# Indexing

Instructor: Dong-Wan Choi (dchoi@inha.ac.kr)
Assistant Professor @ CSE Department, Inha Univ.

# Outline

- Basic Concepts

- Ordered Indices

- $B^+$-Tree Index Files

- B-Tree Index Files

# Motivation

- What is an index in databases?
  - Simply, it is a well-designed data structure aiming to locate records that you want to find.
  - Why so many types of indices?
    - Various data types
    - Various queries
  - Very roughly speaking, what a data engineer does is to design and use appropriate indices on underlying data.

- Why do we need an index anyway?
  - To efficiently perform such a query, for example, as follows:

    **select** name, dept_name, salary
    **from** instructor
    **where** ID = 15151;

# Basic Concepts

- Indexing mechanisms used to speed up access to desired data.
  - E.g., author catalog in library
- **Search Key** - attribute to set of attributes used to look up records in a file.
- An **index file** consists of records (called **index entries**) of the form

| search-key | pointer |
|------------|---------|

- Index files are typically much smaller than the original file
- Two basic kinds of indices:
  - **Ordered indices:**  search keys are stored in sorted order
  - **Hash indices:**  search keys are distributed uniformly across "buckets" using a "hash function".

# Index Evaluation Metrics

- Access types supported efficiently.  E.g.,
  - records with a specified value in the attribute
  - or records with an attribute value falling in a specified range of values.
- Access time
- Insertion time
- Deletion time
- Space overhead

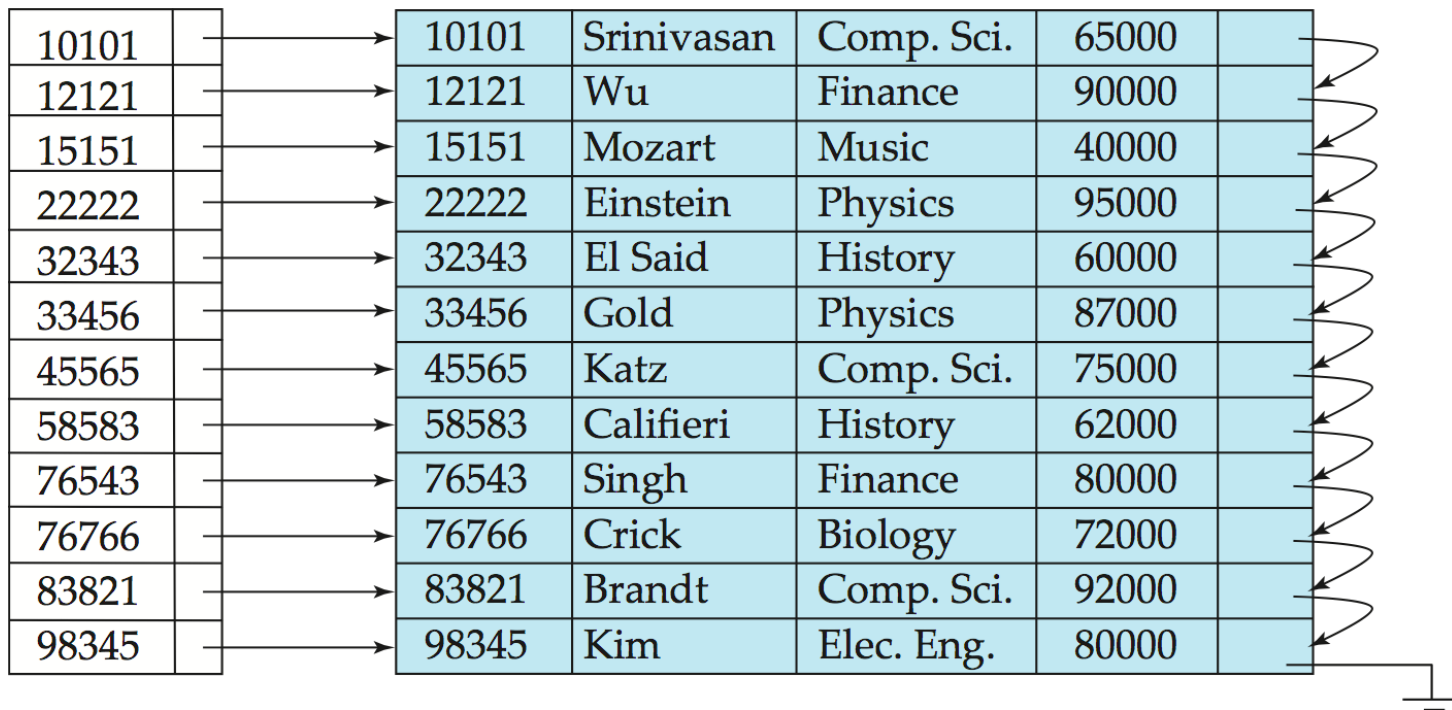**How to analyze the performance of an index** theoretically?
External-memory model (refer to the lecture note)

# Ordered Indices

- In an **ordered index,** index entries are stored sorted on the search key value, e.g., author catalog in library.

