1. BST, AVL trees, and hash tables can all be used to implement a dictionary ADT. $O(\log_2 n)$
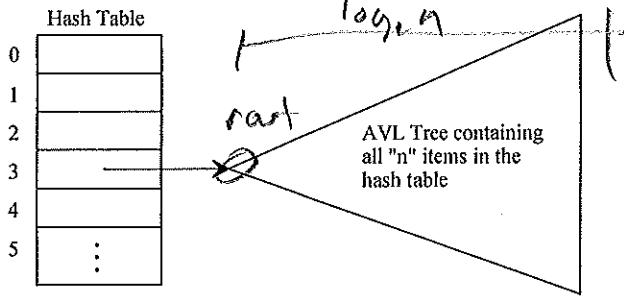
| Dictionary Successful Search Comparisons with 10,000 integer items (Time in seconds) $O(1)$ | | | | | |
|---|---|---|---|---|---|
| | Items added in sorted order | | Items added in random order | | Order did not matter (Hash table sizes $2^{15} = 32K$) |
| | BST | AVL Tree | BST | AVL Tree | Open Addr. (Quadratic) | Closed Addr. (Chaining) |
| Total add/put time | 47.785 | 0.205 | 0.119 | 0.195 | 0.064 | 0.074 |
| Total search time | 38.100 | 0.060 | 0.079 | 0.062 | 0.044 | 0.039 |
| Height of resulting tree | 9,999 | 13 | 30 | 15 | NA | NA |

a) The puts of these 10,000 randomly ordered items into the BST took 0.119 seconds and 0.195 seconds into the AVL tree. Why did the BST puts take less time eventhough the final height was 30 vs. a final AVL tree height of 15?    extra time in AVL to do rebalancing + rotation(s)

b) With a very, very poor hash function or very, very bad choice of keys, then all keys could hash to the same home address.

- What would be the worst-case big-oh of open-address hashing with quadratic probing? $O(n)$

- What would be the worst-case big-oh of chaining using a linked list at each home address (i.e., ChainingDict)? $O(n)$

- What would be the worst-case big-oh of chaining using an AVL tree at each home address?  $O(\log_2 n)$



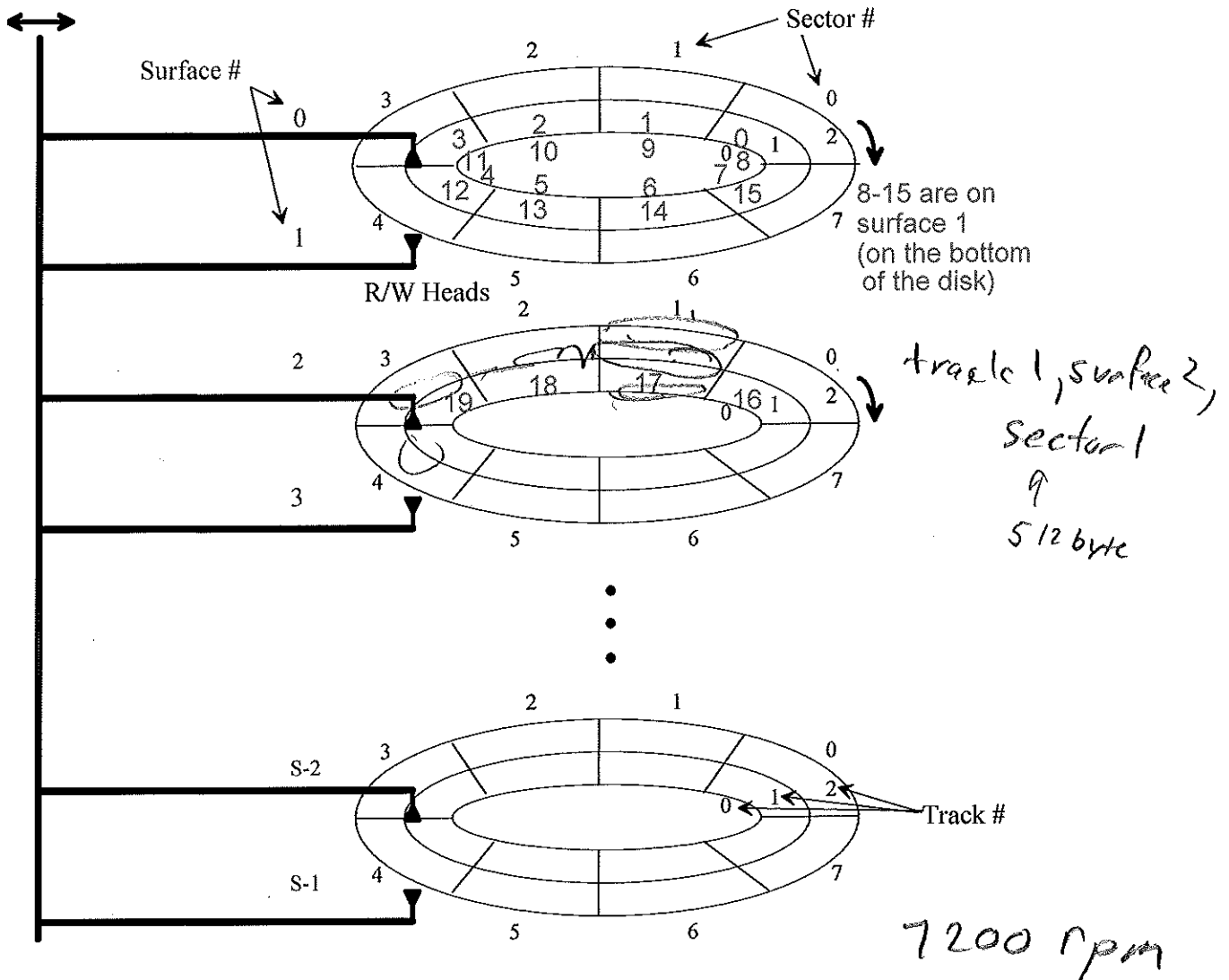Hash Table with AVL Tree containing all "n" items in the hash table; $\log_2 n$ tall, root

