

ADB Lecture 2

Indexing

- In our context, a database file stores data records of the same type (table or relation)
- An index is a file that speeds up searches (SELECT) for a search key on the data file (Field)

STUDENT		
Name	SSN	...

• data file

• want to optimize searches of the form
 SELECT ...
 WHERE SSN = ...

→ Create an index file (index Field = SSN)
 ↑
 Search Key (don't confuse with the keys of the relation)

• Much faster



more

- The index file has two columns

Field Value
 (index field)

Pointer to block

Data File (Sorted)

Name	ID	...
Alex	22	...
Bob	53	...
		...

1st block

Dan	37	...
Emy	62	...
		...

2nd block

Yale	12	...
Zack	30	...
		...

Example)

For each block in data file

Name	PTR
Alex	
Dan	
Yale	

PTR to block

Index File
 (Index Key = Name)
 (Field)

- Now everytime we want to search for a record by name, we don't need to search the data file
 → Search the index file and then follow the pointer to the block
 e.g., `SELECT * WHERE NAME = 'Bob'`
 → Use the index to know in which block Bob will be
- binary search on the index file.
- { . It must be the first because it's clear from the index that the 2nd block starts with Dan ($Dan > Bob$)
 * Hence, perform binary search on 1st block to find Bob

- Notice that this approach has assumed that the data file is sorted by the index field. (will generate soon)

→ Ordered File with

$$r = 3 \times 10^5 \text{ records} \quad (\text{no. of records})$$

$$B = 4096 \text{ bytes} \quad (\text{block size})$$

$$R = 100 \text{ byte} \quad (\text{record length})$$

in this case,

$$bfr = \lfloor \frac{B}{R} \rfloor = 40 \text{ records/block}$$

• blocking factor

$$b = \lceil \frac{r}{bfr} \rceil = 7500 \text{ blocks}$$

• # blocks

Hence, the no. of I/O operations is

$$\# Io = \lceil \log_2 b \rceil = 13$$

(disk) 7500

- Let's see how indexing help
 → Suppose the file is ordered by a key of 9 bytes and that block pointer is 3 bytes

- Then each row in Index File has $8+6$ bytes
→ and it has only $\underbrace{7500}_{b=R_{\text{ind}}}$ rows (one for each block)

- How many blocks will we need to store the index file?

$$bfr_{\text{ind}} = \left\lceil \frac{B}{R_{\text{ind}}} \right\rceil = 273 \text{ index/record}$$

$\textcircled{9}$ blocks

$$b_{\text{index}} = \left\lceil \frac{R_{\text{ind}}}{bfr_{\text{ind}}} \right\rceil = 28 \text{ blocks}$$

• # index blocks

$$\# \text{IO} = \underbrace{\lceil \log_2 b_{\text{ind}} \rceil}_{\text{to find the index block}} + 1 = 6$$

Mem.
 $\textcircled{6}$

to read the block
index points to (data block)

* There are 4 types of indexes

- Primary Index
- Clustering Index
- Secondary Index (Key)
- Secondary Index (NonKey)

- What we've seen so far is the Primary index
- Requires that data is ordered by a Key Field (unique or Primary Key)
 - ↳ Candidate Key ↳ unique by definition

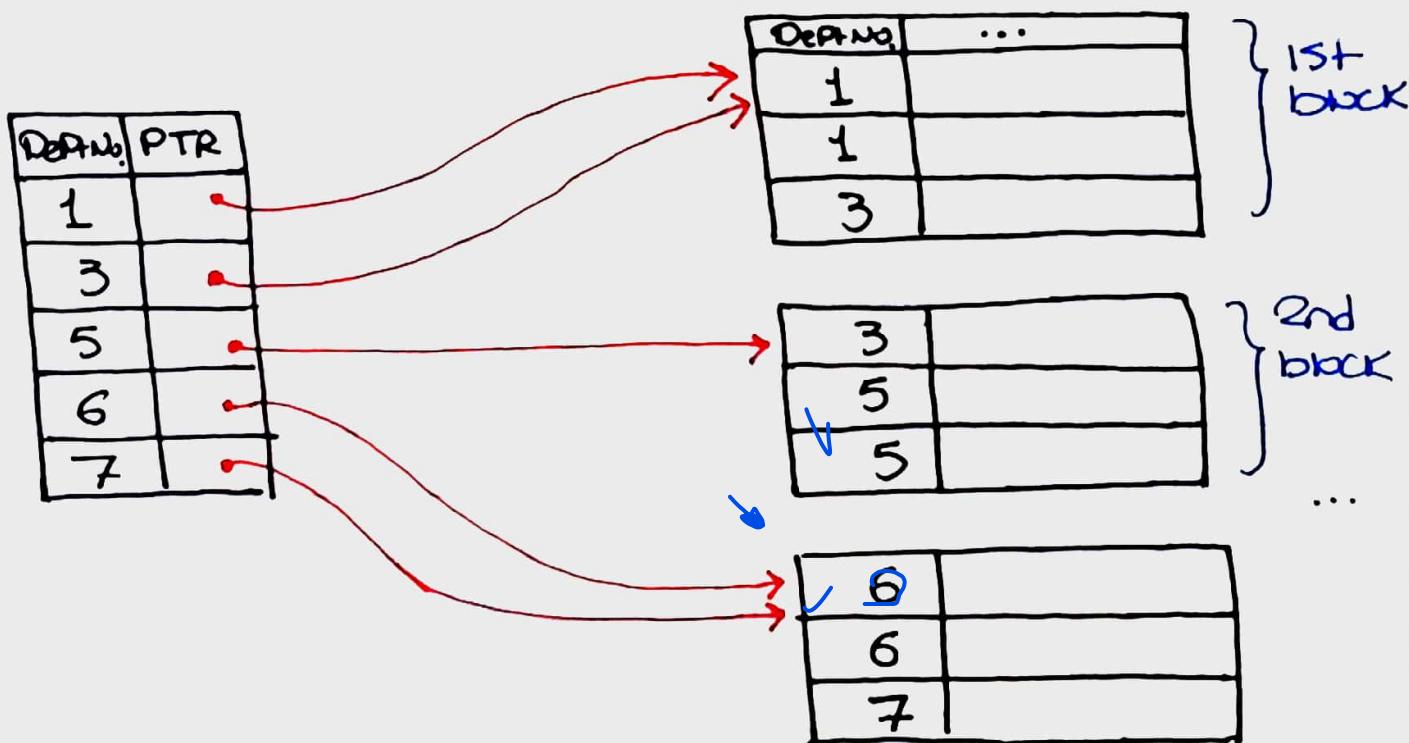
- In this case we create the index using the first element of each block as the value and the block's address (1st record addr.) as pointer "anchor"

- Notice that the index is Sparse (we don't have a pointer for each record in the data file; rather, for each block)

- Which is why the index is smaller (more efficient)

- What if the data is ordered by a field that isn't unique?
→ an index is still possible

→ For each unique value of the field, make a record in the index.