- **Primary index:** in a sequentially ordered file, the index whose search key specifies the sequential order of the file.
  - Also called **clustering index**
  - The search key of a primary index is usually but not necessarily the primary key.

- **Secondary index**: an index whose search key specifies an order different from the sequential order of the file. Also called non-clustering index.

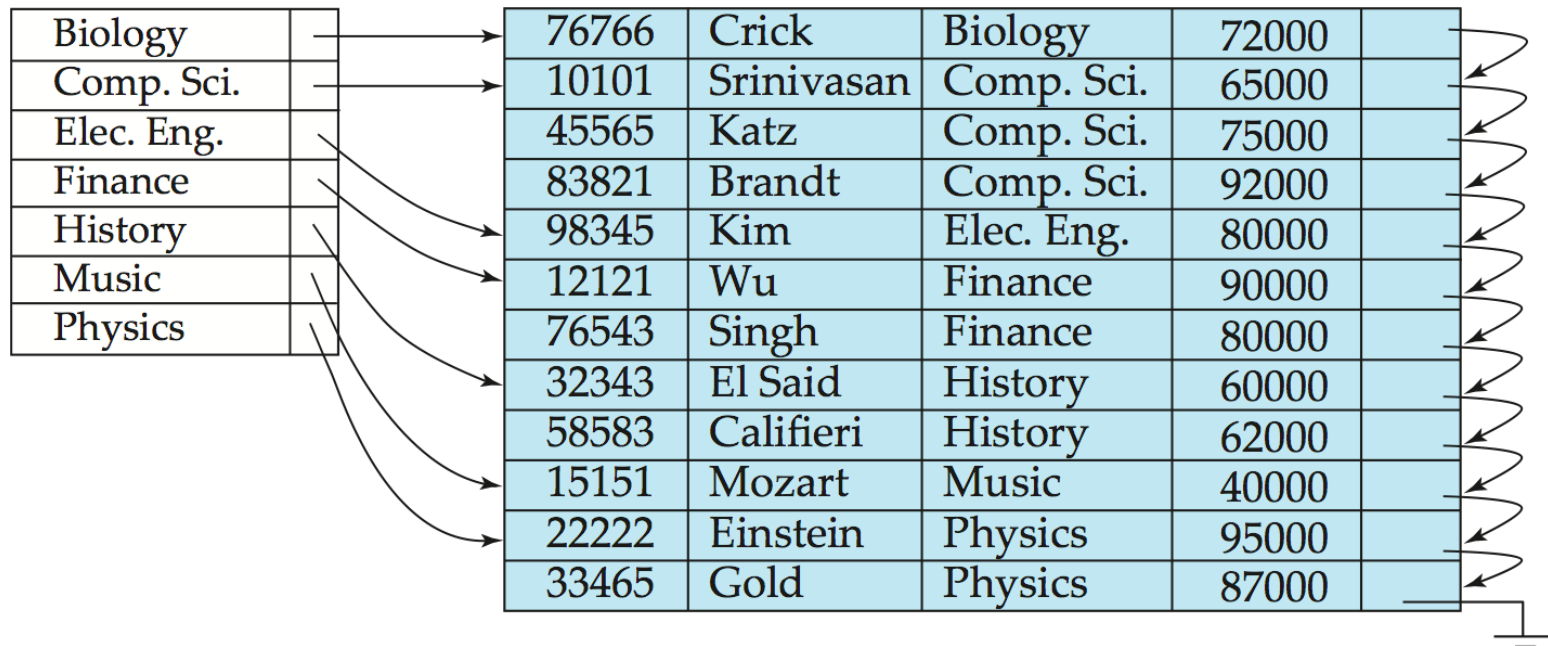- Index-sequential file: ordered sequential file with a primary index.

