

# **Chapter 14: Indexing**

**Sistemas de Bases de Dados 2019/20**

**Capítulo refere-se a: Database System Concepts, 7th Ed**

# Outline

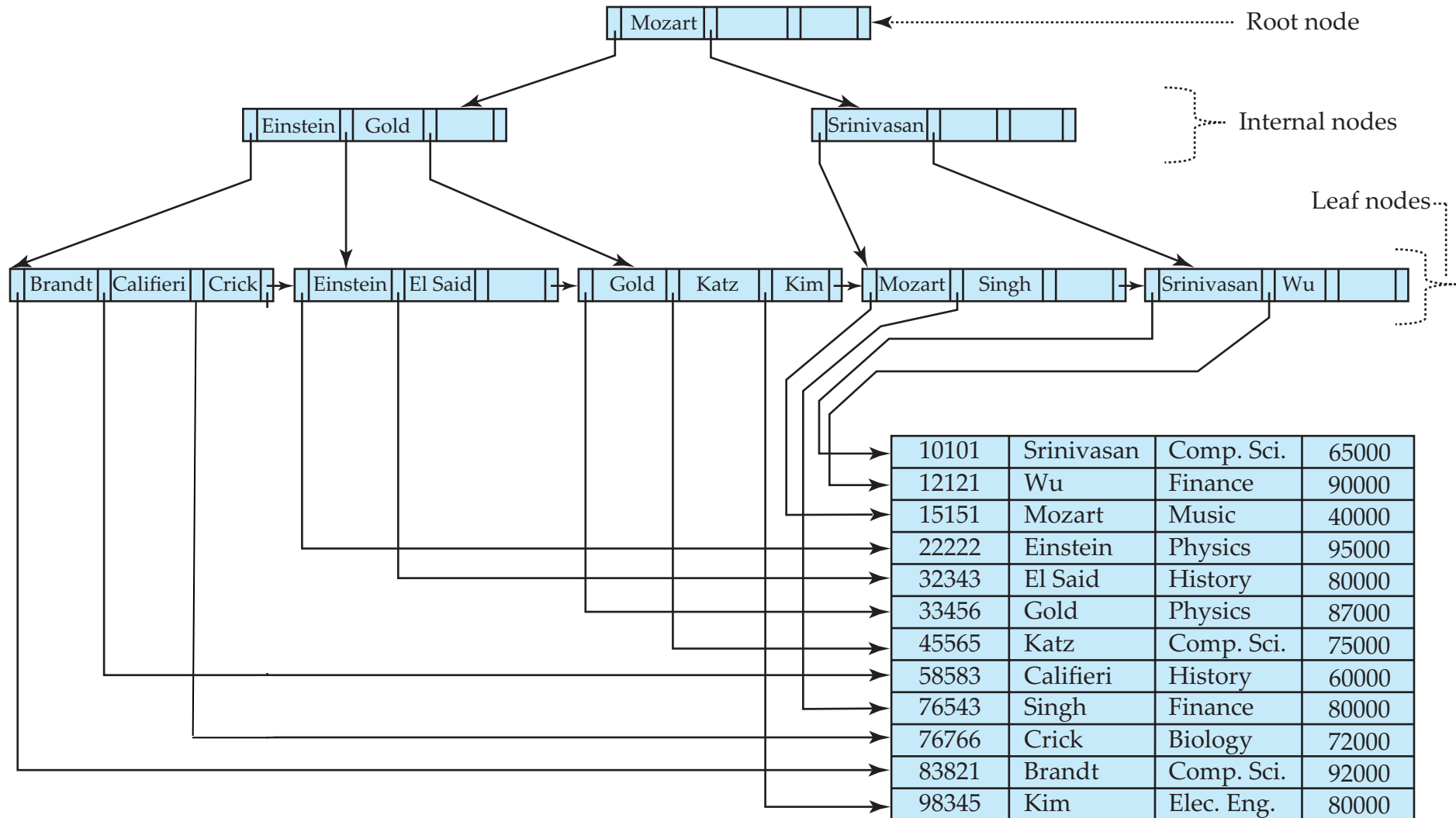
- Basic Concepts
- Ordered Indices
- B+-Tree Index Files
- B-Tree Index Files
- Hashing
- Write-optimized indices
- Spatio-Temporal Indexing

# B+-Tree Index Files

B+-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+-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+-trees outweigh disadvantages
  - B+-trees are used extensively

# Example of B+-Tree



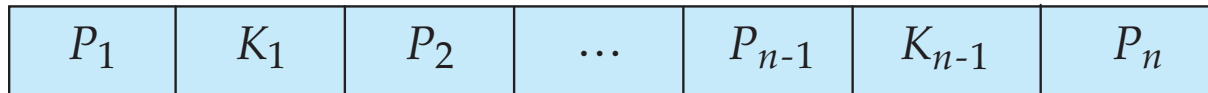
# B<sup>+</sup>-Tree Index Files (Cont.)

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  $\lceil (n-1)/2 \rceil$  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<sup>+</sup>-Tree Node Structure

- Typical node



- $K_i$  are the search-key values
  - $P_i$  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