- For each unique value, the block pointer points to the block where that value has 1st occurred

→ Select *
 WHERE DEPTNO. = 3

→ Select

1. Search in index
 For 3 (binary search)
2. Read the 1st block and find 3 in it
3. Since last record is 3 you have to read next block (might be a continuation)

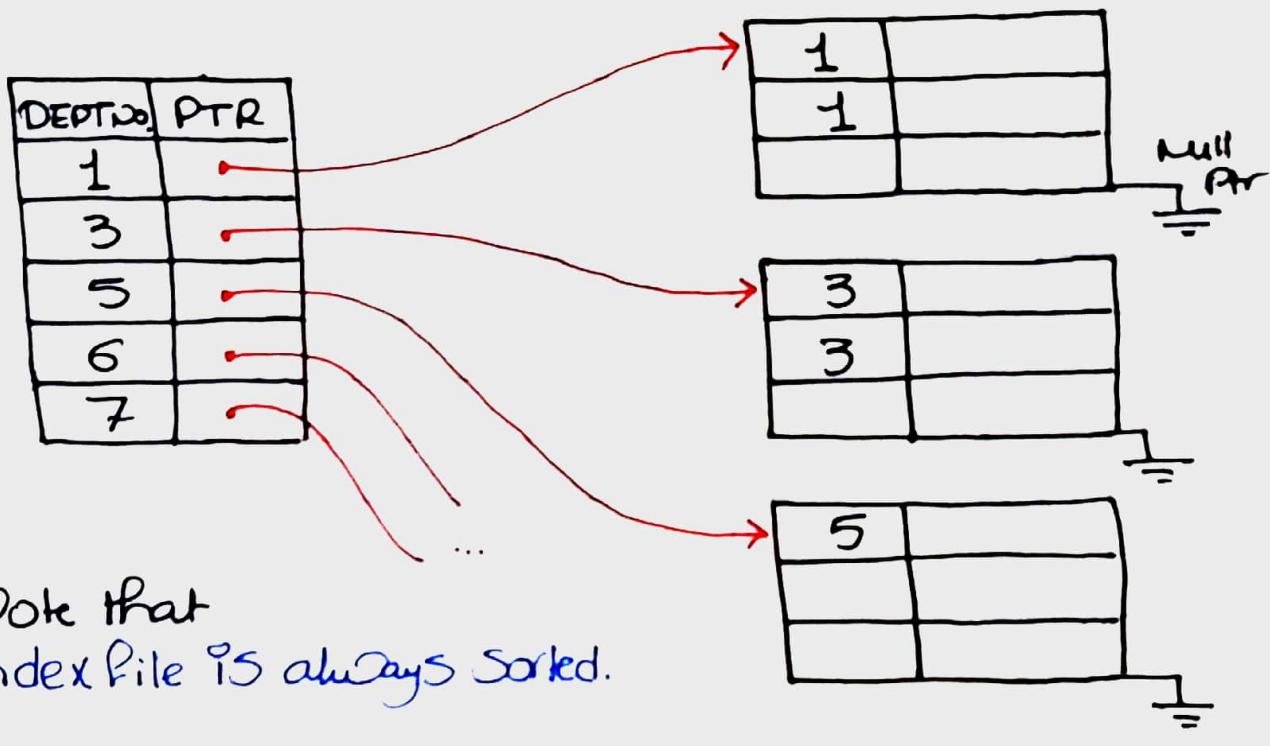
→ Select *
 WHERE DEPT NO. = 5

- Something but
 → 2nd block is returned
 → No need to read next block

- Notice that insertion/deletion in this case isn't so fast (many pointers to move)

Alternative Design

(uses block anchors)



- Note that Index File is always sorted.

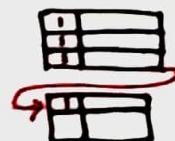
- * In the data file, each group sharing the same value of index (clustering) field should have its own block
 - insert into free space in block (assuming value exists)
 - # That's it
 - # no PTR fix

→ If it doesn't fit

• make a new block

• chain it with the original one (let the original one point to it)

e.g.



* Clearly in this approach

- $\text{Find} = \# \text{unique values} \cdot \text{SparseIndex}$
(#index records)

Example

- Ordered File with

$$\begin{aligned} R &= 3 \times 10^5 \text{ records} \\ B &= 4096 \text{ Bytes} \\ R &= 100 \text{ Bytes} \end{aligned}$$

$$\left. \begin{array}{l} bfr = 40 \text{ records/block} \\ b = 7500 \text{ blocks} \\ \#IO = 13 \end{array} \right\} \quad (\text{last Problem})$$

→ The Ordering attribute is ZipCode

→ Although File has 3×10^5 records, it has only 1000 unique ZipCodes
• i.e., on avg. 300 records have the same ZipCode

* What is #IO for a clustered index.

$$bfr_{ind} = \left\lfloor \frac{B}{R_{ind}} \right\rfloor = 273 \text{ index rec./block}$$

assume 15 (like last Prob.)

$$bind = \left\lceil \frac{r_{ind}}{bfr_{ind}} \right\rceil = 4 \text{ blocks}$$

$$\#IO = \lceil \log_2 bind \rceil + 1 = 3$$

• This will help us only read the 1st block but
→ on average 300 records have same Zip and
each block can hold 40 records. So to read
them all we'll IO in few more blocks.

⑥

* Note that any file can only have one primary index or clustering index (not both)

→ It's either ordered by a unique or non-unique field

- what if we want more?
- what if the file isn't ordered by the search key for which we want queries to be optimized?