2. The data structures we have discussed so far are all in-memory, i.e., data is stored in main/RAM memory. Data can also be stored on secondary storage in a file (e.g., moiveData.txt file). Currently, most secondary storage consists of hard-disks.

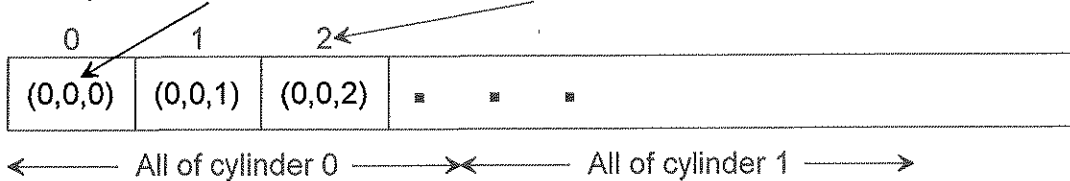a) Complete the following table comparing main/RAM memory vs. hard-disk. $2^{30}$   $2^{40}$

| Criteria | Main/RAM memory | Hard-disk Drive | Solid-State Drive |
|---|---|---|---|
| Size on a typical desktop computer | 8GB | 1 TB | 500GB |
| Average access time | 10 nsec $10 \times 10^{-9}$ sec | 10 X msec $10 \times 10^{-3}$ | 100 $\mu$sec $100 \times 10^{-6}$ |

b) Which criterion seems to be the most important difference between the main and secondary memories?
   — Speed main memory much faster
   — Size – main memory smaller

Logical View of Disk as Linear Collection of Blocks



Sector #

Surface #

0

1

R/W Heads

8-15 are on surface 1 (on the bottom of the disk)

*track 1, surface 2, sector 1*
*9*
*512 byte*

Track #

*7200 rpm*

(track #, surface #, sector #) to Linear block # mapping

| 0 | 1 | 2 | | | | |
|---|---|---|---|---|---|---|
| (0,0,0) | (0,0,1) | (0,0,2) | ▪ | ▪ | ▪ | |

