# 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.

- Disadvantage of B$^+$-trees: extra insertion and deletion overhead, space overhead.

- Advantages of B$^+$-trees outweigh disadvantages, and they are used extensively.

# B$^+$-Tree Index Files (Cont.)

A B$^+$-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$^+$-Tree Node Structure

- Typical node

$$\boxed{P_1} \; \boxed{K_1} \; \boxed{P_2} \; \boxed{\ldots} \; \boxed{P_{n-1}} \; \boxed{K_{n-1}} \; \boxed{P_n}$$
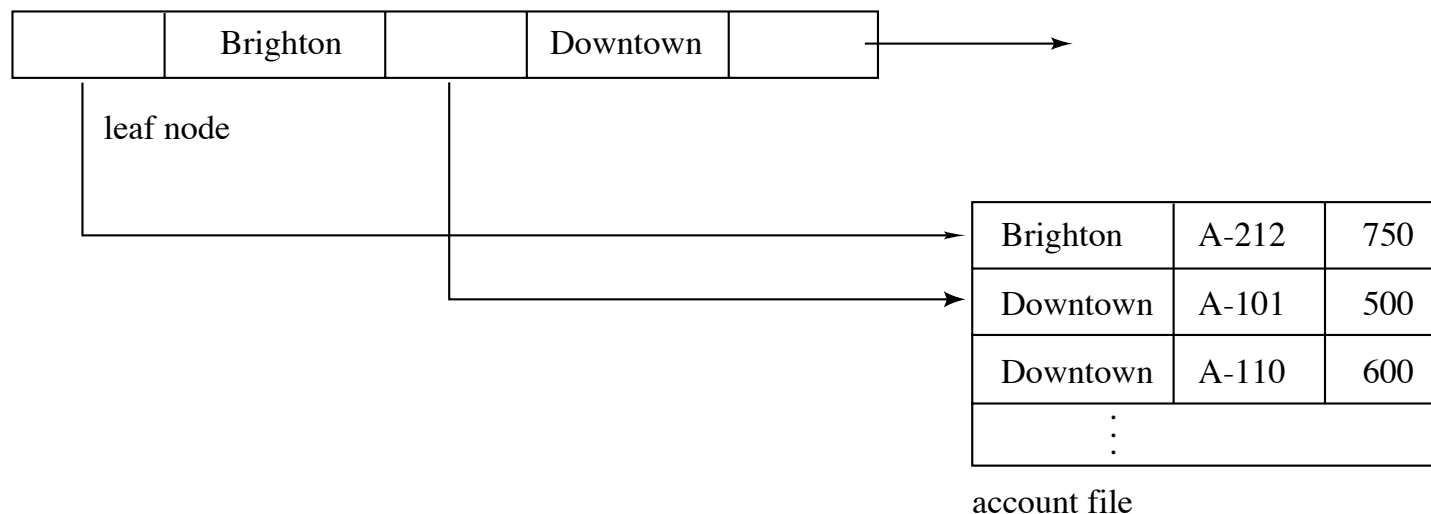
  – $K_i$ are the search-key values

  – $P_i$ 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

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

# Leaf Nodes in B⁺-Trees

Properties of a leaf node:

- For $i = 1, 2, \ldots, n-1$, pointer $P_i$ either points to a file record with search-key value $K_i$, or to a bucket of pointers to file records, each record having search-key value $K_i$. <u>Only need bucket structure if search-key does not form a primary key.</u>
- If $L_i$, $L_j$ are leaf nodes and $i < j$, $L_i$'s search-key values are less than $L_j$'s search-key values
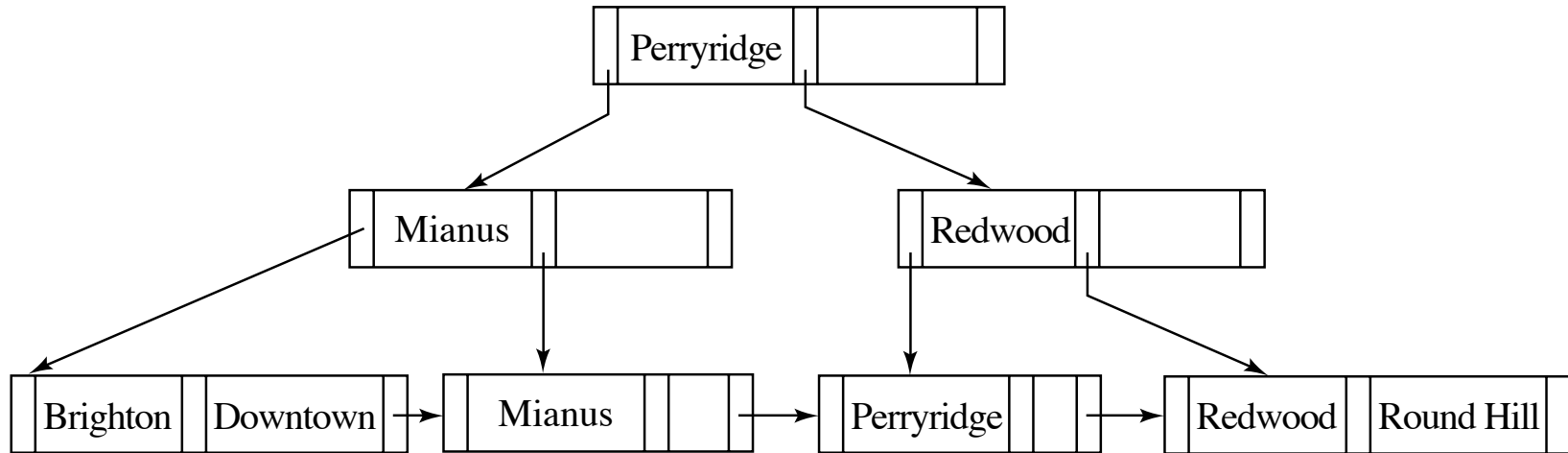- $P_n$ points to next leaf node in search-key order

| | Brighton | | Downtown | | |
|---|---|---|---|---|---|

leaf node

| Brighton | A-212 | 750 |
|---|---|---|
| Downtown | A-101 | 500 |
| Downtown | A-110 | 600 |
| ⋮ | | |

account file

# Non-Leaf Nodes in B$^+$-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 $P_m$ points are greater than or equal to $K_{m-1}$
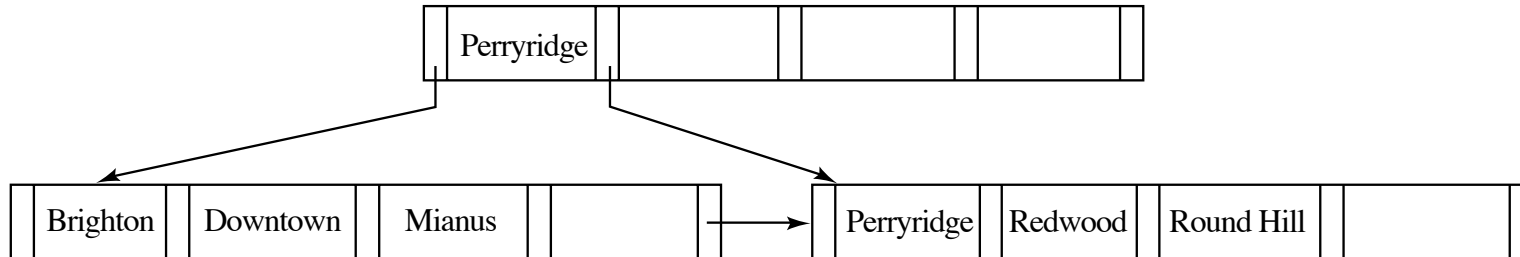
| $P_1$ | $K_1$ | $P_2$ | . . . | $P_{n-1}$ | $K_{n-1}$ | $P_n$ |
|---|---|---|---|---|---|---|

# Example of a B+-tree

```
                              ┌──┬──────────┬──┬──────────┬──┐
                              │  │Perryridge│  │          │  │
                              └──┴──────────┴──┴──────────┴──┘
                               │                  │
              ┌────────────────┘                  └──────────────────┐
              ▼                                                       ▼
      ┌──┬──────┬──┬────────┬──┐                      ┌──┬────────┬──┬────────┬──┐
      │  │Mianus│  │        │  │                      │  │Redwood │  │        │  │
      └──┴──────┴──┴────────┴──┘                      └──┴────────┴──┴────────┴──┘
       │              │                                │              │
   ┌───┘              └──────┐              ┌──────────┘              └────────┐
   ▼                         ▼              ▼                                  ▼
┌──────────┬──────────┬─┐ ┌─┬──────┬──┬─┐ ┌─┬──────────┬──┬─┐ ┌─┬───────┬─────────┬─┐
│ Brighton │ Downtown │ │→│ │Mianus│  │ │→│ │Perryridge│  │ │→│ │Redwood│Round Hill│ │
└──────────┴──────────┴─┘ └─┴──────┴──┴─┘ └─┴──────────┴──┴─┘ └─┴───────┴─────────┴─┘
```

B+-tree for *account* file ($n = 3$)

# Example of a B$^+$-tree

| | Perryridge | | | | | | |
|---|---|---|---|---|---|---|---|

| Brighton | | Downtown | | Mianus | | | | → | Perryridge | | Redwood | | Round Hill | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

B$^+$-tree for *account* file ($n = 5$)

- Leaf nodes must have between 2 and 4 values ($\lceil (n-1)/2 \rceil$ and $n - 1$, with $n = 5$).