* In this case we can use a Secondary Index

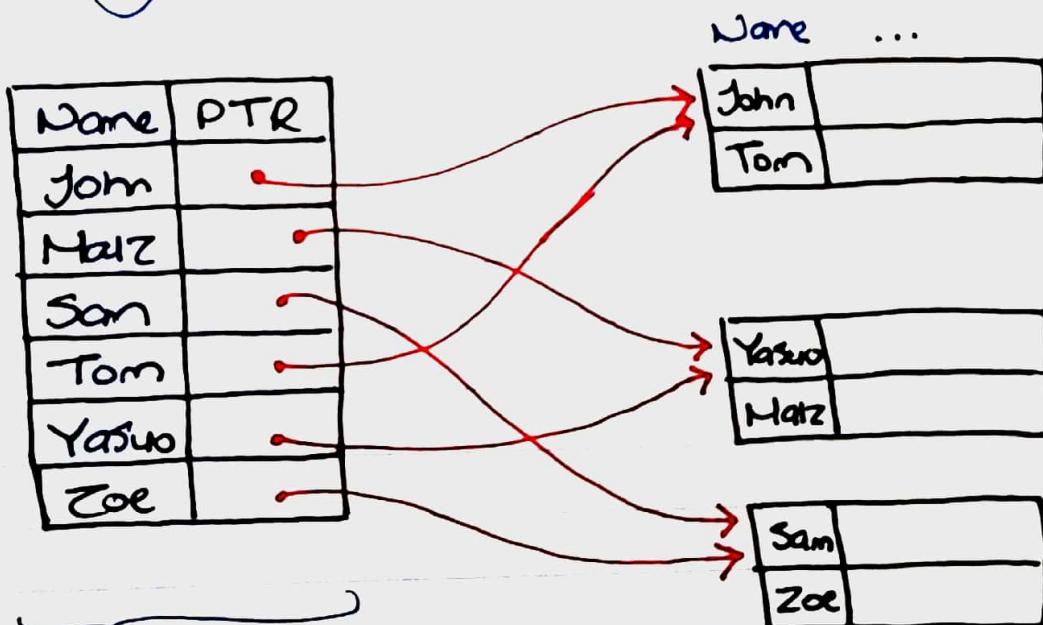
→ in both cases, data file doesn't have to be ordered by index field

- (Key) type
 • index field is unique

(Non-key) type
 • index field is not unique

(i.e., can be created as many times we wish besides the Primary / Clustering index)

Secondary Index (Key)



The index is still sorted

→ has a block pointer for each record in data file (dense)

• Binary Search on it then I/O in the block it points to.

Example

Solved earlier

$$\left. \begin{array}{l} bPr = 40 \text{ records/block} \\ b = 7500 \text{ blocks} \\ \#IO_{\text{sorted}} = 13 \quad (\text{else } 7500/2) \end{array} \right\}$$

$$\left. \begin{array}{l} r = 3 \times 10^5 \text{ records} \\ B = 4096 \text{ Bytes} \\ R = 100 \text{ Bytes} \end{array} \right\}$$

- $R_{\text{ind}} = 15$
- $bPr_{\text{ind}} = \left\lfloor \frac{B}{R_{\text{ind}}} \right\rfloor = 273 \text{ index record/block}$

Solve for # IOs:

$$r_{\text{ind}} = r = 3 \times 10^5 \text{ index records (dense)}$$

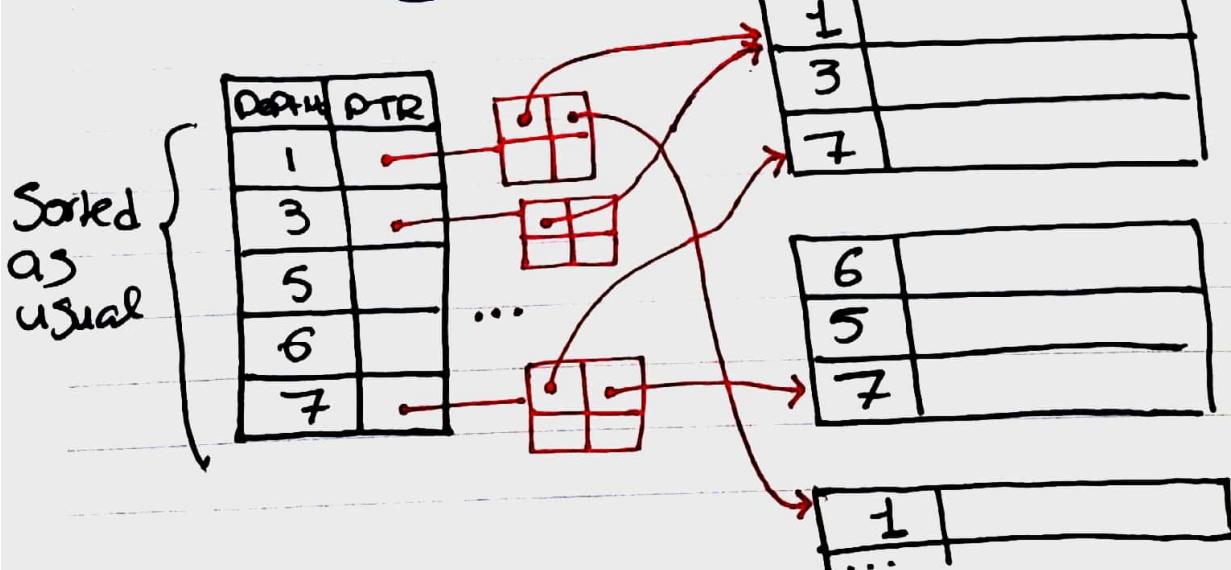
$$b_{\text{ind}} = \lceil \frac{r_{\text{ind}}}{bPr_{\text{ind}}} \rceil = 1099 \text{ block}$$

$$\#IOs = \lceil \log_2 1099 \rceil + 1 = 12$$

. Savings due
to
sort → smaller
record

- What if we want to optimize search queries where
 - data file not sorted
 - query not unique w.r.t Search Key.

Secondary Index (Non-Key)