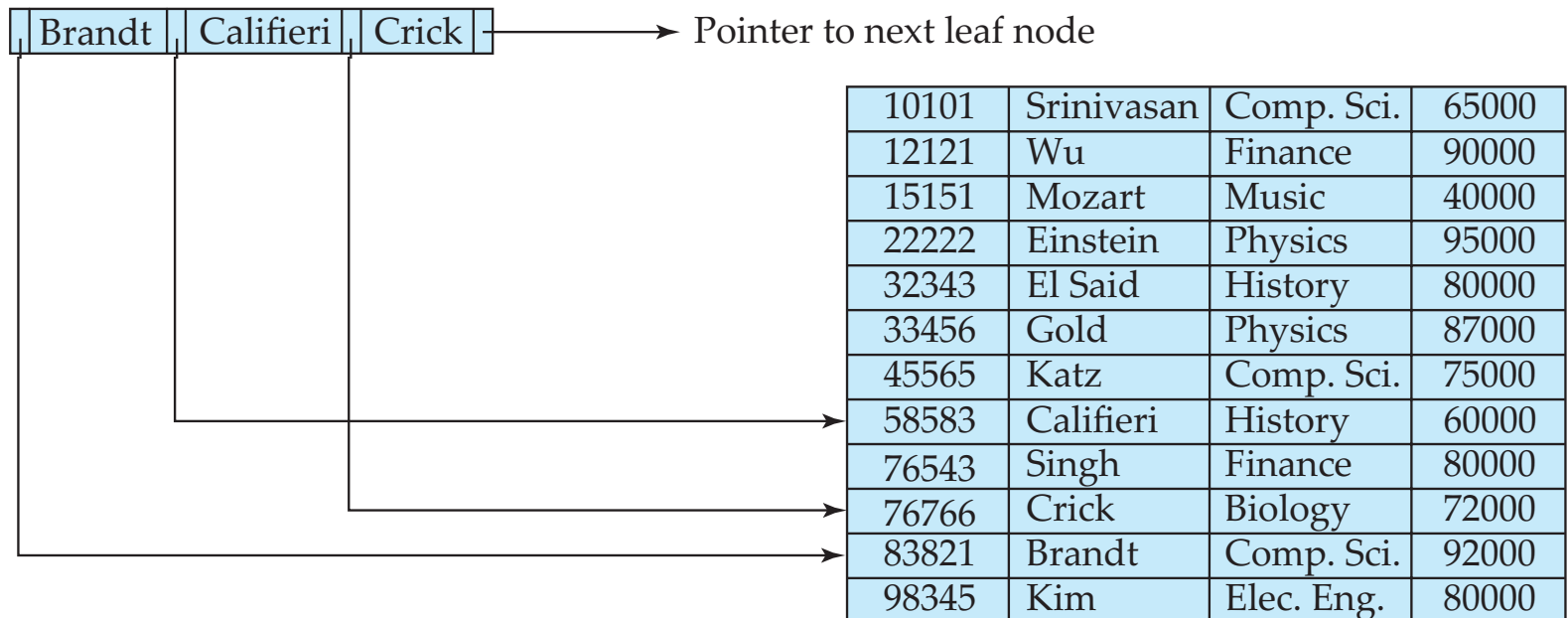
$$K_1 < K_2 < K_3 < \dots < K_{n-1}$$

(Initially assume no duplicate keys, address duplicates later)

# Leaf Nodes in B+-Trees

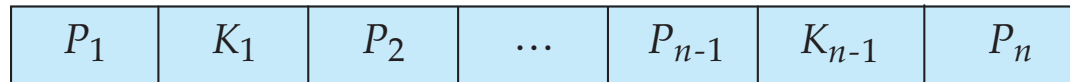
Properties of a leaf node:

- For  $i = 1, 2, \dots, n-1$ , 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  $L_j$ 's search-key values
- $P_n$  points to next leaf node in search-key order  
leaf node



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

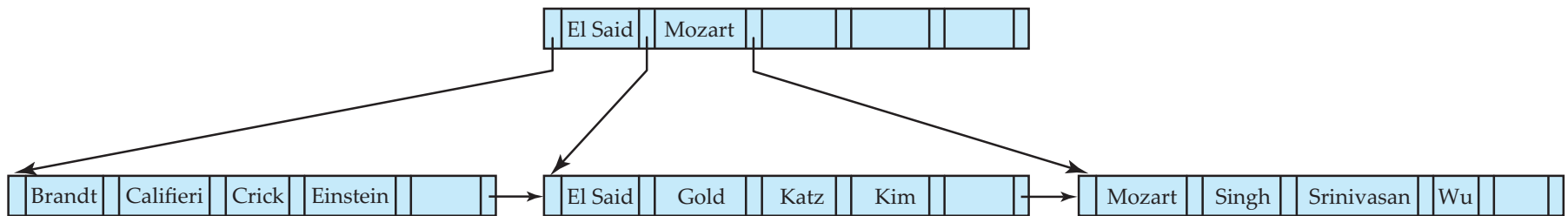
- Non leaf nodes form a multi-level sparse index on the leaf nodes. For a non-leaf node with  $m$  pointers:
  - All the search-keys in the subtree to which  $P_1$  points are less than  $K_1$
  - For  $2 \leq i \leq n - 1$ , all the search-keys in the subtree to which  $P_i$  points have values greater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





# Example of B<sup>+</sup>-tree

- 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 ( $\lceil n/2 \rceil$  and  $n$  with  $n = 6$ ).
- Root must have at least 2 children.

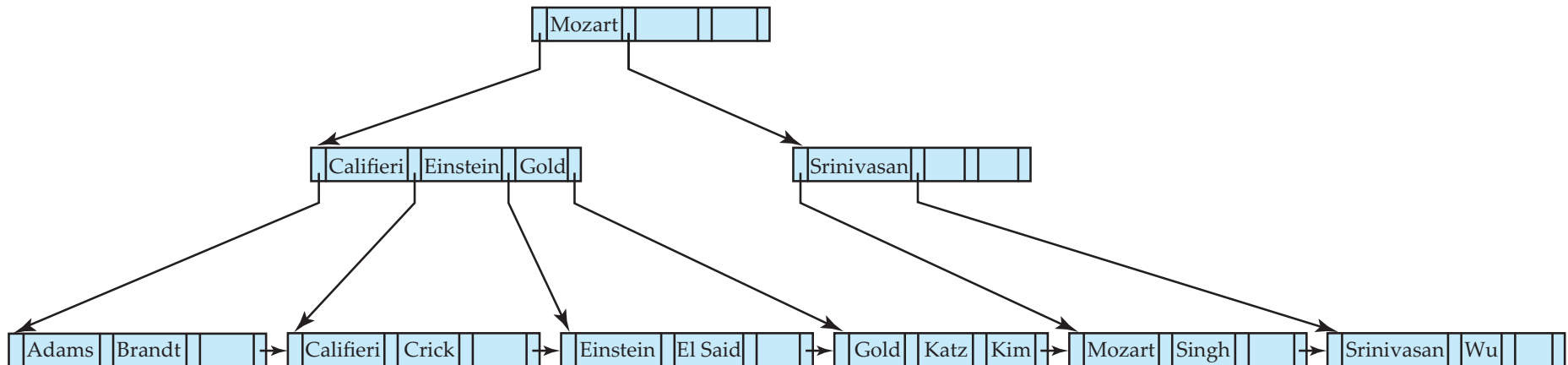
# Observations about B<sup>+</sup>-trees