- Non-leaf nodes other than root must have between 3 and 5 children ($\lceil n/2 \rceil$ and $n$ with $n = 5$).

- Root must have at least 2 children

# Observations about B$^+$-trees

- Since the inter-node connections are done by pointers, there is no assumption that in the B$^+$-tree, the "logically" close blocks are "physically" close.

- The non-leaf levels of the B$^+$-tree form a hierarchy of sparse indices.

- The B$^+$-tree contains a relatively small number of levels (logarithmic in the size of the main file), 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 see).

# Queries on B⁺-Trees

- Find all records with a search-key value of $k$.

  - Start with the root node
    * Examine the node for the smallest search-key value $> k$.
    * If such a value exists, assume it is $K_i$. Then follow $P_i$ to the child node
    * Otherwise $k \geq K_{m-1}$, where there are $m$ pointers in the node. Then follow $P_m$ to the child node.

  - If the node reached by following the pointer above is not a leaf node, repeat the above procedure on the node, and follow the corresponding pointer.

  - Eventually reach a leaf node. If key $K_i = k$, follow pointer $P_i$ to the desired record or bucket. Else no record with search-key value $k$ exists.

# Queries on B$^+$-Trees (Cont.)

- In processing a query, a path is traversed in the tree from the root to some leaf node.

- If there are $K$ search-key values in the file, the path is no longer 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.

- 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 30 millisecond!
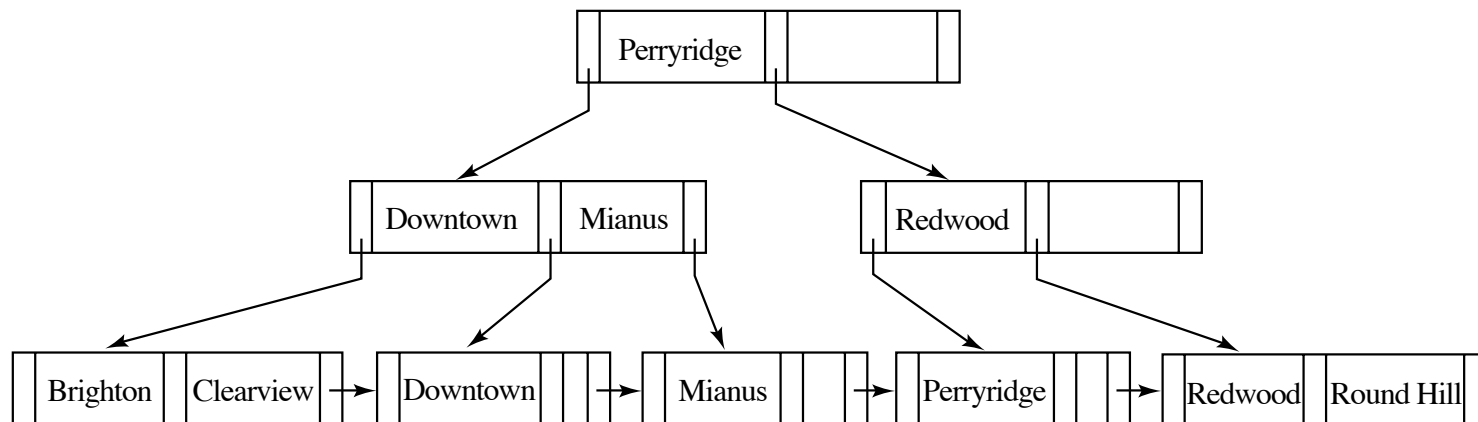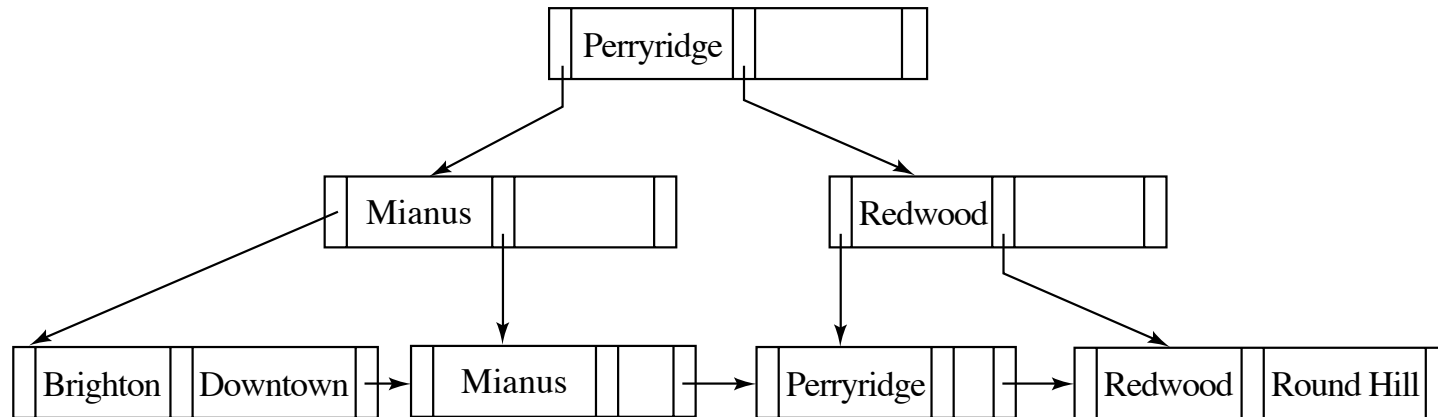
# Updates on B$^+$-Trees: Insertion