# Dense Index Files

- **Dense index** — Index record appears for every search-key value in the file.
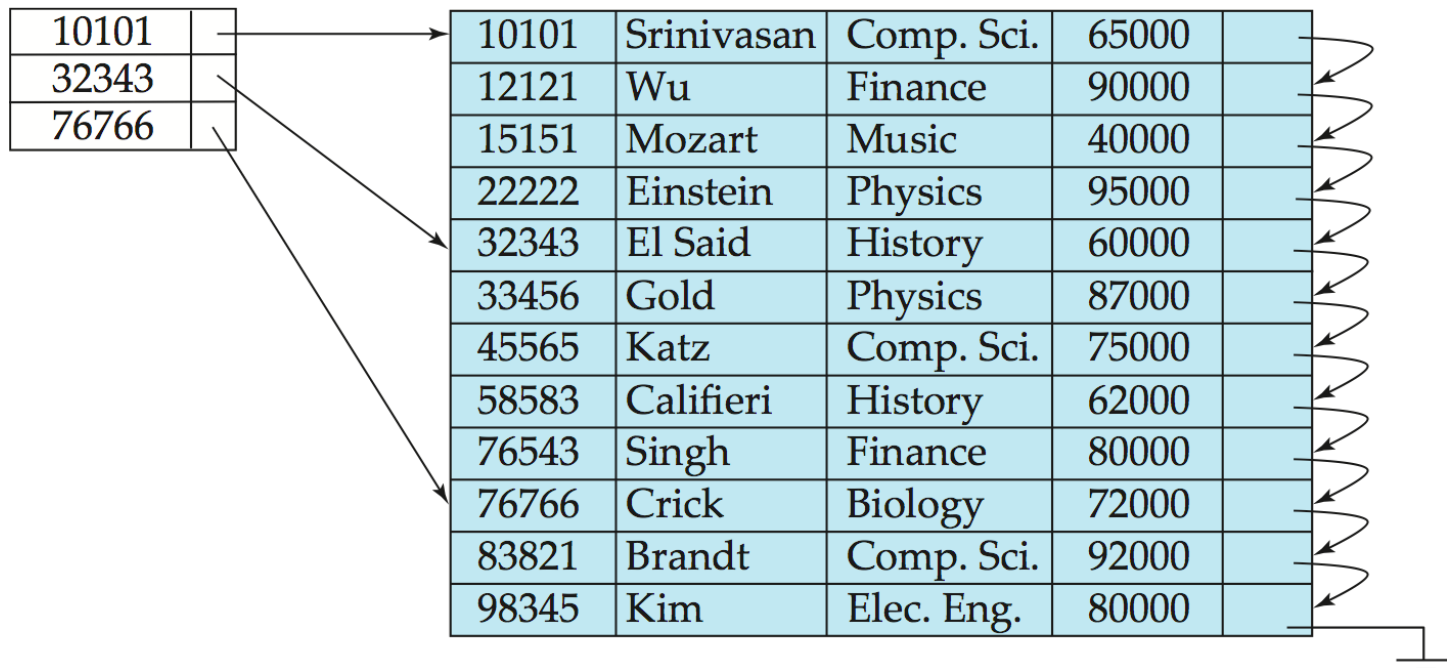- E.g. index on *ID* attribute of *instructor* relation

| | | | | |
|---|---|---|---|---|
| 10101 | 10101 | Srinivasan | Comp. Sci. | 65000 |
| 12121 | 12121 | Wu | Finance | 90000 |
| 15151 | 15151 | Mozart | Music | 40000 |
| 22222 | 22222 | Einstein | Physics | 95000 |
| 32343 | 32343 | El Said | History | 60000 |
| 33456 | 33456 | Gold | Physics | 87000 |
| 45565 | 45565 | Katz | Comp. Sci. | 75000 |
| 58583 | 58583 | Califieri | History | 62000 |
| 76543 | 76543 | Singh | Finance | 80000 |
| 76766 | 76766 | Crick | Biology | 72000 |
| 83821 | 83821 | Brandt | Comp. Sci. | 92000 |
| 98345 | 98345 | Kim | Elec. Eng. | 80000 |

# Dense Index Files (Cont.)

- Dense index on *dept_name*, with *instructor* file sorted on *dept_name*

| | | | | |
|---|---|---|---|---|
| Biology | 76766 | Crick | Biology | 72000 |
| Comp. Sci. | 10101 | Srinivasan | Comp. Sci. | 65000 |
| Elec. Eng. | 45565 | Katz | Comp. Sci. | 75000 |
| Finance | 83821 | Brandt | Comp. Sci. | 92000 |
| History | 98345 | Kim | Elec. Eng. | 80000 |
| Music | 12121 | Wu | Finance | 90000 |
| Physics | 76543 | Singh | Finance | 80000 |
| | 32343 | El Said | History | 60000 |
| | 58583 | Califieri | History | 62000 |
| | 15151 | Mozart | Music | 40000 |
| | 22222 | Einstein | Physics | 95000 |
| | 33465 | Gold | Physics | 87000 |