- 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<sup>+</sup>-tree contains a relatively small number of levels
  - Level below root has at least  $2 * \lceil n/2 \rceil$  values
  - Next level has at least  $2 * \lceil n/2 \rceil * \lceil n/2 \rceil$  values
  - .. etc.
- If there are  $K$  search-key values in the file, the tree height is no more than  $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$
- thus searches (and also insertions and deletion, as we shall see) can be efficient.
- Most space occupied by a B<sup>+</sup>-tree is in the leaves:
  - E.g. a B<sup>+</sup>-tree with  $n=100$ , 1M values, each one (+pointer) with 40 bytes
    - Each node is at most 4K (guess why I chose these numbers :-)
    - The maximum height of the tree is  $\lceil \log_{50}(1000000) \rceil = 4$
    - Leaves occupy between 1Mb and 2Mb
    - Level 3 occupies at most  $2\text{Mb}/50 = 40\text{Kb}$
    - The whole intermediate levels have 48Kb

# Queries on B<sup>+</sup>-Trees

**function** *find*(*v*)

1. *C* = *root*
2. **while** (*C* is not a leaf node)
  1. Let *i* be least number s.t.  $V \leq K_i$ .
  2. **if** there is no such number *i* **then**
  3.     Set *C* = *last non-null pointer in C*
  4. **else if** ( $v = C.K_i$ ) Set *C* = *P<sub>i+1</sub>*
  5. **else set** *C* = *C.P<sub>i</sub>*
3. **if** for some *i*,  $K_i = V$  **then** return *C.P<sub>i</sub>*
4. **else** return null /\* no record with search-key value *v* exists. \*/



# Queries on B<sup>+</sup>Trees (Cont.)

- If there are  $K$  search-key values in the file, the height of the tree is no more than  $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$ .
- 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  $\log_{50}(1,000,000) = 4$  nodes are accessed in a lookup traversal from root to leaf.
- 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

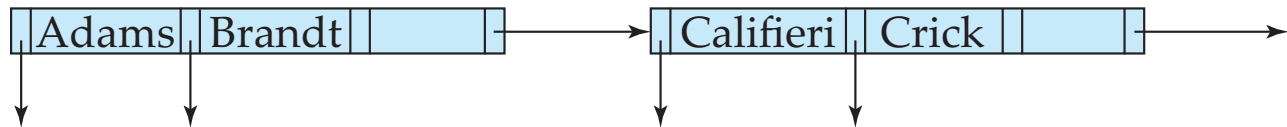
Assume record already added to the file. Let

- $pr$  be pointer to the record, and let
- $v$  be the search key value of the record

1. Find the leaf node in which the search-key value would appear
  1. If there is room in the leaf node, insert  $(v, pr)$  pair in the leaf node
  2. Otherwise, split the node (along with the new  $(v, pr)$  entry) as discussed in the next slide, and propagate updates to parent nodes.

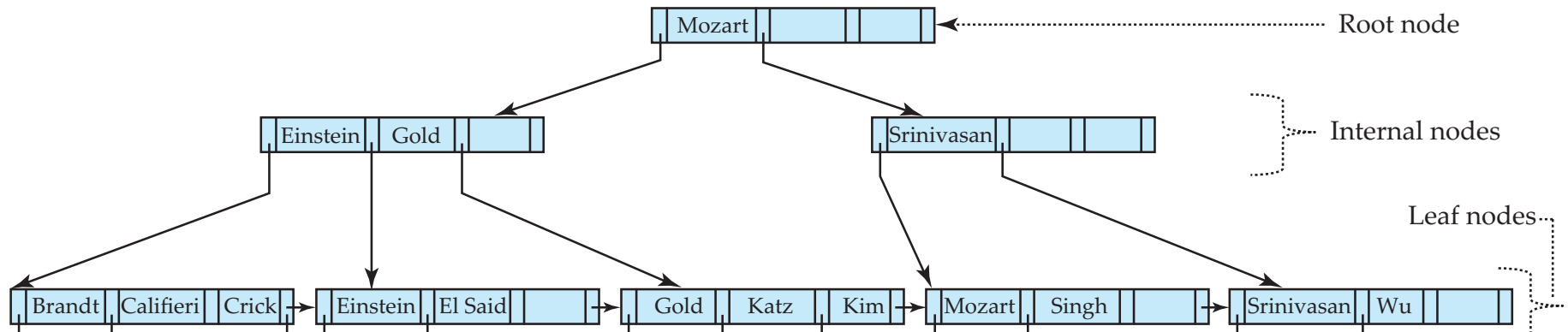
# Updates on B<sup>+</sup>-Trees: Insertion (Cont.)

- Splitting a leaf node:
  - take the  $n$  (search-key value, pointer) pairs (including the one being inserted) in sorted order. Place the first  $\lceil n/2 \rceil$  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.