←——— All of cylinder 0 ———→←——— All of cylinder 1 ———→

*seek = move R/w heads*
*largest*
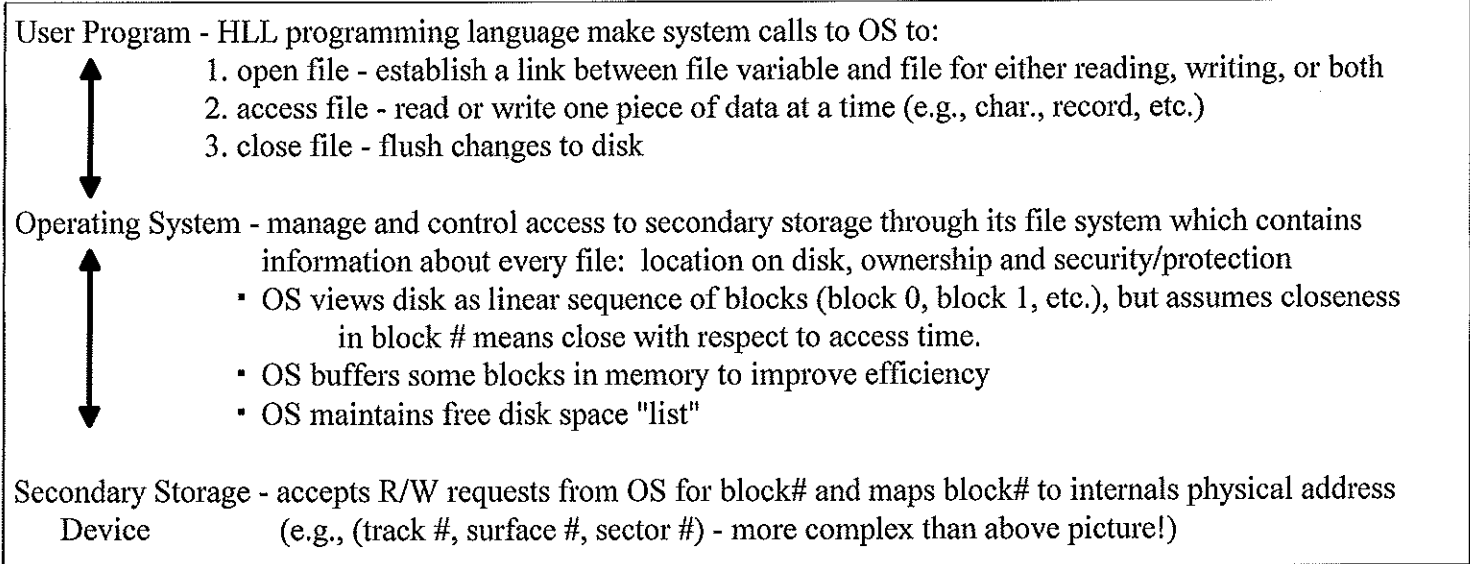
Bits of linear block # :  | track # | surface # | sector # |

3. Disk-access time = (seek time) + (rotational delay) + (date transfer time).   How is each component of the disk-access time effected by increasing the disk's RPMs (revolutions per minute)?

*rotation delay + data transfer decreased, but not seek time*

b) If we want fast access to a collection of sectors, where can we place them to minimize seek time and rotational delay?    *Same cylinder/track to minimize seek time, and consecutive sectors to minimize rotational delay*

User Program - HLL programming language make system calls to OS to:
1. open file - establish a link between file variable and file for either reading, writing, or both
2. access file - read or write one piece of data at a time (e.g., char., record, etc.)
3. close file - flush changes to disk

Operating System - manage and control access to secondary storage through its file system which contains information about every file:  location on disk, ownership and security/protection
- OS views disk as linear sequence of blocks (block 0, block 1, etc.), but assumes closeness in block # means close with respect to access time.
- OS buffers some blocks in memory to improve efficiency
- OS maintains free disk space "list"