# Sparse Index Files

- **Sparse Index**: contains index records for only some search-key values.
  - Applicable when records are sequentially ordered on search-key
- To locate a record with search-key value $K$ we:
  - Find index record with largest search-key value < $K$
  - Search file sequentially starting at the record to which the index record points
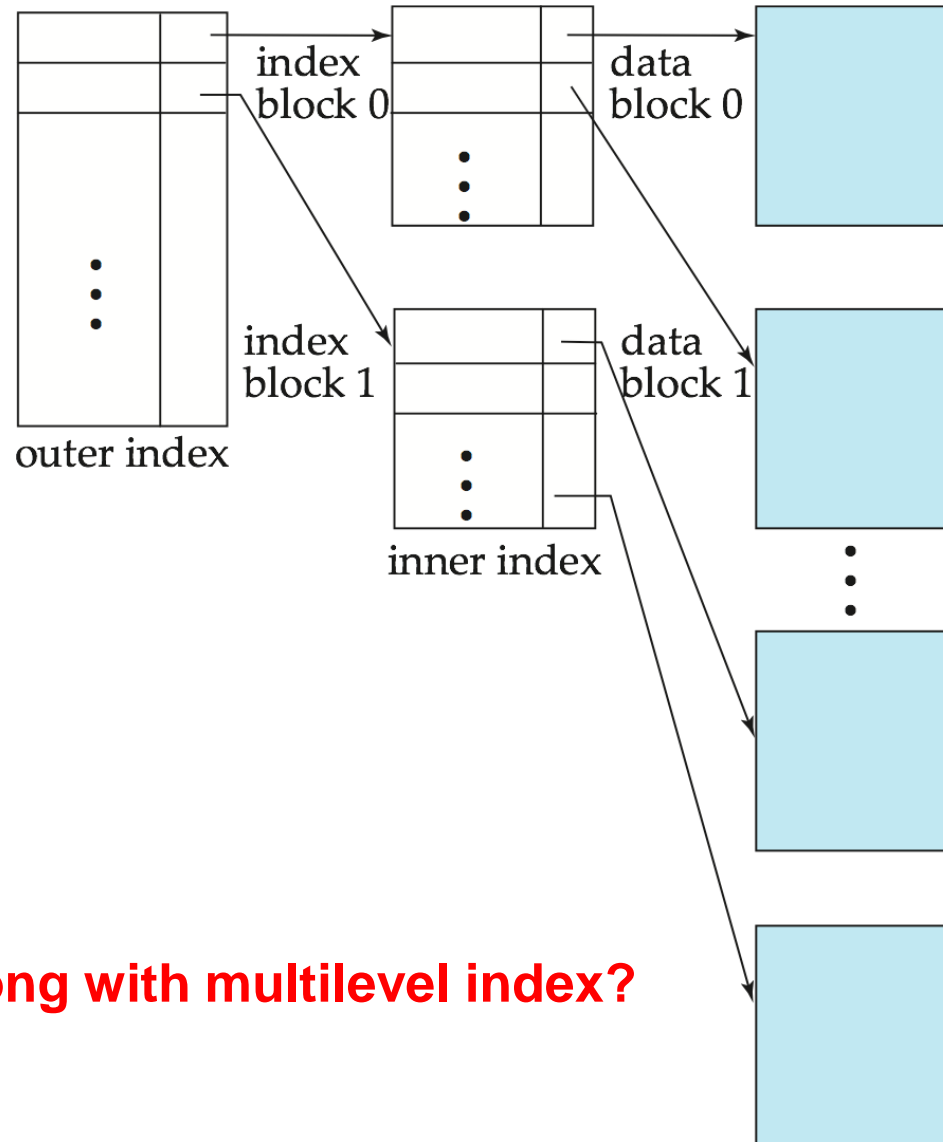
# Sparse Index Files (Cont.)

- Compared to dense indices:
  - Less space and less maintenance overhead for insertions and deletions.
  - Generally slower than dense index for locating records.
- **Good tradeoff**: sparse index with an index entry for every block in file, corresponding to least search-key value in the block. **Why?**

# Multilevel Index

- If primary index does not fit in memory, access becomes expensive.

- Solution: treat primary index kept on disk as a sequential file and construct a sparse index on it.
  - outer index – a sparse index of primary index
  - inner index – the primary index file

- If even outer index is too large to fit in main memory, yet another level of index can be created, and so on.

- Indices at all levels must be updated on insertion or deletion from the file.

# Multilevel Index (Cont.)



**What is wrong with multilevel index?**

# Index Update:  Deletion

| 10101 | · |
|-------|---|
| 32343 | · |
| 76766 | · |

| 10101 | Srinivasan | Comp. Sci. | 65000 | |
|-------|-----------|-----------|-------|---|
| 12121 | Wu | Finance | 90000 | |
| 15151 | Mozart | Music | 40000 | |
| 22222 | Einstein | Physics | 95000 | |
| 32343 | El Said | History | 60000 | |
| 33456 | Gold | Physics | 87000 | |
| 45565 | Katz | Comp. Sci. | 75000 | |
| 58583 | Califieri | History | 62000 | |
| 76543 | Singh | Finance | 80000 | |
| 76766 | Crick | Biology | 72000 | |
| 83821 | Brandt | Comp. Sci. | 92000 | |
| 98345 | Kim | Elec. Eng. | 80000 | |

- If deleted record was the only record in the file with its particular search-key value, the search-key is deleted from the index also.