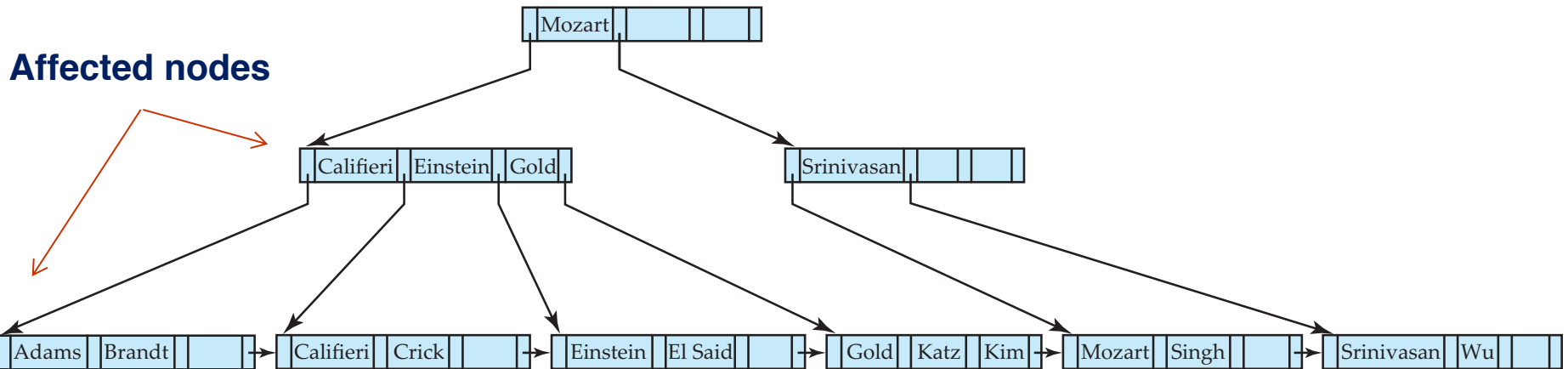


Result of splitting node containing Brandt, Califieri and Crick on inserting Adams  
Next step: insert entry with (Califieri, pointer-to-new-node) into parent