Secondary Storage - accepts R/W requests from OS for block# and maps block# to internals physical address
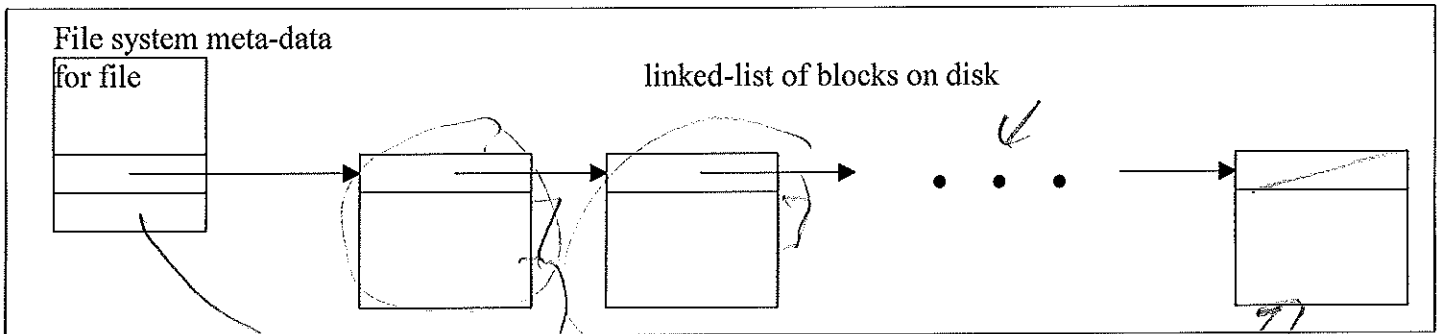Device          (e.g., (track #, surface #, sector #) - more complex than above picture!)

Kinds of File Access:
- serial/sequential files - open at the beginning and read sequentially from beginning to end linearly
- random-access files - "*seek*" to any position by specifying a byte-offset from the beginning of the file, record #, etc.
- random-access of a record by key

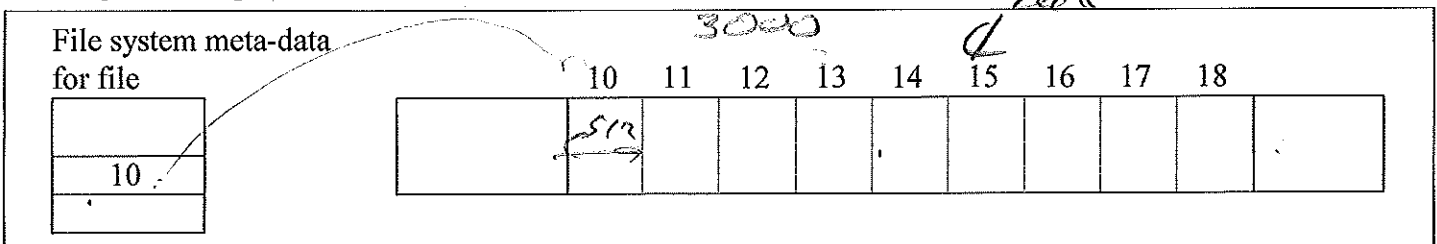Implementation of Files on Disk- how are blocks allocated?
4. non-contiguous - scattered across linear address space of OS and disk

File system meta-data
for file                    linked-list of blocks on disk

a) What types of file access are supported efficiently? *Serial/seq.*

b) How easy is it for the file to grow in size? *easy*

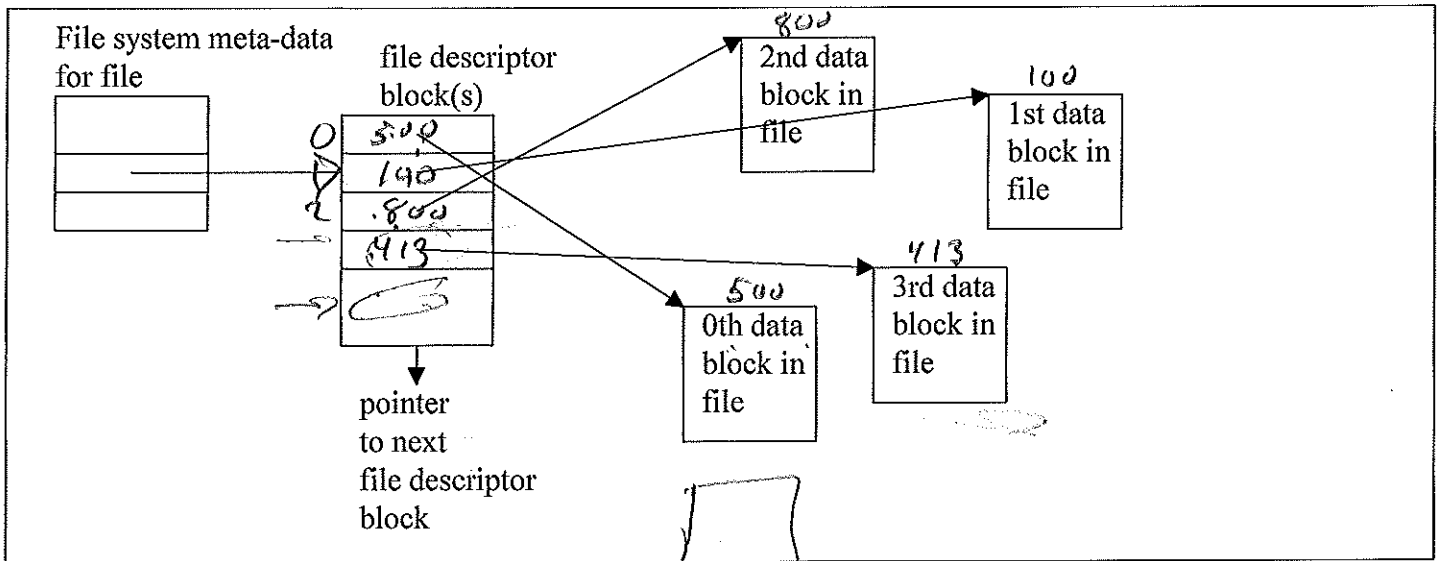5. contiguous - sequential collection of blocks from OS linear view of disk   *seek*