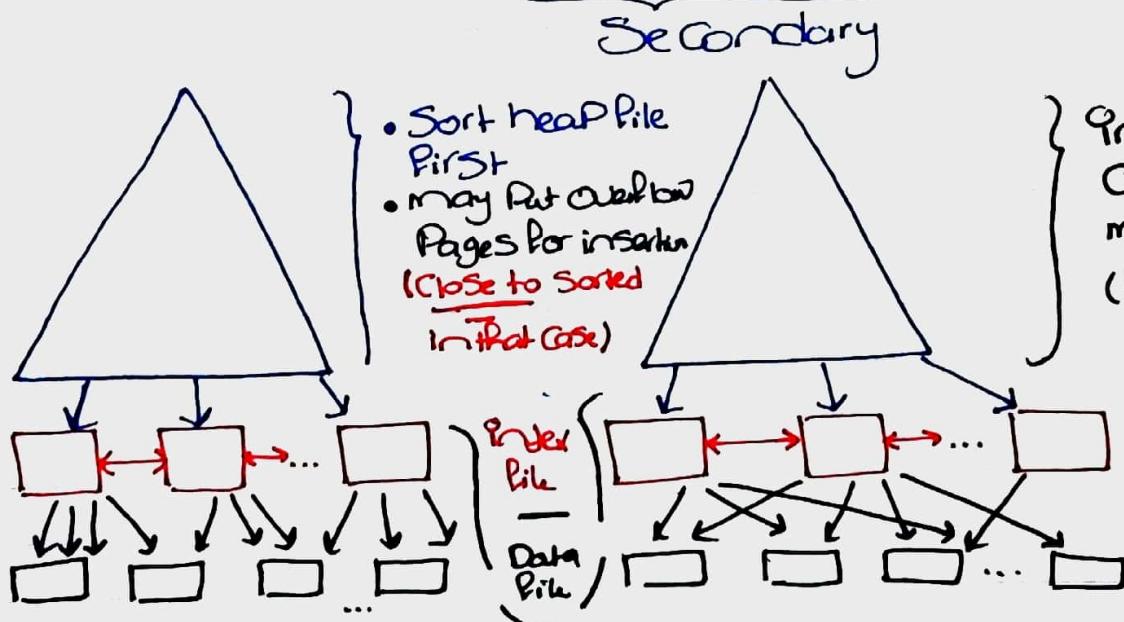


- Each unique index points to a block of pointers
→ The block has pointers to all data blocks where the index value occurs

• Summary

	Kind	Sparsity	Block Anchor
need sort { Primary key	b	SParse	✓
Clustering non-key	#values	SParse	✓ or ✗
Don't need sort { Secondary (Key)	r	Dense	X
Secondary (Non-Key)	#values	SParse (other variant → r)	SParse (other var → Dense)

• Clustered VS. Unclustered Index



- Some Order ↘ different order
- At most one clustered index for any file (one search key)

• Multi-level Index

→ Single-level index is an ordered file

→ Can hence always create a Primary index on the index itself

- Upper-level would be called Second-level index

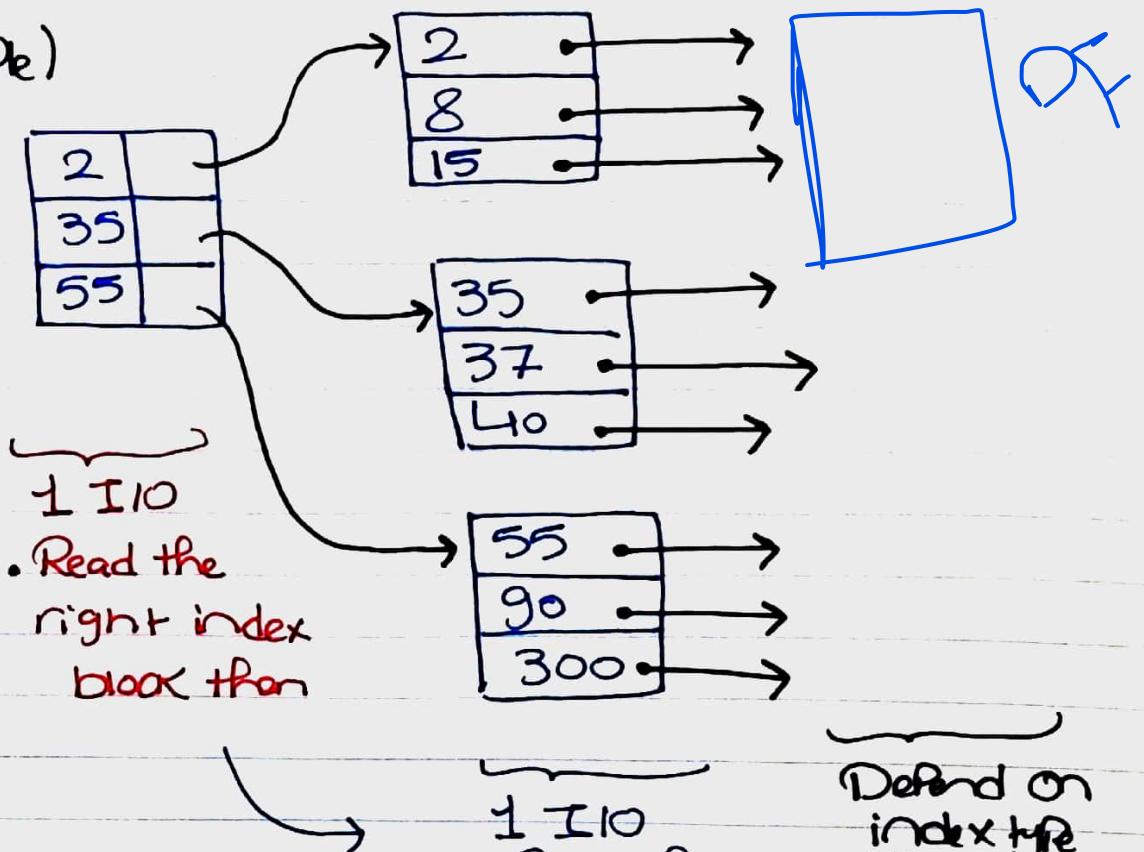
* Don't have to stop at 2nd level

→ Can keep going until top most level is one block

- Search becomes $O(\# \text{ levels})$ rather than blocks

- Type of 1st level doesn't matter (only should be more than one block)

Example)



ISAM

(indexed Seq. access method)

Problems with Index

→ Deletion
→ Use markers (Soft + delete)

Insertion

- Insert in right position in index
- May need to move records in data file, may affect anchor records (intuitively)

• Linked list or overflow records may help

→ Need a more dynamic solution.

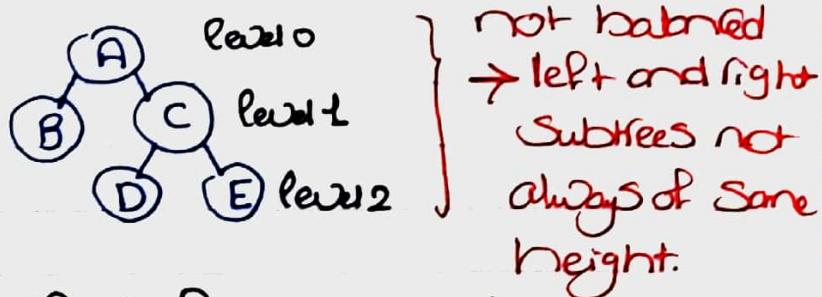
• A dynamic index isn't statically allocated and can grow/shrink as needed due to insertions/delets



Can either be tree based or hash based.

• Tree-based Dynamic Index (Multi-level)

→ Example of a tree



• ~~other tree satisfies that For any node~~

* Search Tree

→ Each node consists of P-1 Keys

→ To the left and to the right of each key there's a pointer to a node in the lower level (or nullptr).

• Hence, P Pointers for each node, each pointing to some sub-tree(node).

p pointers
p-1 keys.