- Single-level index entry deletion:
  - Dense indices – deletion of search-key is similar to file record deletion.
  - Sparse indices
    - If an entry for the search key exists in the index, it is deleted by replacing the entry in the index with the next search-key value in the file (in search-key order).
    - If the next search-key value already has an index entry, the entry is deleted instead of being replaced.

# Index Update:  Insertion

- **Single-level index insertion:**
  - Perform a lookup using the search-key value appearing in the record to be inserted.
  - **Dense indices** – if the search-key value does not appear in the index, insert it.
  - **Sparse indices** – if index stores an entry for each block of the file, no change needs to be made to the index unless a new block is created.
    - *If a new block is created, the first search-key value appearing in the new block is inserted into the index.*
- **Multilevel insertion and deletion:**  algorithms are simple extensions of the single-level algorithms

# Secondary Indices

- Frequently, one wants to find all the records whose values in a certain field (which is not the search-key of the primary index) satisfy some condition.
    - Example 1: In the *instructor* relation stored sequentially by ID, we may want to find all instructors in a particular department
    - Example 2: as above, but where we want to find all instructors with a specified salary or with salary in a specified range of values
- We can have a secondary index with an index record for each search-key value

# Secondary Indices Example



**Secondary index on *salary* field of *instructor***

- Index record points to a bucket that contains pointers to all the actual records with that particular search-key value.
- Secondary indices have to be dense

# Primary and Secondary Indices

- Indices offer substantial benefits when searching for records.
- BUT: Updating indices imposes overhead on database modification --when a file is modified, every index on the file must be updated,
- Sequential scan using primary index is efficient, but a sequential scan using a secondary index is expensive
  - Each record access may fetch a new block from disk
  - Block fetch requires about 5 to 10 milliseconds, versus about 100 nanoseconds for memory access

# 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

# Why B$^+$-Tree?

- Every relational database system supports B$^+$-Trees

  - IBM DB2, Informix, MS SQL, Oracle, Sybase, SQLite, MySQL, PostgreSQL, Tibero, etc

  - They support other types of indices, too, but the most widely used one is the B$^+$-Tree and its variants (B-Tree, B-*Tree, etc)

- Almost every file system uses B$^+$-Trees

  - NTFS, EXT4, ReiserFS, NSS, XFS, JFS, ReFS, BFS, etc

# 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

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

  - $K_i$ 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

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

(Initially assume no duplicate keys, address duplicates later)

# Leaf Nodes in B⁺-Trees

## Properties of a leaf node:

- For $i = 1, 2, ..., n-1$, 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

# 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 $P_n$ points have values greater than or equal to $K_{n-1}$

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

# Example of B⁺-tree



B⁺-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⁺-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⁺-tree form a hierarchy of sparse indices.

- The B⁺-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 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 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. $V \leq K_i$.
     ② If no such exists, set *C = last non-null pointer in C*
     ③ 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.

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

- If there are $K$ search-key values in the file, the height of the tree is no more than $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$.

- 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 20 milliseconds

# Updates on B⁺-Trees:  Insertion

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

II.   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.

III.  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⁺-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+-Tree Insertion



