 6 blocks (2 for R , 4 for S).

(a) Partitioning for S

Each chunk of S uses 4 blocks.

$$\text{Total outer chunks: } \frac{4000}{4} = 1000.$$

(b) Scanning S for Each Chunk of R

For each outer chunk (1000 total) and scan the entire R (2000 blocks).

Total cost for R : $1000 \times 2000 = 2,000,000$ blocks.

(c) Total Cost

$$\text{Total Block Accesses} = 4,000 (S) + 2,000,000 (R) = \boxed{2,004,000}$$