

### **B+-Tree Index Files**

B<sup>+</sup>-tree indices are an alternative to indexed-sequential files.

- Disadvantage of indexed-sequential files
  - performance degrades as file grows, since many overflow blocks get created.
  - Periodic reorganization of entire file is required.
- Advantage of B<sup>+</sup>-tree index files:
  - automatically reorganizes itself with small, local, changes, in the face of insertions and deletions.
  - Reorganization of entire file is not required to maintain performance.
- (Minor) disadvantage of B+-trees:
  - extra insertion and deletion overhead, space overhead.
- Advantages of B<sup>+</sup>-trees outweigh disadvantages
  - B+-trees are used extensively in DBMSs



#### **Example of B+-Tree**





# **B+-Tree Index Files (Cont.)**

A B<sup>+</sup>-tree is a rooted tree satisfying the following properties:

- All paths from root to leaf are of the same length
- Each node that is not a root or a leaf has between  $\lceil n/2 \rceil$  and *n* children.
- A leaf node has between [(n-1)/2] and n-1 values
- Special cases:
  - If the root is not a leaf, it has at least 2 children.
  - If the root is a leaf (that is, there are no other nodes in the tree), it can have between 0 and (n-1) values.



### **B+-Tree Node Structure**

#### Typical node

$P_1$ $P_1$	$K_1 \qquad P_2$	•••	<i>P</i> <sub><i>n</i>-1</sub>	<i>K</i> <sub><i>n</i>-1</sub>	$P_n$
-------------	------------------	-----	--------------------------------	--------------------------------	-------

- K<sub>i</sub> are the search-key values
- P<sub>i</sub> are pointers to children (for non-leaf nodes) or pointers to records or buckets of records (for leaf nodes).
- The search-keys in a node are ordered

$$K_1 < K_2 < K_3 < \ldots < K_{n-1}$$

(Initially assume no duplicate keys, address duplicates later)



### Leaf Nodes in B<sup>+</sup>-Trees

Properties of a leaf node:

- For i = 1, 2, ..., n-1, pointer  $P_i$  points to a file record with search-key value  $K_i$ ,
- If  $L_i$ ,  $L_j$  are leaf nodes and i < j,  $L_i$ 's search-key values are less than or equal to  $L_i$ 's search-key values
- P<sub>n</sub> points to next leaf node in search-key order leaf node

Brandt	Brandt Califieri Crick Pointer to next leaf node						
			10101	Srinivasan	Comp. Sci.	65000	
			12121	Wu	Finance	90000	
			15151	Mozart	Music	40000	
			22222	Einstein	Physics	95000	
			32343	El Said	History	80000	
			33456	Gold	Physics	87000	
			45565	Katz	Comp. Sci.	75000	
			58583	Califieri	History	60000	
			76543	Singh	Finance	80000	
	l		76766	Crick	Biology	72000	
		)	83821	Brandt	Comp. Sci.	92000	
			98345	Kim	Elec. Eng.	80000	



## Non-Leaf Nodes in B<sup>+</sup>-Trees

- Non leaf nodes form a multi-level sparse index on the leaf nodes. For a non-leaf node with *m* pointers:
  - All the search-keys in the subtree to which  $P_1$  points are less than  $K_1$
  - For  $2 \le i \le n 1$ , all the search-keys in the subtree to which  $P_i$  points have values  $>= K_{i-1}$  and  $< K_i$
  - All the search-keys in the subtree to which  $P_n$  points have values greater than or equal to  $K_{n-1}$



B<sup>+</sup>-tree for *instructor* file (n = 6)

- Leaf nodes must have between 3 and 5 values  $(\lceil (n-1)/2 \rceil$  and n-1, with n = 6).
- Non-leaf nodes other than root must have between 3 and 6 children ([(n/2] and n with n =6).
- Root must have at least 2 children.



#### **Observations about B+-trees**

- Since the inter-node connections are done by pointers, "logically" close blocks need not be "physically" close.
- The non-leaf levels of the B<sup>+</sup>-tree form a hierarchy of sparse indices.
- The B+-tree contains a relatively small number of levels
  - Level below root has at least 2\* [n/2] values
  - Next level has at least 2\* [n/2] \* [n/2] values
  - .. etc, level at depth D has at least 2\* [n/2]<sup>D</sup>
  - If there are K search-key values in the file, the tree height is no more than [ log<sub>[n/2]</sub>(K)]
  - thus searches can be conducted efficiently.
- Insertions and deletions to the main file can be handled efficiently, as the index can be restructured in logarithmic time (as we shall briefly see, but more can be found in the book).