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

# B⁺-Tree Insertion
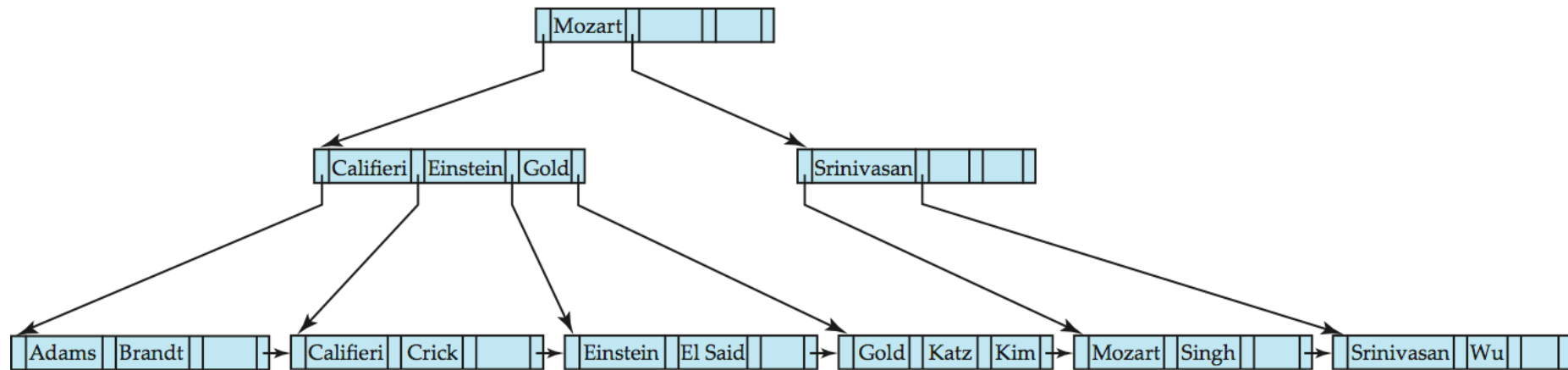


B⁺-Tree before and after insertion of "Lamport"
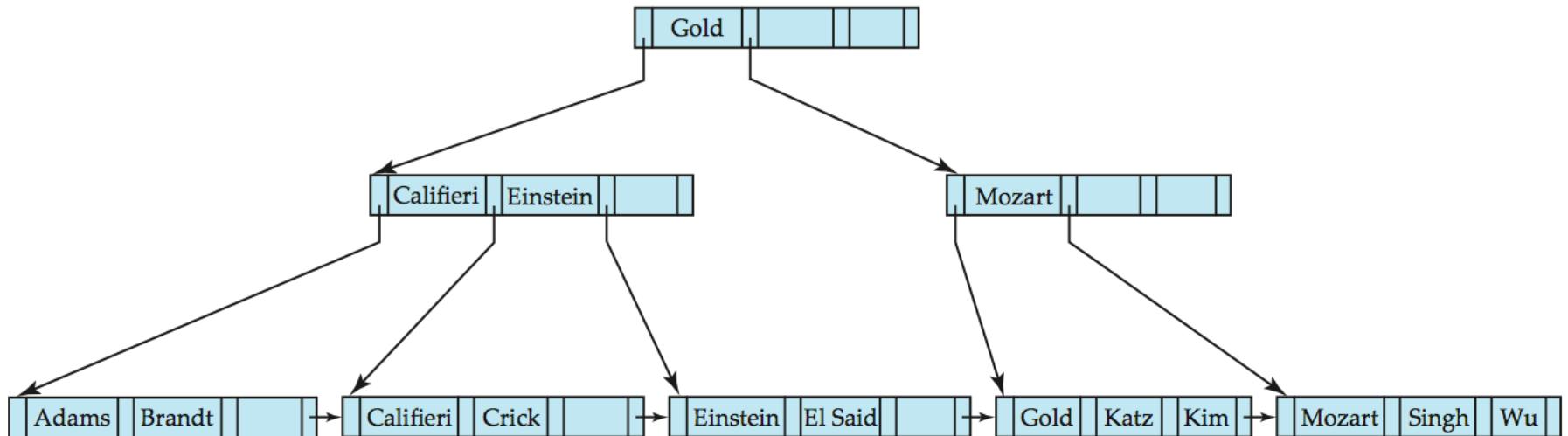
# Insertion in B⁺-Trees (Cont.)

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

# Examples of B⁺-Tree Deletion



Before and after deleting "Srinivasan"