# B<sup>+</sup>-Tree Insertion

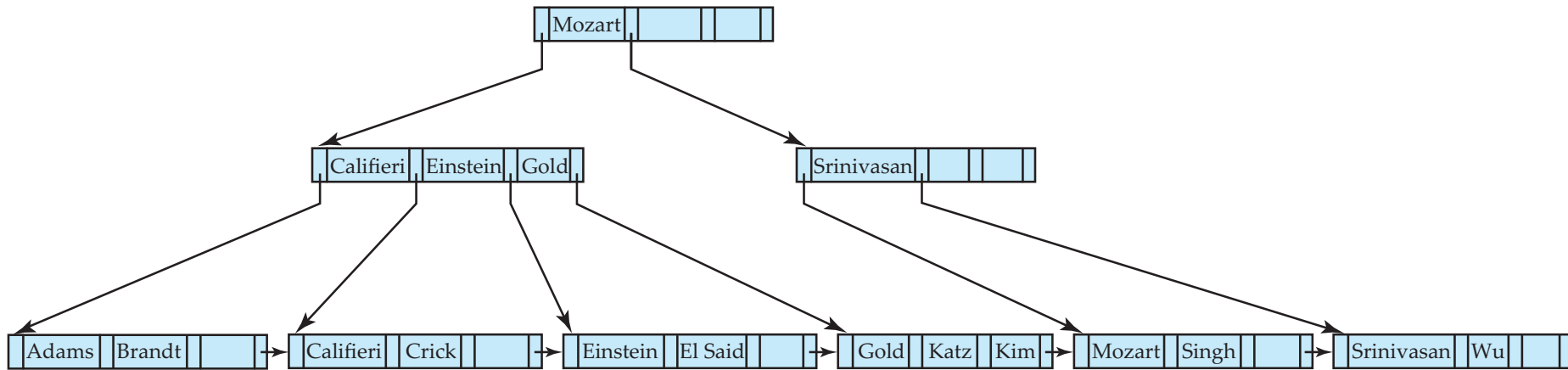


## B<sup>+</sup>-Tree before and after insertion of “Adams”

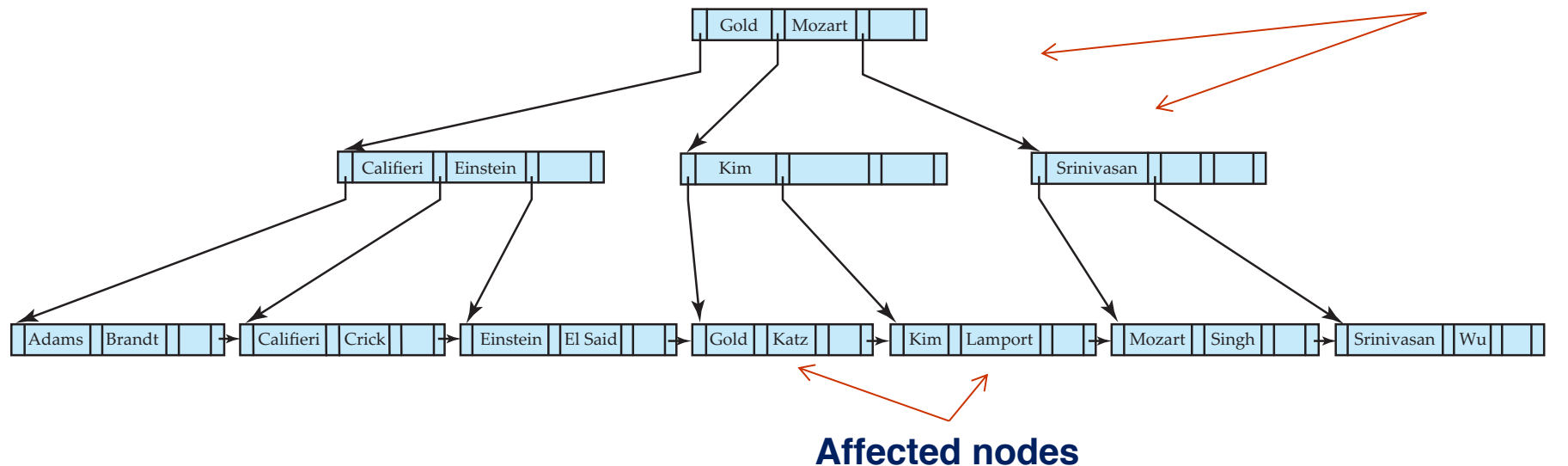


## B<sup>+</sup>-Tree before and after insertion of “Adams”

# B+-Tree Insertion



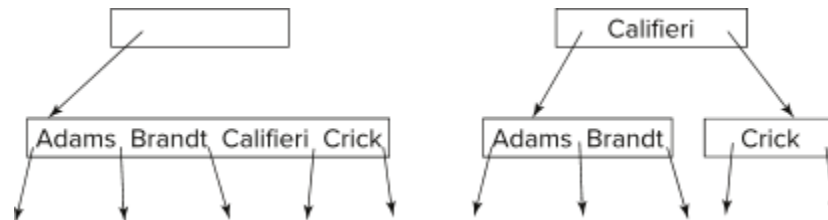
**B+-Tree before and after insertion of “Lampport”**





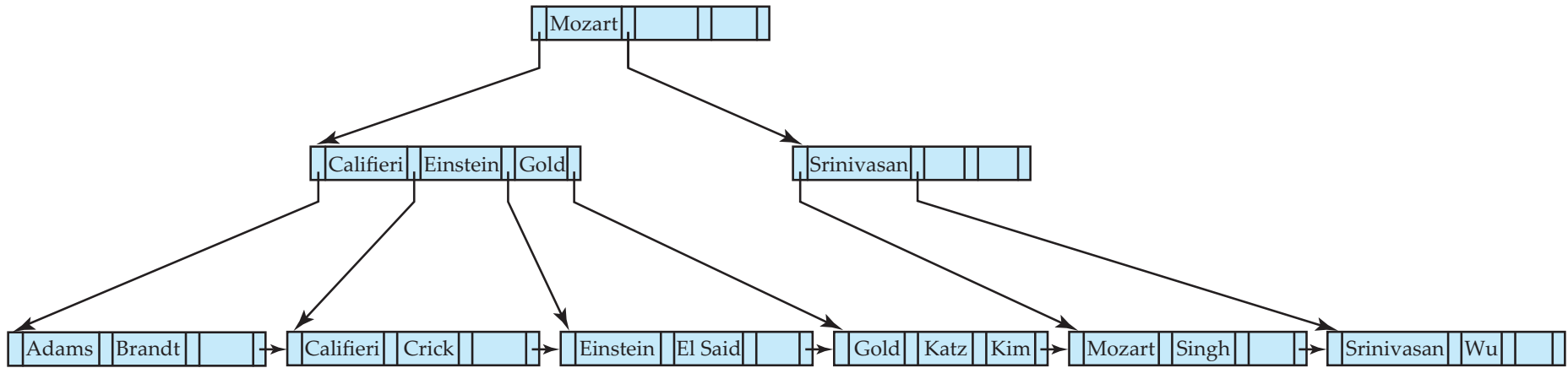
# Insertion in B<sup>+</sup>-Trees (Cont.)

- 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  $P_1, K_1, \dots, K_{\lceil n/2 \rceil - 1}, P_{\lceil n/2 \rceil}$  from M back into node N
  - Copy  $P_{\lceil n/2 \rceil + 1}, K_{\lceil n/2 \rceil + 1}, \dots, K_n, P_{n+1}$  from M into newly allocated node N'
  - Insert  $(K_{\lceil n/2 \rceil}, N')$  into parent N
- Example