## **Queries on B+-Trees**

Find record with search-key value V.

- 1. C=root
- 2. While C is not a leaf node {
  - 1. Let *i* be least value s.t.  $V \le K_i$ .
  - 2. If no such exists, set C = last non-null pointer in C

B. Else { if 
$$(V = K_i)$$
 Set  $C = P_{i+1}$  else set  $C = P_i$ 

- 3. Let *i* be least value s.t.  $K_i = V$
- 4. If there is such a value *i*, follow pointer  $P_i$  to the desired record.
- 5. Else no record with search-key value k exists.



Database System Concepts - 6<sup>th</sup> Edition

©Silberschatz, Korth and Sudarshan



# **Handling Duplicates**

- With duplicate search keys
  - In both leaf and internal nodes,
    - we cannot guarantee that  $K_1 < K_2 < K_3 < \ldots < K_{n-1}$
    - but can guarantee  $K_1 \le K_2 \le K_3 \le \ldots \le K_{n-1}$
  - Search-keys in the subtree to which P<sub>i</sub> points
    - are  $\leq K_{i,}$ , but not necessarily  $< K_{i,}$
    - To see why, suppose same search key value V is present in two leaf node L<sub>i</sub> and L<sub>i+1</sub>. Then in parent node K<sub>i</sub> must be equal to V



# Handling Duplicates

- We modify find procedure as follows
  - traverse  $P_i$  even if  $V = K_i$
  - As soon as we reach a leaf node C check if C has only search key values less than V
    - if so set C = right sibling of C before checking whether C contains V
- Procedure printAll
  - uses modified find procedure to find first occurrence of V
  - Traverse through consecutive leaves to find all occurrences of V

<sup>\*\*</sup> Errata note: modified find procedure missing in first printing of 6<sup>th</sup> edition



# **Queries on B+-Trees (Cont.)**

- If there are K search-key values in the file, the height of the tree is no more than  $\lceil \log_{\lfloor n/2 \rceil}(K) \rceil$ .
- A node is generally the same size as a disk block, typically 4 kilobytes
  - and *n* is typically around 100 (40 bytes per index entry).
- With 1 million search key values and n = 100
  - at most log<sub>50</sub>(1,000,000) = 4 nodes are accessed in a lookup, i.e. at most 4 accesses to disk blocks are needed
- Contrast this with a balanced binary tree with 1 million search key values around 20 nodes are accessed in a lookup
  - above difference is significant since every node access may need a disk I/O, costing around 20 milliseconds



## **Updates on B<sup>+</sup>-Trees: Insertion**

- 1. Find the leaf node in which the search-key value would appear
- 2. If the search-key value is already present in the leaf node
  - 1. Add record to the file
  - 2. If necessary add a pointer to the bucket.
- 3. If the search-key value is not present, then
  - 1. add the record to the main file (and create a bucket if necessary)
  - 2. If there is room in the leaf node, insert (key-value, pointer) pair in the leaf node
  - 3. Otherwise, split the node (along with the new (key-value, pointer) entry) as discussed in the next slide.



# Updates on B<sup>+</sup>-Trees: Insertion (Cont.)

- Splitting a leaf node:
  - take the n (search-key value, pointer) pairs (including the one being inserted) in sorted order. Place the first [n/2] in the original node, and the rest in a new node.
  - let the new node be p, and let k be the least key value in p. Insert
     (k,p) in the parent of the node being split.
  - If the parent is full, split it and **propagate** the split further up.
- Splitting of nodes proceeds upwards till a node that is not full is found.
  - In the worst case the root node may be split increasing the height of the tree by 1.