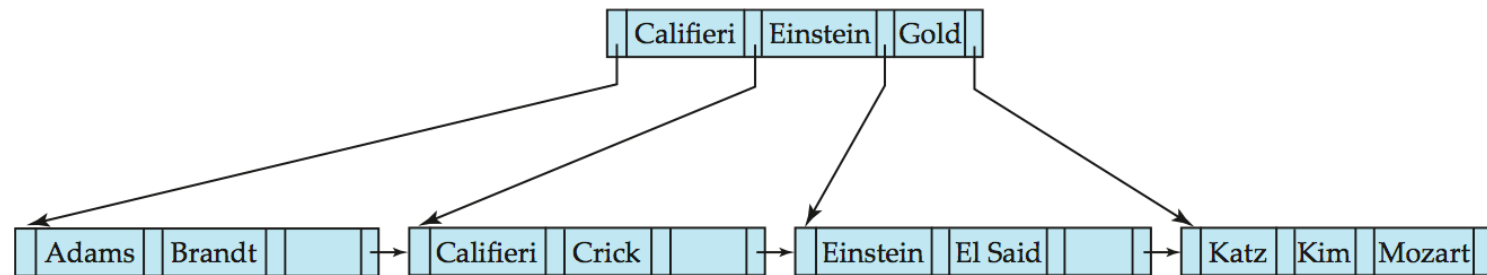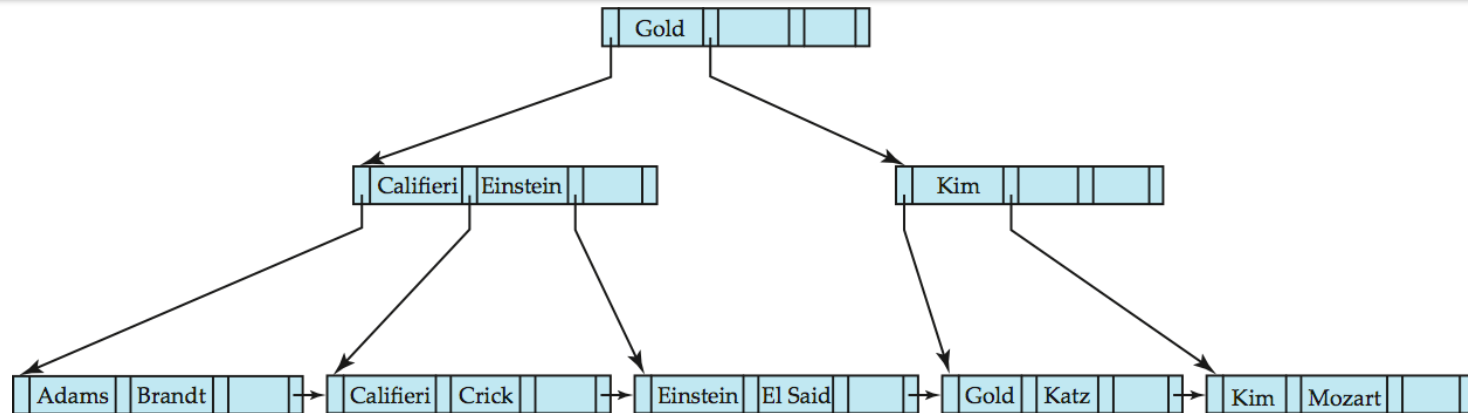

Deleting "Srinivasan" causes merging of under-full leaves

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

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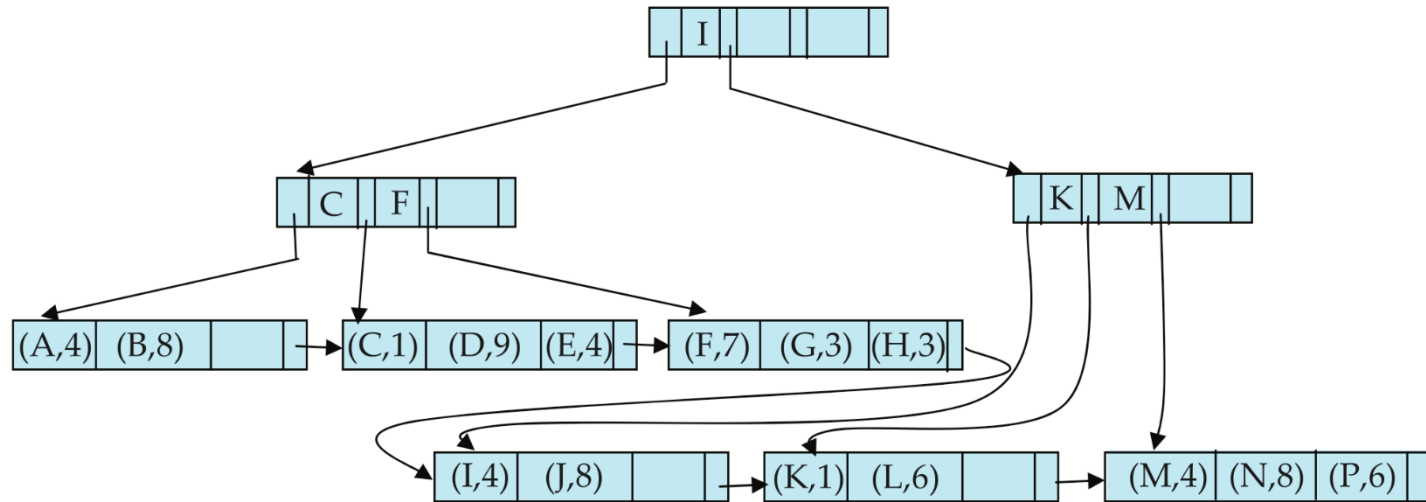
# Updates on B⁺-Trees:  Deletion

- 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.