- Find the leaf node in which the search-key value would appear

- If the search-key value is already there in the leaf node, record is added to file and if necessary pointer is inserted into bucket.

- If the search-key value is not there, then add the record to the main file and create bucket if necessary. Then:

  - if there is room in the leaf node, insert (search-key value, record/bucket pointer) pair into leaf node at appropriate position.

  - if there is no room in the leaf node, split it and insert (search-key value, record/bucket pointer) pair as discussed in the next slide.

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

- Splitting a 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.

- The 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.

# Updates on B+-Trees: Insertion (Cont.)
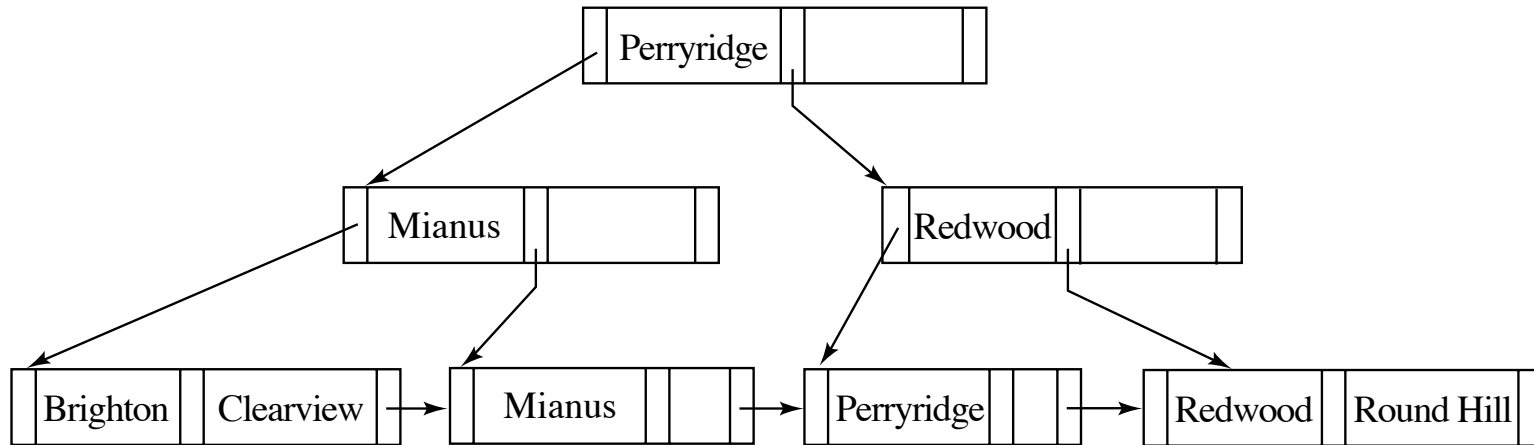


B+-Tree before and after insertion of "Clearview"

# 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, and the entries in the node and a sibling fit into a single node, then

  - 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$^+$-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