Result of splitting node containing Brandt, Califieri and Crick on inserting Adams Next step: insert entry with (Califieri,pointer-to-new-node) into parent



#### **B+-Tree Insertion**



B<sup>+</sup>-Tree before and after insertion of "Adams"

Einstein El Said

┢

Srinivasan

Gold Katz Kim +

Mozart Singh

**Database System Concepts - 6th Edition** 

Adams Brandt

Califieri Einstein Gold

Califieri Crick

Srinivasan

Wu



#### **B+-Tree Insertion**



B<sup>+</sup>-Tree before and after insertion of "Lamport"



# Insertion in B+-Trees (Cont.)

- Splitting a non-leaf node: when inserting (k,p) into an already full internal node N
  - Copy N to an in-memory area M with space for n+1 pointers and n keys
  - Insert (k,p) into M
  - Copy  $P_1, K_1, ..., K_{[(n+1)/2]-1}, P_{[(n+1)/2]}$  from M back into node N
  - Copy  $P_{[(n+1)/2]+1}$ ,  $K_{[(n+1)/2]+1}$ ,...,  $K_n$ ,  $P_{n+1}$  from M into newly allocated node N'
  - Insert ( $K_{[(n+1)/2]}$ ,N') into parent N
- Read pseudocode in book ERRATA !





## **Examples of B<sup>+</sup>-Tree Deletion**



#### Before and after deleting "Srinivasan"



Deleting "Srinivasan" causes merging of under-full leaves





Deletion of "Singh" and "Wu" from result of previous example

- Leaf containing Singh and Wu became underfull, and borrowed a value Kim from its left sibling
- Search-key value in the parent changes as a result



Before and after deletion of "Gold" from earlier example

- Node with Gold and Katz became underfull, and was merged with its sibling
- Parent node becomes underfull, and is merged with its sibling
  - Value separating two nodes (at the parent) is pulled down when merging
  - Root node then has only one child, and is deleted

**Database System Concepts - 6th Edition** 



## **Updates on B+-Trees: Deletion**

- Find the record to be deleted, and remove it from the main file and from the bucket (if present)
- Remove (search-key value, pointer) from the leaf node if there is no bucket or if the bucket has become empty
- If the node has too few entries due to the removal, but the entries in the node and a sibling do not fit into a single node, then redistribute pointers:
  - Redistribute the pointers between the node and a sibling such that both have more than the minimum number of entries.
  - Update the corresponding search-key value in the parent of the node.



## **Updates on B<sup>+</sup>-Trees: Deletion**

- Otherwise, if the node has too few entries due to the removal, and the entries in the node and a sibling fit into a single node, then *merge siblings:* 
  - Insert all the search-key values in the two nodes into a single node (the one on the left), and delete the other node.
  - Delete the pair  $(K_{i-1}, P_i)$ , where  $P_i$  is the pointer to the deleted node, from its parent, recursively using the above procedure.
- The node deletions may cascade upwards till a node which has  $\lceil n/2 \rceil$  or more pointers is found.
- If the root node has only one pointer after deletion, it is deleted and the sole child becomes the root.



# **Non-Unique Search Keys**

- Alternatives to scheme described earlier
  - Buckets on separate block (bad idea)
  - List of tuple pointers with each key
    - Extra code to handle long lists
    - Deletion of a tuple can be expensive if there are many duplicates on search key (why?)
    - Low space overhead, no extra cost for queries
  - Make search key unique by adding a record-identifier
    - Extra storage overhead for keys
    - Simpler code for insertion/deletion
    - Widely used (e.g. Oracle always assumes this by adding row-id)



## **B+-Tree File Organization**

- Index file degradation problem is solved by using B+-Tree indices.
- Data file degradation problem is solved by using B<sup>+</sup>-Tree File Organization.
- The leaf nodes in a B<sup>+</sup>-tree file organization store records, instead of pointers.
- Leaf nodes are still required to be half full
  - Since records are larger than pointers, the maximum number of records that can be stored in a leaf node is less than the number of pointers in a nonleaf node.
- Insertion and deletion are handled in the same way as insertion and deletion of entries in a B<sup>+</sup>-tree index.
- May be used to store big objects (those that do not fit into a single record)



## **B+-Tree File Organization (Cont.)**



Example of B<sup>+</sup>-tree File Organization

- Good space utilization important since records use more space than pointers.
- To improve space utilization, involve more sibling nodes in redistribution during splits and merges
  - Involving 2 siblings in redistribution (to avoid split / merge where possible) results in each node having at least |2n/3| entries