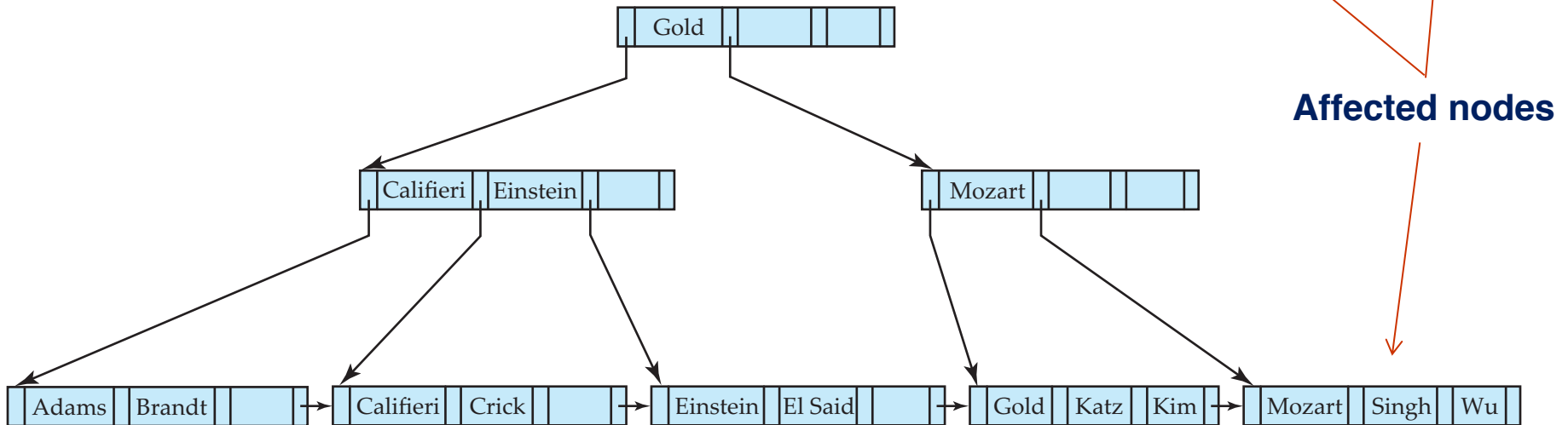


- **Read pseudocode in book!**

# Examples of B+-Tree Deletion



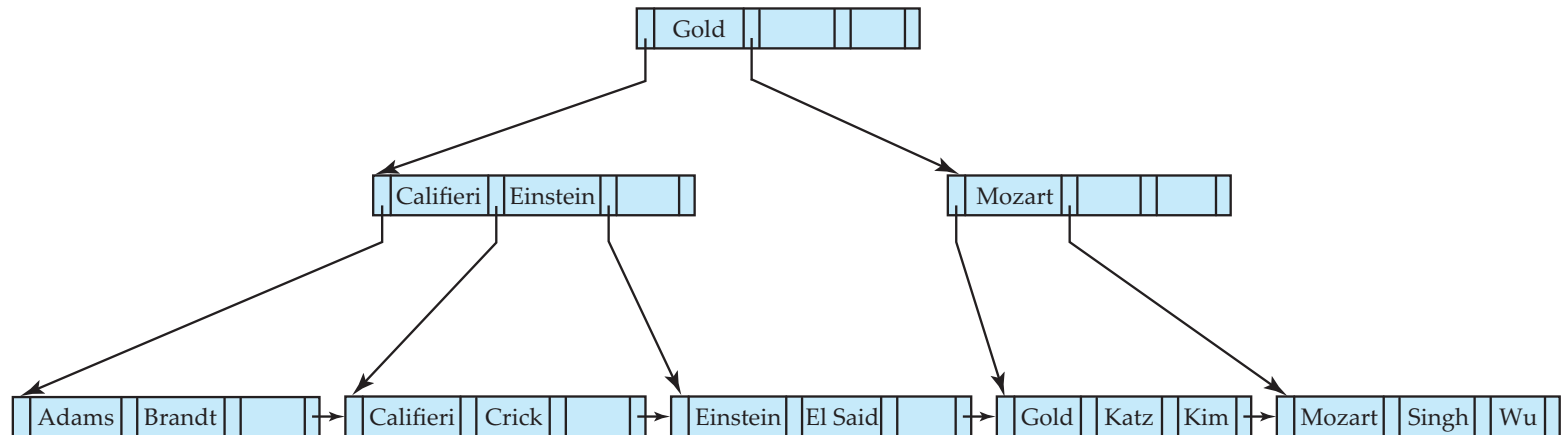
Before and after deleting “Srinivasan”



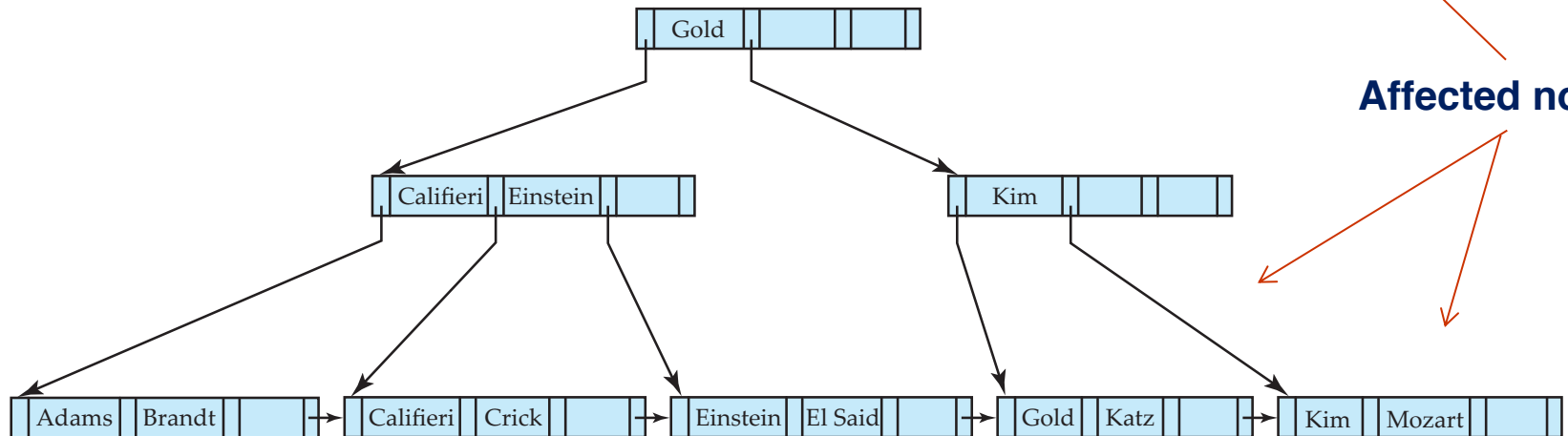
Affected nodes

- Deleting “Srinivasan” causes **merging** of under-full leaves

# Examples of B+-Tree Deletion (Cont.)

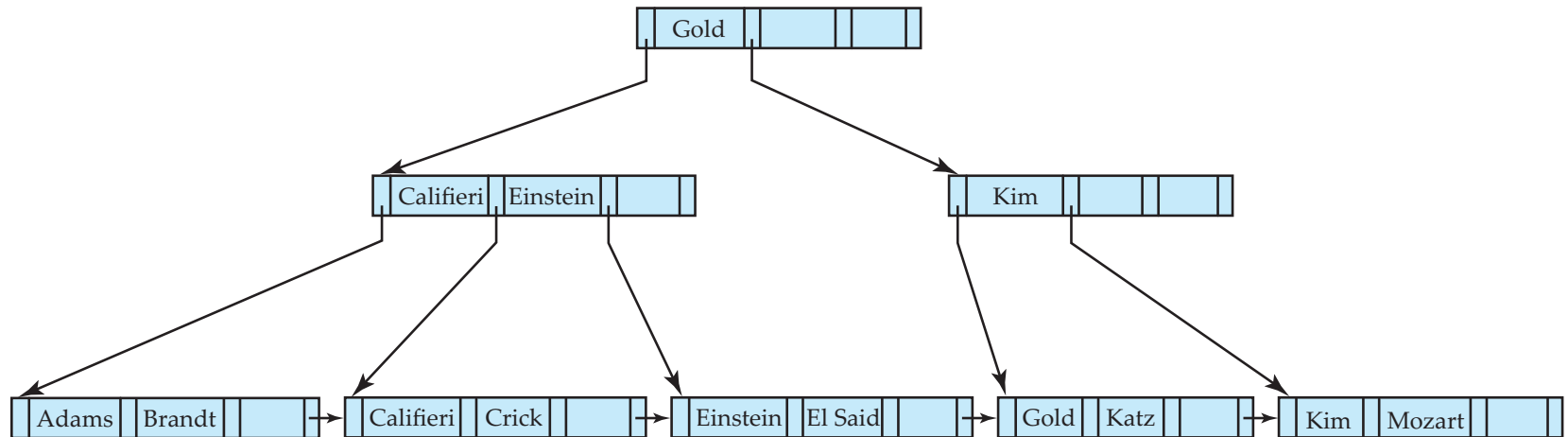


Before and after deleting “Singh” and “Wu”