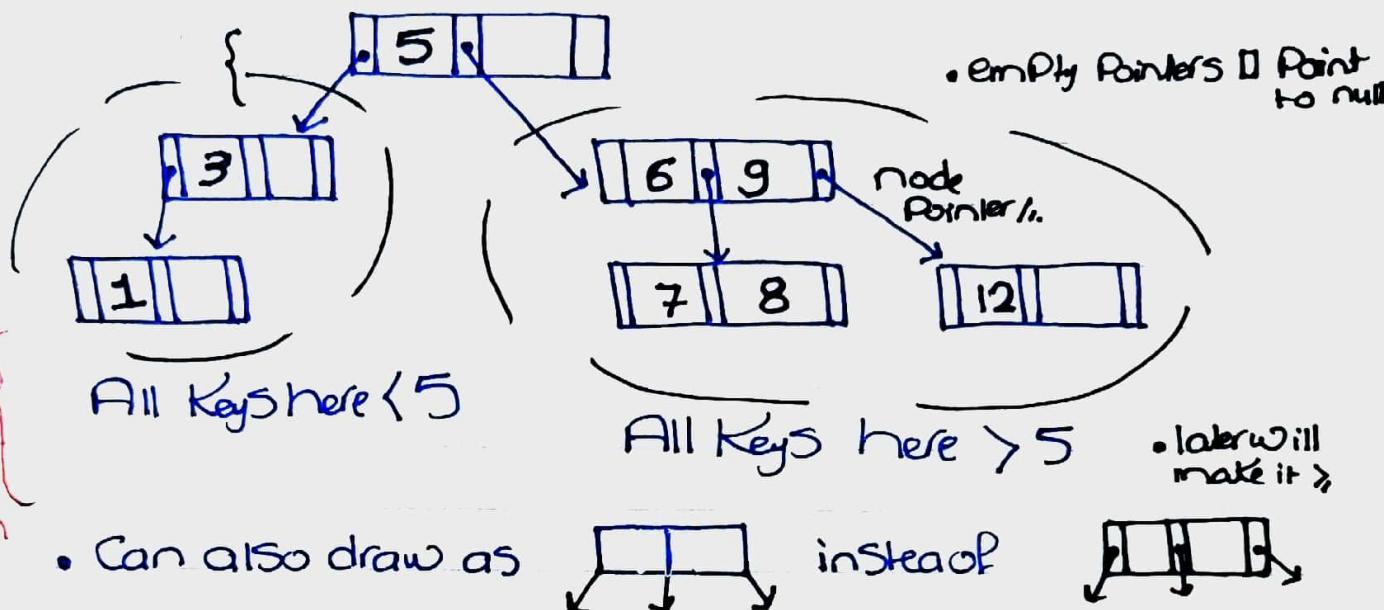
→ Condition:

For any Key in any node:

- It's greater than any key in the left subtree
- It's smaller than any key in the right subtree

Example ($P=3$)

(2 Keys & 3 Pts Per Node)



* B and B+ Trees

→ Variations of Search tree that make it Perfectly balanced (all leaves @ same level)

- Hence, Insertion and deletion are efficient (and nontrivial)
- also,

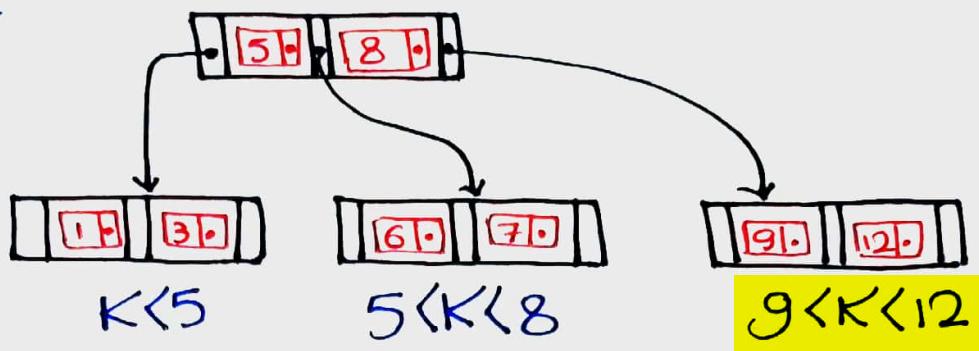
* Each Node Corresponds to a disk block (often has space for new entries)

Condition: Each node (disk block) must be kept between half-full and completely full

* Most widely used multi-level index.

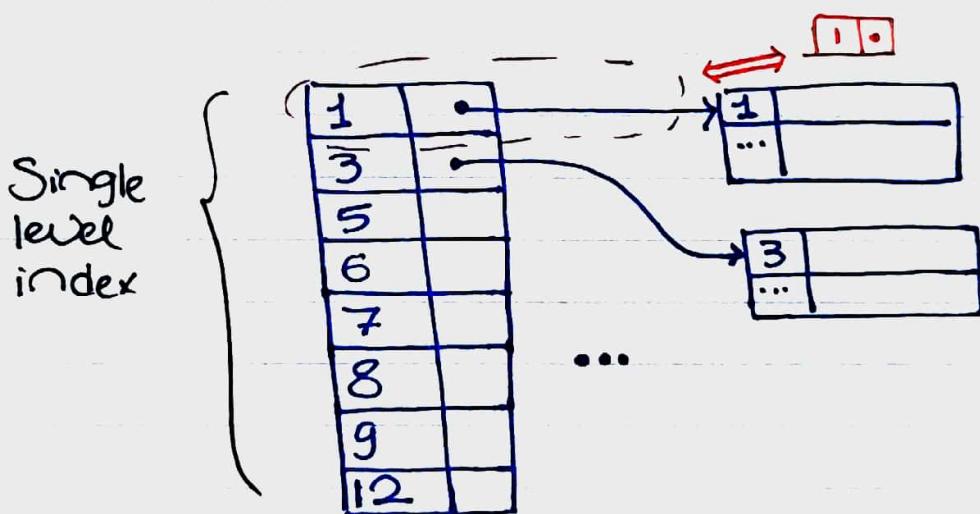
① degan

B-Tree



Don't
Confuse with
node pointer.