  - 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.
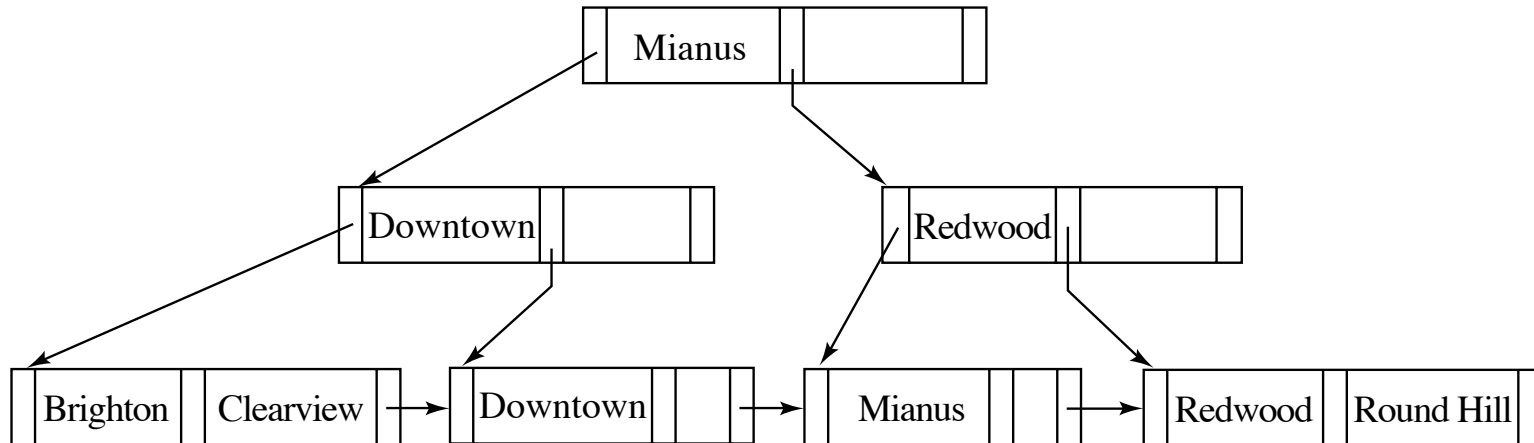
# Examples of B⁺-Tree Deletion



Result after deleting "Downtown" from *account*

- The removal of the leaf node containing "Downtown" did not result in its parent having too little pointers. So the cascaded deletions stopped with the deleted leaf node's parent.

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

```
                              ┌──────┬──┬──────┐
                              │Mianus│  │      │
                              └──────┴──┴──────┘
                   ┌─────────────┘        └──────────┐
          ┌────────┬──┬────┐            ┌───────┬──┬────┐
          │Downtown│  │    │            │Redwood│  │    │
          └────────┴──┴────┘            └───────┴──┴────┘
      ┌──────┘        └────┐        ┌──────┘        └──────┐
```

┌────────┬─┬─────────┐   ┌────────┬─┬─┐   ┌──────┬─┬─┐   ┌───────┬─┬──────────┐
│Brighton│ │Clearview│ → │Downtown│ │ │→│Mianus│ │ │→│Redwood│ │Round Hill│
└────────┴─┴─────────┘   └────────┴─┴─┘   └──────┴─┴─┘   └───────┴─┴──────────┘

Deletion of "Perryridge" instead of "Downtown"

- The deleted "Perryridge" node's parent became too small, but its sibling did not have space to accept one more pointer. So redistribution is performed. Observe that the root node's search-key value changes as a result.

# B$^+$-Tree File Organization

- Index file degradation problem is solved by using B$^+$-Tree Indices. Data file degradation problem is solved by using B$^+$-Tree File Organization.

- The leaf nodes in a B$^+$-tree file organization store records, instead of pointers.

- Since records are large 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.

- Leaf nodes are still required to be half full.

- Insertion and deletion are handled in the same way as insertion and deletion of entries in a B$^+$-tree index.

- Good space utilization is important since records use more space than pointers. To improve space utilization, involve more sibling nodes in redistribution during splits and merges.

# 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

| $P_1$ | $K_1$ | $P_2$ | . . . | $P_{n-1}$ | $K_{n-1}$ | $P_n$ |
|---|---|---|---|---|---|---|

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

| $P_1$ | $B_1$ | $K_1$ | $P_2$ | $B_2$ | $K_2$ | . . . | $P_{m-1}$ | $B_{m-1}$ | $K_{m-1}$ | $P_m$ |
|---|---|---|---|---|---|---|---|---|---|---|

# 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 $B^+$-Tree
  - Insertion and deletion more complicated than in $B^+$-Trees
  - Implementation is harder than $B^+$-Trees.

- Typically, advantages of B-Trees do not outweigh disadvantages.