## **Other Issues in Indexing**

#### Record relocation and secondary indices

- If a record moves, all secondary indices that store record pointers have to be updated
- Node splits in B<sup>+</sup>-tree file organizations become very expensive
- Solution: use primary-index search key instead of record pointer in secondary index
  - > Extra traversal of primary index to locate record
    - Higher cost for queries, but node splits are cheap
  - Add record-id if primary-index search key is non-unique



# **Indexing Strings**

- Variable length strings as keys
  - Variable fanout
  - Use space utilization as criterion for splitting, not number of pointers
- Prefix compression
  - Key values at internal nodes can be prefixes of full key
    - Keep enough characters to distinguish entries in the subtrees separated by the key value
      - E.g. "Silas" and "Silberschatz" can be separated by "Silb"
  - Keys in leaf node can be compressed by sharing common prefixes



# **Bulk Loading and Bottom-Up Build**

- Inserting entries one-at-a-time into a B<sup>+</sup>-tree requires  $\geq$  1 IO per entry
  - assuming leaf level does not fit in memory
  - can be very inefficient for loading a large number of entries at a time (bulk loading)
- Efficient alternative 1:
  - sort entries first (using efficient external-memory sort algorithms discussed later in Section 12.4)
  - insert in sorted order
    - insertion will go to existing page (or cause a split)
    - much improved IO performance, but most leaf nodes half full
- Efficient alternative 2: **Bottom-up B+-tree construction** 
  - As before sort entries
  - And then create tree layer-by-layer, starting with leaf level
    - details as an exercise
  - Implemented as part of bulk-load utility by most database systems



### **B-Tree Index Files**

- Similar to B+-tree, but B-tree allows search-key values to appear only once; eliminates redundant storage of search keys.
- Search keys in nonleaf nodes appear nowhere else in the Btree; an additional pointer field for each search key in a nonleaf node must be included.
- Generalized B-tree leaf node



(b)

Nonleaf node – pointers Bi are the bucket or file record pointers.

# **B-Tree Index File Example**



#### B-tree (above) and B+-tree (below) on same data



**Database System Concepts - 6th Edition** 



# **B-Tree Index Files (Cont.)**

- Advantages of B-Tree indices:
  - May use less tree nodes than a corresponding B+-Tree.
  - Sometimes possible to find search-key value before reaching leaf node.
- Disadvantages of B-Tree indices:
  - Only small fraction of all search-key values are found early
  - Non-leaf nodes are larger, so fan-out is reduced. Thus, B-Trees typically have greater depth than corresponding B<sup>+</sup>-Tree
  - Insertion and deletion more complicated than in B+-Trees
  - Implementation is harder than B+-Trees.
- Typically, advantages of B-Trees do not out weigh disadvantages.



# **Multiple-Key Access**

- Use multiple indices for certain types of queries.
- Example:

select ID

from instructor

where dept\_name = "Finance" and salary = 80000

- Possible strategies for processing query using indices on single attributes:
  - 1. Use index on *dept\_name* to find instructors with department name Finance; test *salary = 80000*
  - 2. Use index on salary to find instructors with a salary of \$80000; test dept\_name = "Finance".
  - 3. Use *dept\_name* index to find pointers to all records pertaining to the "Finance" department. Similarly use index on *salary*. Take intersection of both sets of pointers obtained.



# **Indices on Multiple Keys**

- Alternatively, use composite search keys that are search keys containing more than one attribute
  - E.g. (*dept\_name, salary*)
- Lexicographic ordering:  $(a_1, a_2) < (b_1, b_2)$  if either
  - a<sub>1</sub> < b<sub>1</sub>, or
  - $a_1 = b_1$  and  $a_2 < b_2$



## **Indices on Multiple Attributes**

Suppose we have an index on combined search-key (*dept\_name, salary*).

With the **where** clause

where dept\_name = "Finance" and salary = 80000
the index on (dept\_name, salary) can be used to fetch only records
that satisfy both conditions.

 Using separate indices in less efficient — we may fetch many records (or pointers) that satisfy only one of the conditions.

- Can also efficiently handle where dept\_name = "Finance" and salary < 80000</p>
- But cannot efficiently handle
  - where *dept\_name* < "Finance" and *balance* = 80000
  - May fetch many records that satisfy the first but not the second condition



#### **Other Features**

#### Covering indices

- Add extra attributes to index so (some) queries can avoid fetching the actual records
  - Particularly useful for secondary indices
    - Why?
- Can store extra attributes only at leaf