- It's a Search tree (we know already)
 - Each Key is also associated with data Pointer
- E.g., above tree is an Optimization over

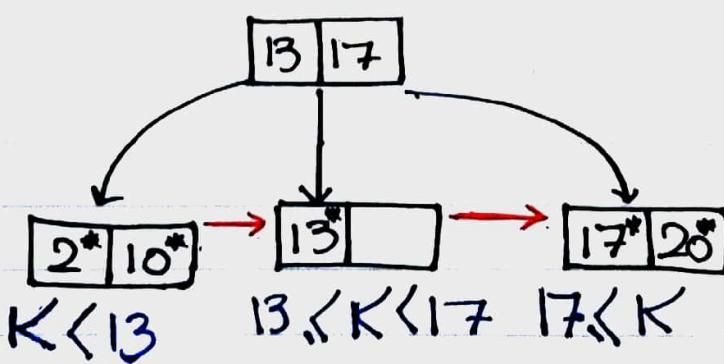


B+ Tree

- Restricts that data Pointers must only occur at leaves (only node pointers have that).
- Thus, more space to node Pointers; higher branching Factor, higher Capacity of Search Values (i.e., less depth or more search keys than B-trees)

• Has 2 node Pointers
• 13 vs. 13*

Has
data Pointer



Optimization over



- Notice that all leaves are also linked

- Convention for might be confused in slides

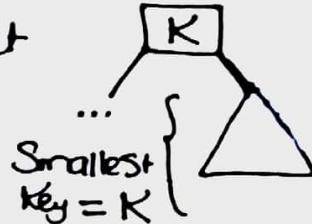
* In both B and B+ Trees

- Node not full \rightarrow efficient insertion $\boxed{AB} \rightarrow \boxed{\underline{ABC}}$
else may need to split, can affect other leafs.
- Node won't become less than half full after deletion \rightarrow efficient deletion $\boxed{AB} \rightarrow \boxed{\underline{A}}$
else may need to merge with other nodes. just

• We'll Focus on B+ Trees

Recap and Notation

- Each leaf node holds up to q entries (at least $\lceil q/2 \rceil$ to satisfy the condition)
- All leaves are at same depth
- Each internal node has P Pointers (& thus $P-1$ Keys)
(at least $\lceil P/2 \rceil$ Pointers, $\lceil \frac{P}{2} \rceil - 1$ Keys to satisfy the condition)
- It always holds that



also P or L

- else, something in the steps is possibly wrong

- Root is initially a leaf node (becomes internal node once it has children; generally)
- Only root can violate the minimum occupying condition ($\lceil \frac{q}{2} \rceil$ or $\lceil \frac{P}{2} \rceil - 1$) but $P \geq 2$ must hold.

Some Practical tips

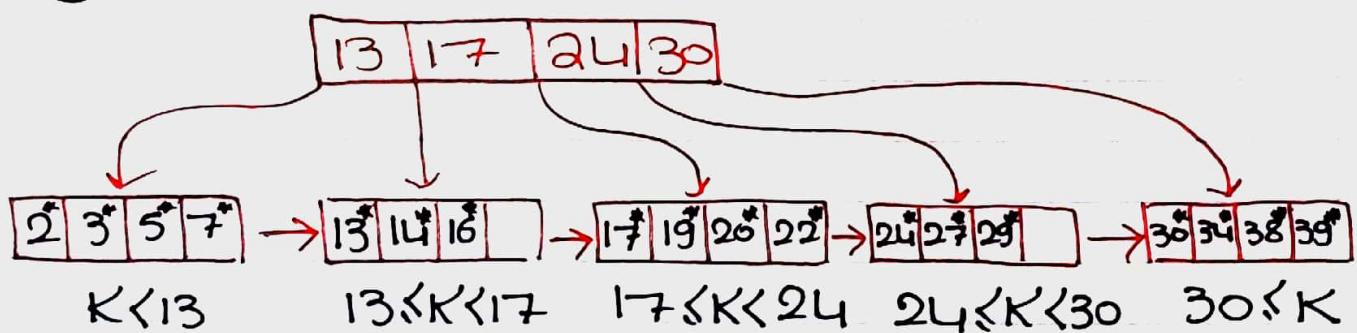
- Node = block of data on disk
- To speed up search can buffer 1st one or two levels on main memory

→ Most of the time, internal nodes aren't completely full

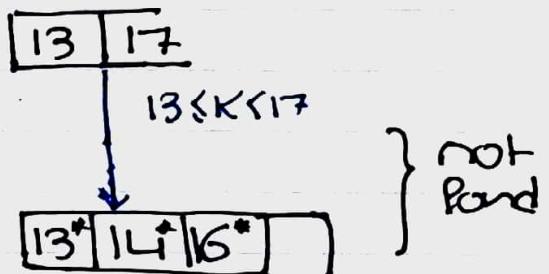
- Big waste for main memory (it's small)
- Disk is okay with that (big) (efficient in I/O)

* Will consider searching, insertion & deletion for B+ tree.

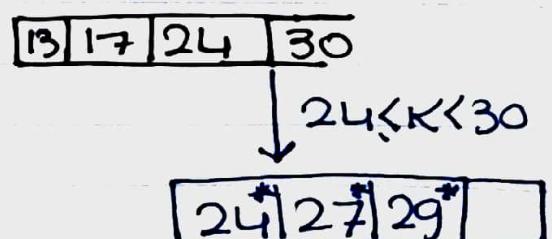
• Range Search & Equality Search

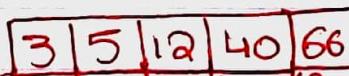


• Equality Search for 15



• Equality Search for 27



• Given , can you from one sight tell which pointer will help us find 53? how about 40?

• Range Search for values in [15, 28]

→ Get the leaf block where 15 is expected

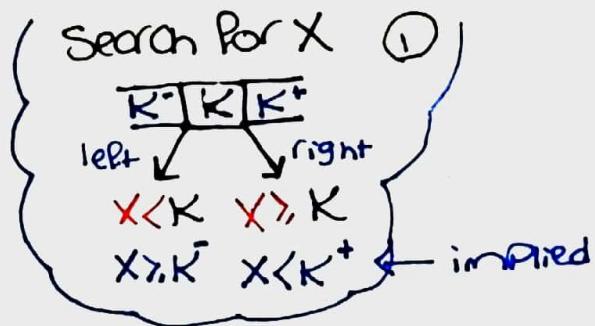
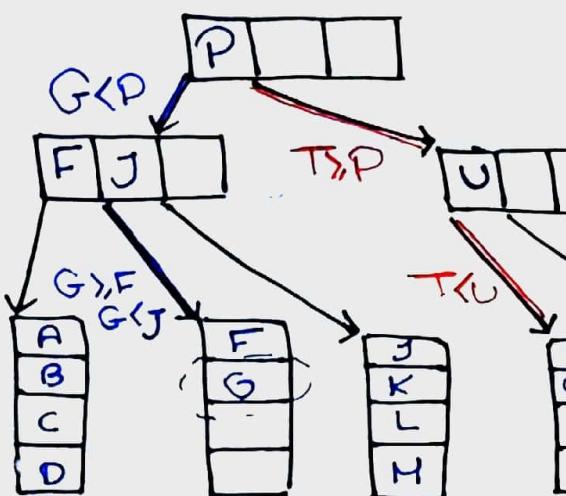


→ Loop over keys starting from that block and take all that satisfy condition

16 17 19 20 22 27

stop, it's sorted.