File system meta-data
for file

| | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | |

10

a) What types of file access are supported efficiently? *seq., seek, random key*
binary search $O(\log_2 n)$
hashing $O(1)$

b) How easy is it for the file to grow in size? *hard. might not be room after file to grow file.*

6. file descriptor blocks - list of blocks hold the address of the physical location of data blocks
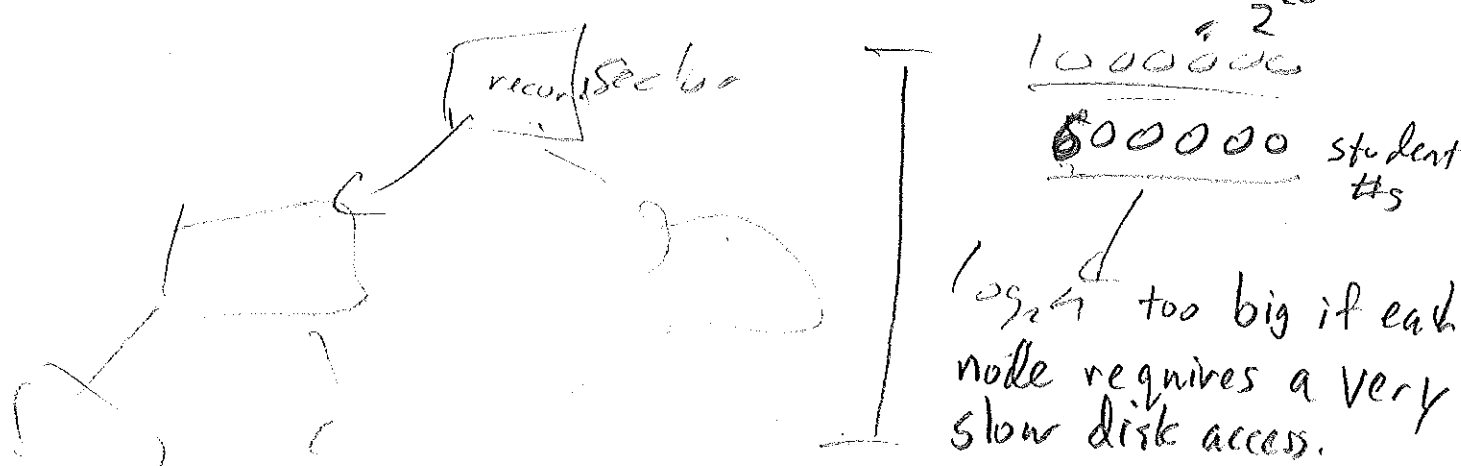


a) What types of file access are supported efficiently? *seq., seek, random-by-key*
   *- binary search*
   *- hashing*