# 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.

- 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⁺-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

# 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$^+$-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
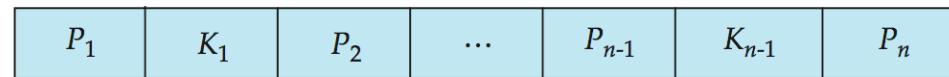
# B-Tree Index Files
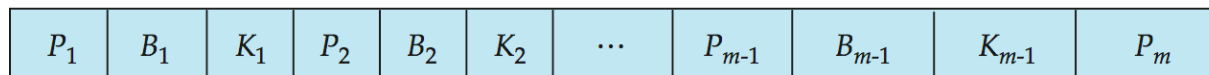


**Rudolf Bayer**
(1939 ~ present)

Bayer, R.; McCreight, E. (1970), "Organization and Maintenance of Large Ordered Indexes"

# 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.
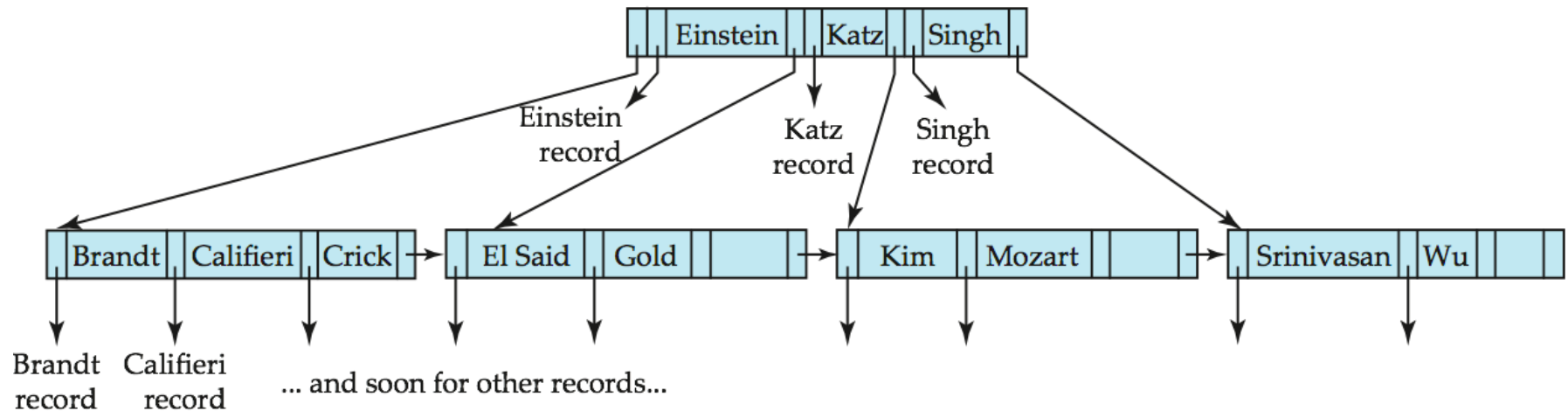
- Generalized B-tree leaf node (a)

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

(a)

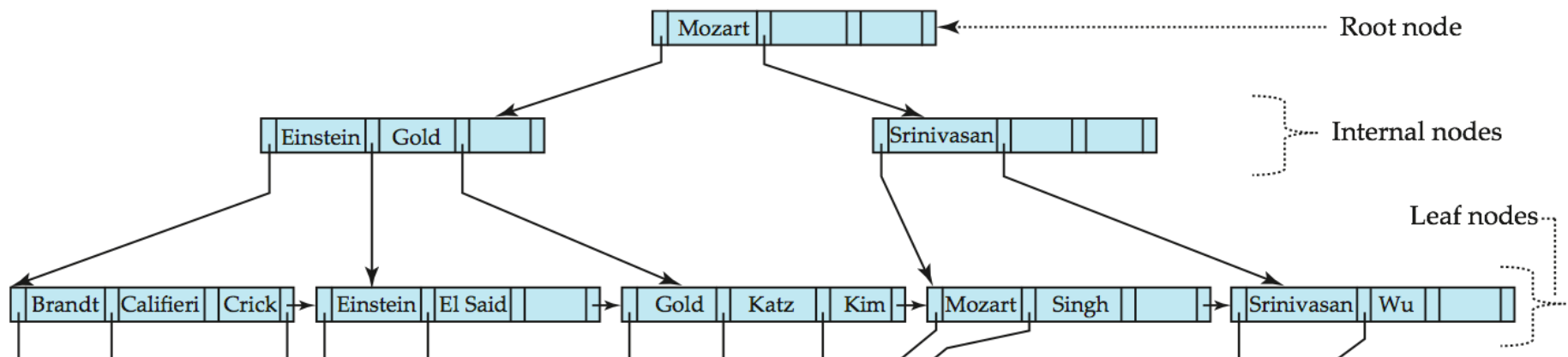| $P_1$ | $B_1$ | $K_1$ | $P_2$ | $B_2$ | $K_2$ | ... | $P_{m-1}$ | $B_{m-1}$ | $K_{m-1}$ | $P_m$ |
|---|---|---|---|---|---|---|---|---|---|---|

(b)

- Nonleaf node (b) – pointers Bi 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$^+$-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 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

- **Composite search keys** 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
    - $a_1 < b_1$, 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
- But cannot efficiently handle
    **where** *dept_name* < "Finance" **and** *balance* = 80000
    - May fetch many records that satisfy the first but not the second condition