Generalizing:



- Search For G
- Search For T
→ not found

$$P = q = 4$$

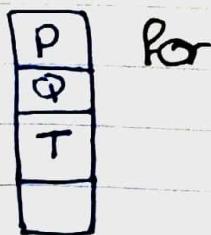
* Can also write an algo. out of this

Start { → Check if you should take extreme Ptrs of node
($X <$ 1st Key or $X >$ last Key)
From root & Keep taking Ptrs until ⇒ Once a leaf is found
→ Use linear or binary search to find when ① is satisfied
→ Use linear or binary search to find the value there (report not found if needed)

* B+ Tree Insertion

- Search for key to find the right block
- If not full, insert it
e.g., to insert T we have the 4th leaf

to validate { . Keep in mind leaves are sorted
(on both the block and leaves level)



- If Full, Create a new Sibling and redistribute then Fix Parent Pointers.

e.g., to insert E

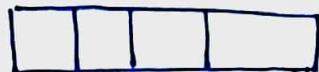
→ E < P

→ E < F

hence should go in



- Create a new Sibling (Splitting)



E

- Redistribute such that

q is even {

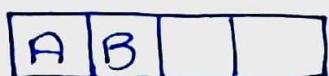
→ Middle Value is 1st element of new Sibling

q is odd {

→ half of the values go to 1st Sibling, 2nd half goes to 2nd Sibling

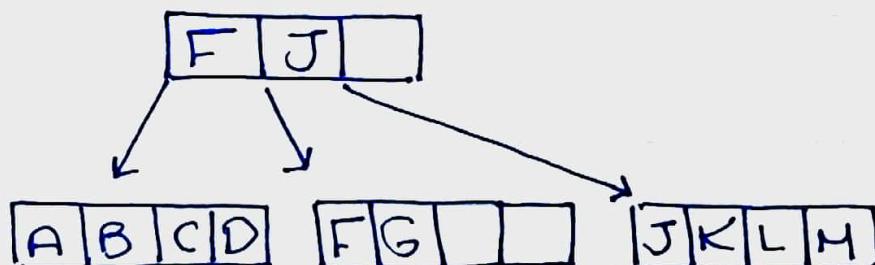
9
mark

Thus, in this case

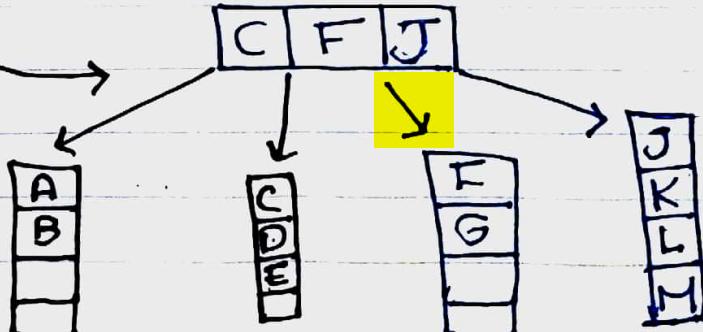


• C is middle value

- Now Fix Parent Pointers (adding 1st element of new Sibling 1st)



becomes

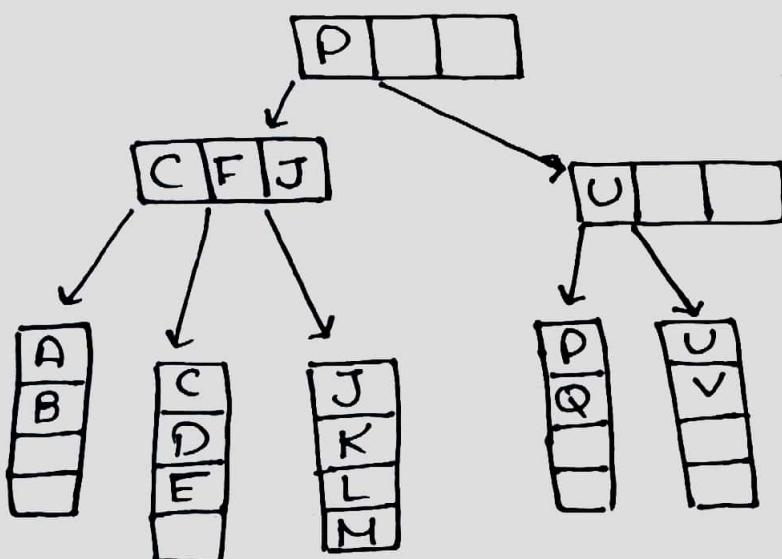


} Parent Fixed.

- What if the Parent doesn't have space for 1st element of 2nd Sibling?

- In that case we'd need to apply a different type of splitting on Parent
(Splitting leaf \neq Splitting internal node)
 - * what you saw last page
 - * Happens when leaf has no space
 - * may or may not happen after splitting leaf
 - \rightarrow In case Parent can't accommodate new child, will happen.

- To split an internal node
 - \rightarrow Create a new node (now have two)
 - \rightarrow move the $\lceil \frac{P}{2} \rceil$ th key to Parent
 - \rightarrow Redistribute the keys before and after it to 1st and 2nd node respectively