b) How easy is it for the file to grow in size? *easy*

7. To implement "random-access of a record by key" in a file how might we use hashing?

With contiguous or file-descriptor allocation we can run key through hash function to get home address (logical block # in contiguous or index into file-descriptor block) Hopefully, we find key there otherwise need "rehashing", e.g., linear probing.
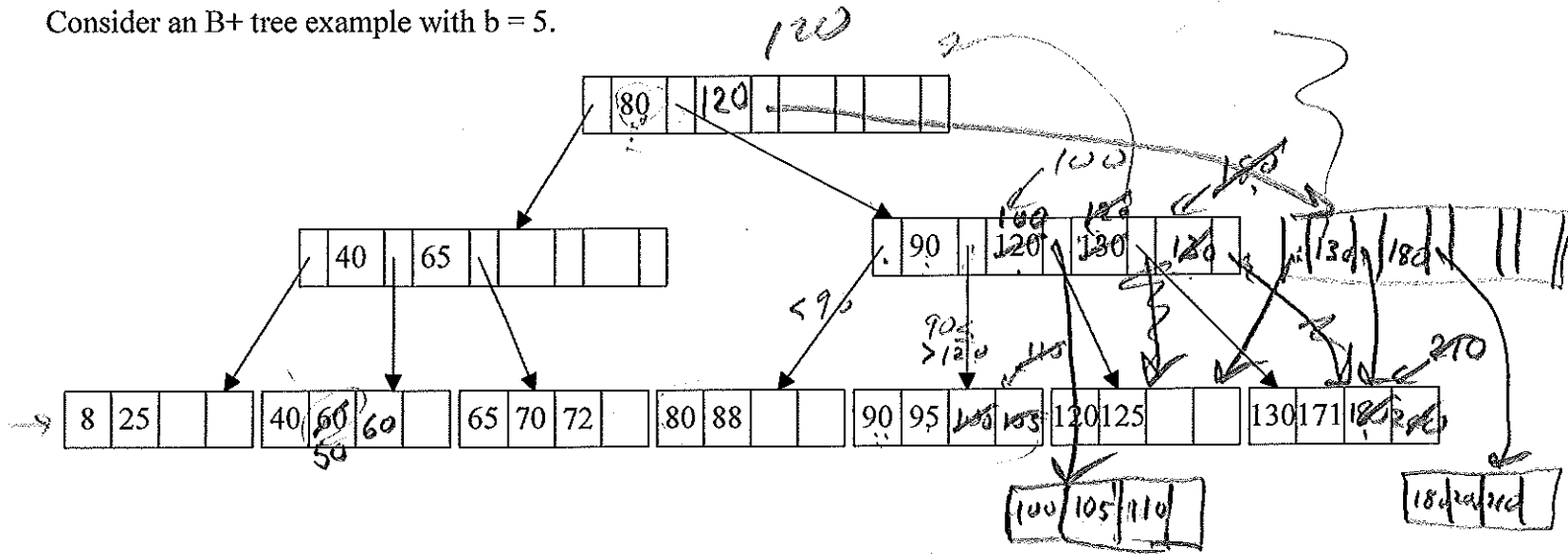
8. To implement "random-access of a record by key" in a file why would an AVL tree not work well?

$2^{20}$
$1000000$
$\overline{600000}$ students

$\log_2^d$ too big if each node requires a very slow disk access.