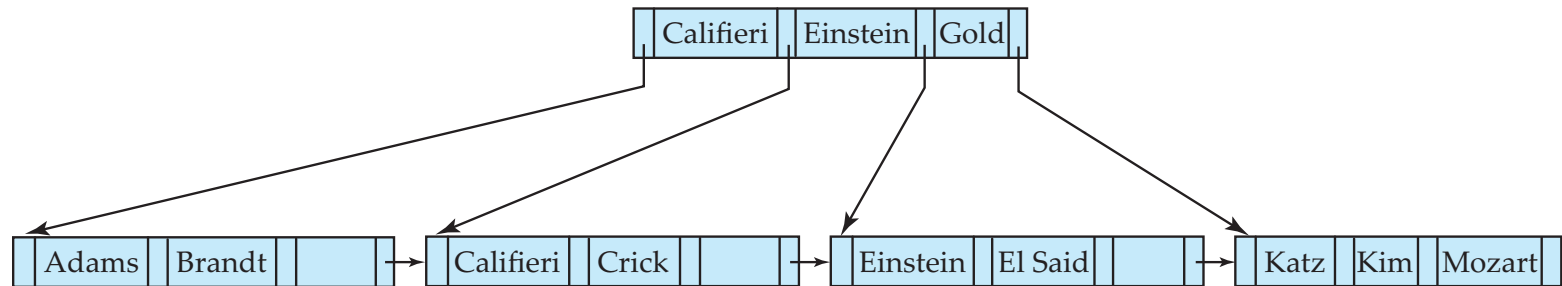


- Leaf containing Singh and Wu became underfull, and **borrowed a value** Kim from its left sibling
- Search-key value in the parent changes as a result

# Example of B<sup>+</sup>-tree Deletion (Cont.)



**Before and after deletion of “Gold”**



- 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

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

Assume record already deleted from file. Let  $V$  be the search key value of the record, and  $Pr$  be the pointer to the record.

- Remove  $(Pr, V)$  from the leaf node
- 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.

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

- Otherwise, 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.
- 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.
- You may experiment at:  
<https://www.cs.usfca.edu/~galles/visualization/BPlusTree.html>

# Complexity of Updates

- Cost (in terms of number of I/O operations) of insertion and deletion of a single entry proportional to height of the tree
  - With  $K$  entries and maximum fanout of  $n$ , worst case complexity of insert/delete of an entry is  $O(\log_{\lceil n/2 \rceil}(K))$
- In practice, number of I/O operations is less:
  - Internal nodes tend to be in buffer
  - Splits/merges are rare, most insert/delete operations only affect a leaf node
- Average node occupancy depends on insertion order
  - 2/3rds with random,  $\frac{1}{2}$  with insertion in sorted order

# 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?)
      - Worst case complexity may be linear!
    - 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

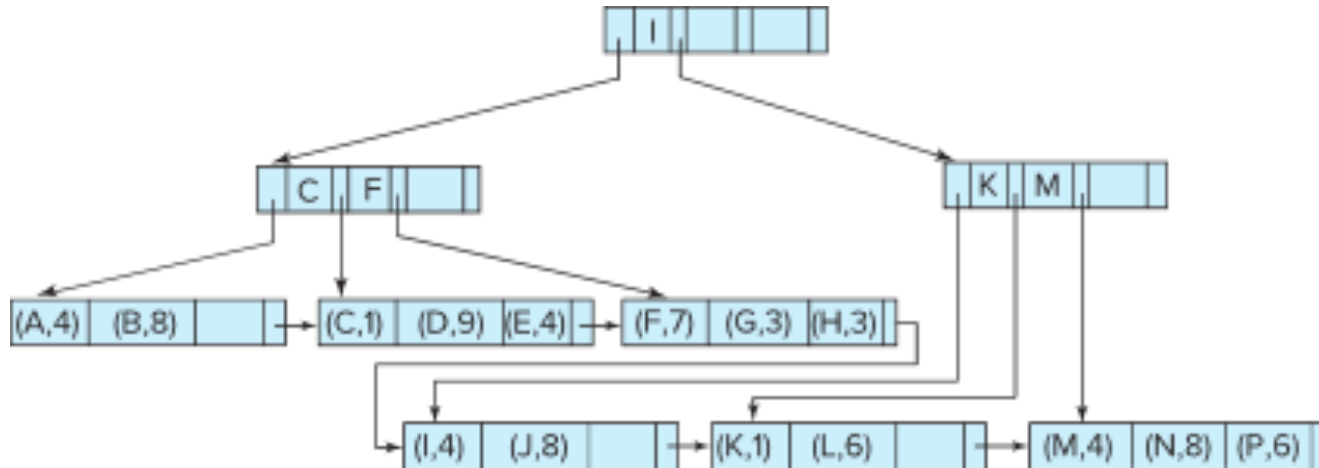


# B<sup>+</sup>-Tree File Organization

- B<sup>+</sup>-Tree File Organization:
  - Leaf nodes in a B<sup>+</sup>-tree file organization store records, instead of pointers
  - Helps keep data records clustered even when there are insertions/deletions/updates
- 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.