• To insert N

\rightarrow NKP

\rightarrow N, F
N, K, J

• goes to 3rd leaf

Parent C | F | J has no space for L (C | F | J | L) {

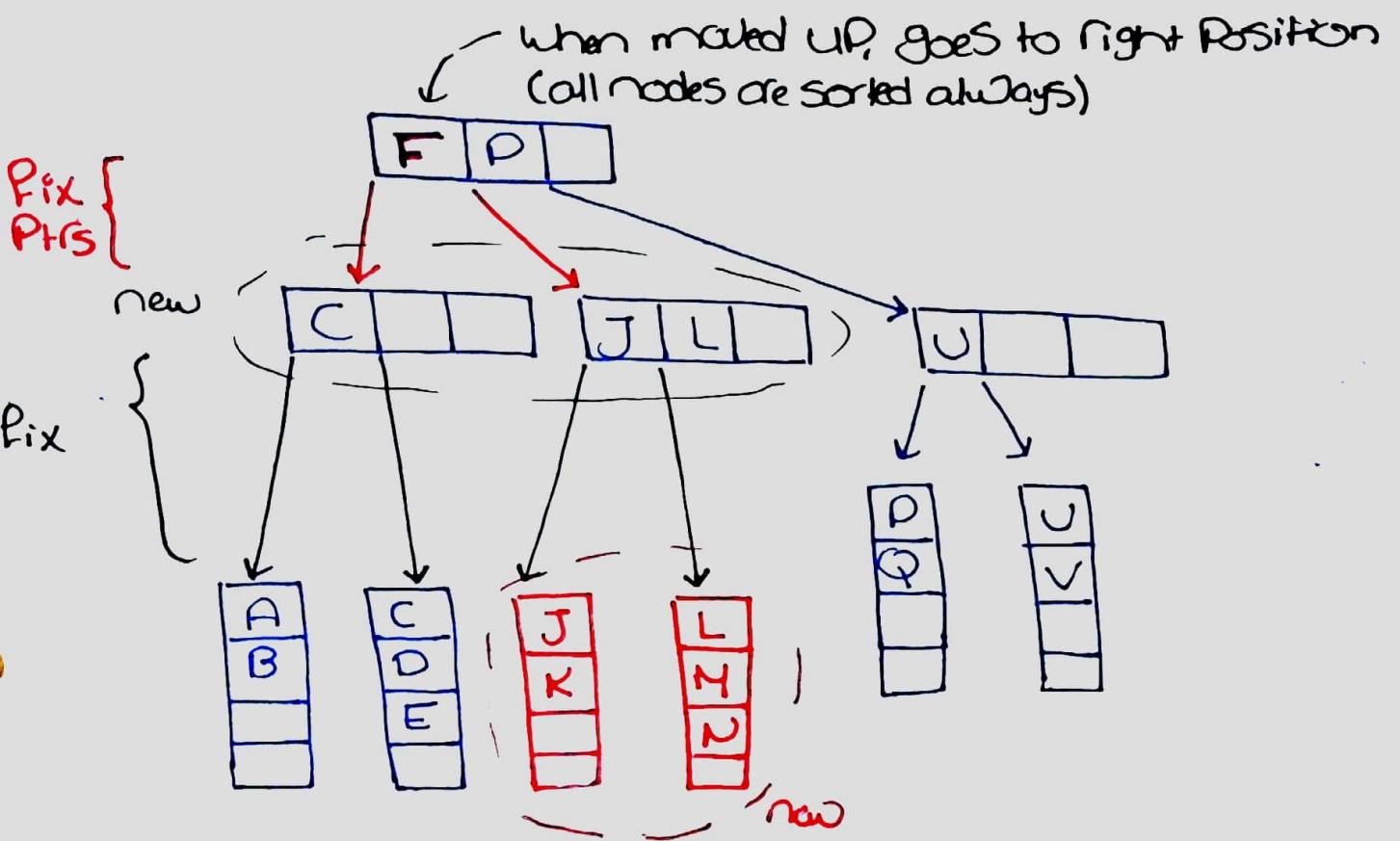
J	K
---	---

L	M	N
---	---	---

• Split Parent
 $P = 4$
 $\lceil \frac{P}{2} \rceil_{th} = F$ (goes up)



\rightarrow Now fix Parent and Children Pointers of internal node

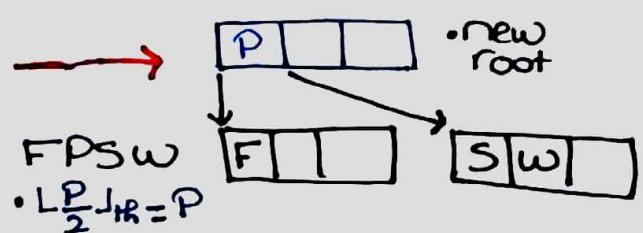


→ If the upper level's Parent (was $\boxed{P \quad \quad}$) can't accommodate the new entry that should move up, it will be split in the same way. ☺

- until we reach the root

→ When the root is split we make the key that should move up the only element in the new root.

e.g. $\boxed{F \quad P \quad S}$

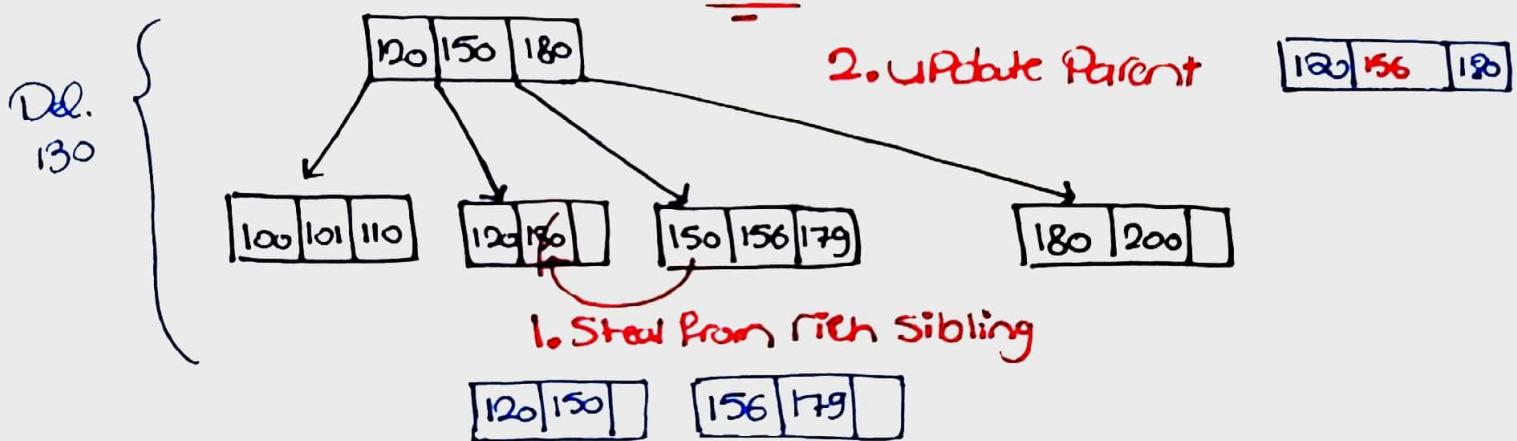


- Notice that depth must increase by 1 whenever root is split

B+ Tree Deletion

• In B+ Tree
 → Cousins are also Siblings

- If the leaf node will still have $\lceil \frac{q}{2} \rceil$ or more
 - Then everything is okay, just delete.
- Otherwise if after deleting you have less than $\lceil \frac{q}{2} \rceil$ (i.e., $\lceil \frac{q}{2} \rceil - 1$)
 - See if Stealing max element from left sibling or min element from right sibling doesn't put them at risk ($\lfloor \frac{q}{2} \rfloor$).
 - Update Parent Key to match new 1st element



- If both siblings are poor (Stealing violates their $\lceil \frac{q}{2} \rceil$) then merge with one of them
 - guaranteed to fit in a one node
 - Update Parent by removing key