9. A B+ Tree is a multi-way tree (typically in the order of 100s children per node) used primarily as a file-index structure to allow fast search (as well as insertions and deletions) for a target key on disk. Two types of *pages* (B+ tree "nodes") exist:
- Data pages - which always appear as leaves on the same level of a B+ tree (usually a doubly-linked list too)
- Index pages - the root and other interior nodes above the data page leaves. Index nodes contain some minimum and maximum number of keys and pointers bases on the B+ tree's *branching factor (b)* and *fill factor*. A 50% fill factor would be the minimum for any B+ tree. All index pages must have $\lceil b/2 \rceil \le$ # child $\le b$, except the root which must have at least two children.
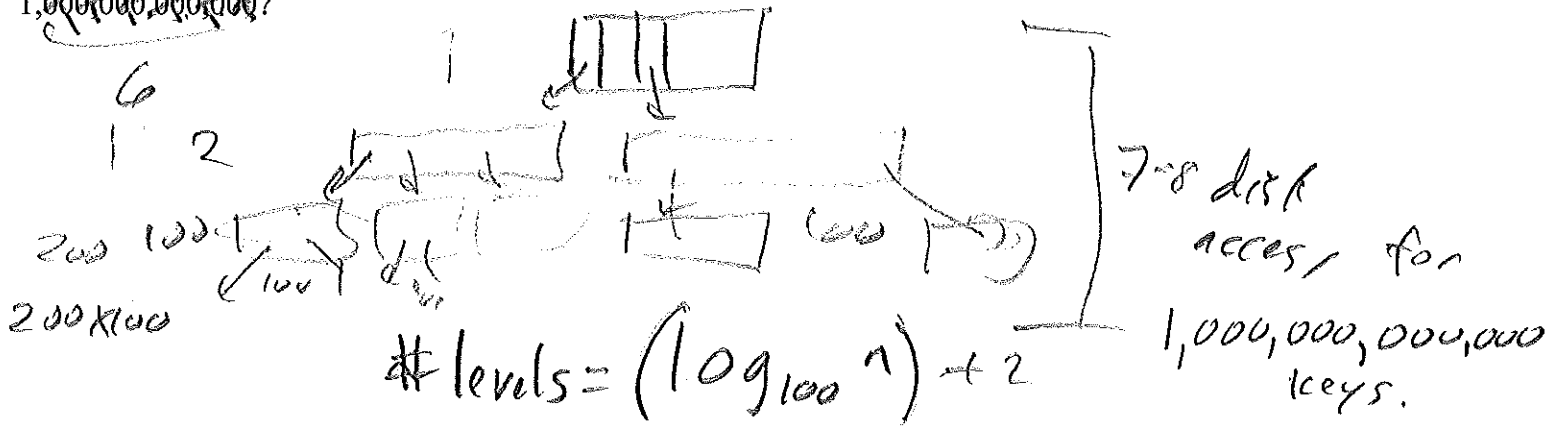
Consider an B+ tree example with b = 5.



a) How would you find 88?  *Start at root 88 > 80, so follow pointer to index page 88 < 90 so follow point left of 90 to data page with 88*

b) The insert algorithm for a B+ tree is summarized by the below table. Where would you insert 50, 100, 105, 110, 180, 200, 210?

| Situation | | insertion Algorithm |
|---|---|---|
| Data Page Full? | Parent Index Page Full? | |
| No | No | Place record in sorted position in the appropriate data page. |
| Yes | No | 1. Split data page with records < middle key going in left data page and records ≥ middle key going in right data page. <br> 2. Place middle key in index page in sorted order with the pointer immediately to its left pointing to the left data page and the pointer immediately to its right pointing to the right data page. |
| Yes | Yes | 1. Split data page with records < middle key going in left data page and records ≥ middle key going in right data page. <br> 2. Adding middle key to parent index page causes it to split with keys < middle key going into the left index page, keys > middle key going in right index page, **and** the middle key inserted into the next higher level index page. If the next higher index page is full continue to splitting index pages up the B+ tree as necessary. |

c) For a B+ tree with a branch factor 201, what would be the worst case height of the tree if the number of keys was 1,000,000,000,000?

*[handwritten work]*

6

1   2

200  100

200 X 100

$$\# levels = \left( \log_{100} ^{\wedge} \right) + 2$$

7-8 disk access for 1,000,000,000,000 keys.

10. The deletion algorithm for a B+ tree is summarized by the below table.

| Situation | | deletion Algorithm |
|---|---|---|
| Data Page Below Fill Factor? | Parent Index Page Below Fill Factor? | |
| No | No | Delete record from the data page. Shifting records with larger keys to left to fill in the hole. If the deleted key appears in the index page, use the next key to replace it. |
| Yes | No | 1. Combine data page and its sibling. Change the index page to reflect the change. |
| Yes | Yes | 1. Combine data page and its sibling. <br> 2. Adjusting the index page to reflect the change causes it to drop below the fill factor, so combine the index page with its sibling. <br> 3. Continue combining the next higher level index pages until you reach an index page with the correct fill factor or you reach the root index page. |

Consider an B+ tree example with b = 5 and 50% fill factor. Delete 89, 65, and 88. What is the resulting B+ tree?



new root level after deletion