# B+-Tree File Organization (Cont.)

- Example of B+-tree File Organization



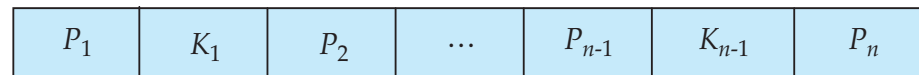
- Good space utilization important since records use more space than pointers.
- To improve space utilization, involve more sibling nodes in redistribution during splits and merges
  - Involving 2 siblings in redistribution (to avoid split / merge where possible) results in each node having at least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

# 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

# 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 B-tree; an additional pointer field for each search key in a nonleaf node must be included.
- Generalized B-tree leaf node



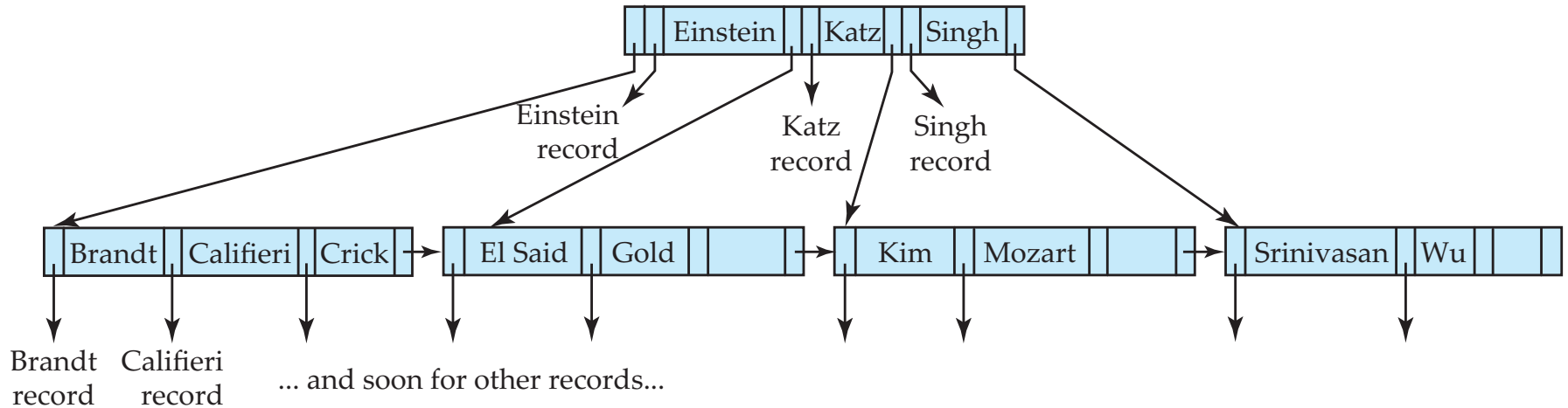
(a)



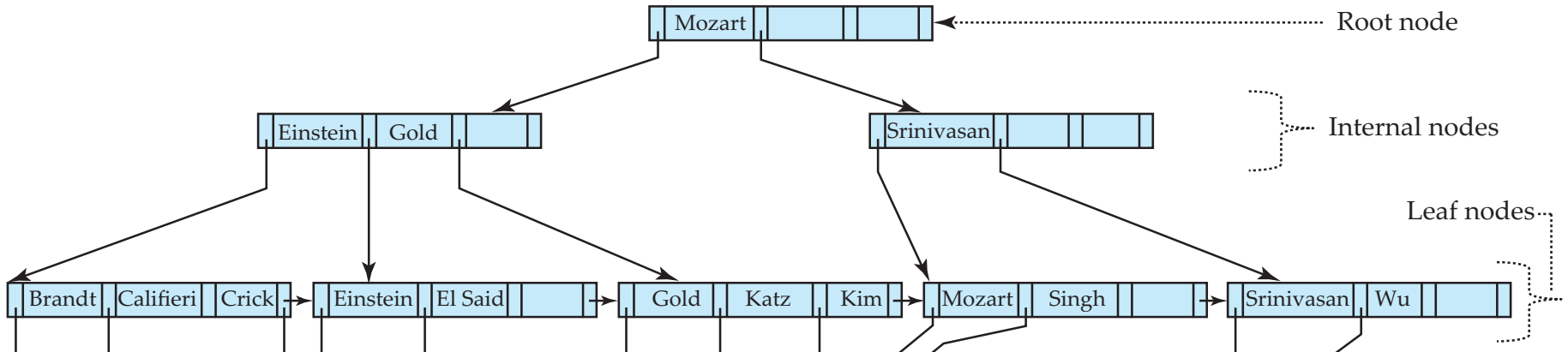
(b)

- Nonleaf node – pointers  $B_i$  are the bucket or file record pointers.

# B-Tree Index File Example



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



# B-Tree Index Files (Cont.)

- Advantages of B-Tree indices:
  - May use less tree nodes than a corresponding B<sup>+</sup>-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<sup>+</sup>-Trees
  - Implementation is harder than B<sup>+</sup>-Trees.
- Typically, advantages of B-Trees do not outweigh disadvantages.
  - Note that only a **very very small part** of the search-keys are in intermediate levels!
- DBMS always use B<sup>+</sup>-trees (though they sometimes refer to them as B-Trees)

# Indexing on Flash

- Random I/O cost much lower on flash
  - 20 to 100 microseconds for read/write
- Writes are not in-place, and (eventually) require a more expensive erase
- Optimum page size therefore much smaller
- Bulk-loading still useful since it minimizes page erases
- Write-optimized tree structures (to be discussed next week) have been adapted to minimize page writes for flash-optimized search trees

# Indexing in Main Memory

- Random access in memory
  - Much cheaper than on disk/flash
  - But still expensive compared to cache read
  - Data structures that make best use of cache preferable
  - Binary search for a key value within a large B<sup>+</sup>-tree node results in many cache misses
- B<sup>+</sup>- trees with small nodes that fit in cache line are preferable to reduce cache misses
- Key idea: use large node size to optimize disk access, but structure data within a node using a tree with small node size, instead of using an array.