



Indexing and Hashing

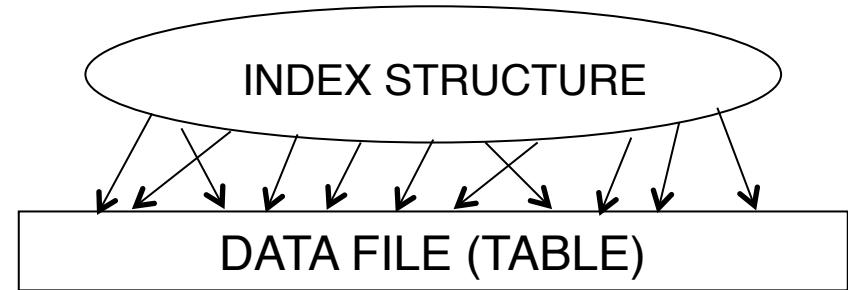
- Basic Concepts
- Ordered Indices
- B⁺-Tree Index Files
- B-Tree Index Files
- Static Hashing
- Dynamic Hashing
- Comparison of Ordered Indexing and Hashing
- Index Definition in SQL



Basic Concepts

- An index is a data structure that supports efficient retrieval of “certain types of tuples” from a table

- E.g., author catalog in library
- With no index:
 - ▶ Scan entire table
 - ▶ Binary search if sorted



- An **index file** consists of records (called **index entries**) of this form



- **Search Key** – attribute (or attributes) used to look up records in a file.
- 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”.



Typical Problem

“Find records of a certain type quickly”

- Equality Queries

$$\sigma_{\text{ssn} = 619441789}(\text{Employees})$$

- Range Queries

$$\sigma_{10 \leq \text{age} \leq 20}(\text{Employees})$$

These are the most interesting types of conditions

Plus AND and OR of several such conditions

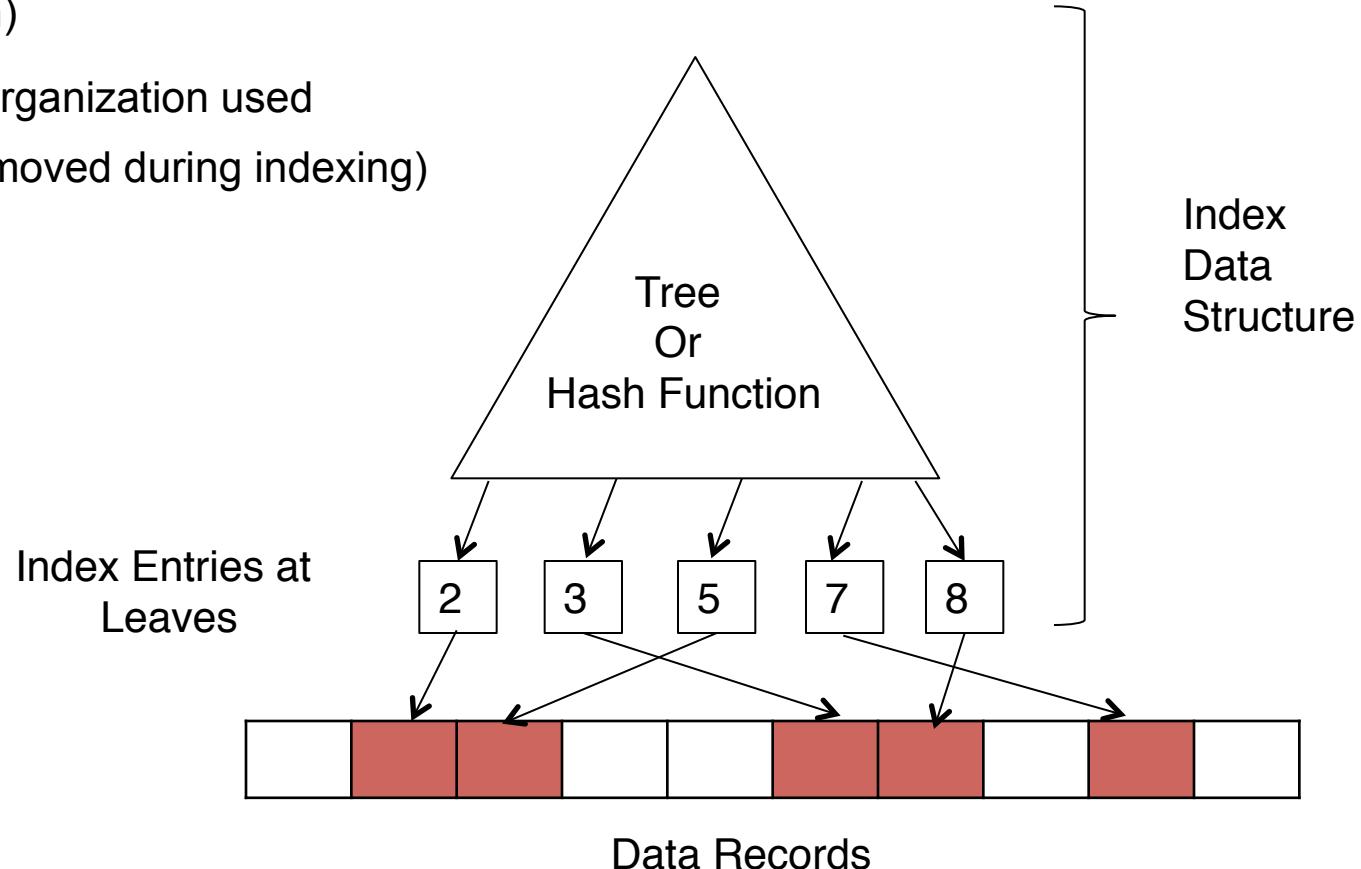
Goal: faster processing of SQL queries (selects and joins)



Indexing: Basic Setup

Indexing Hides:

- Type of structure used
(tree or hash)
- File/record organization used
(data is not moved during indexing)





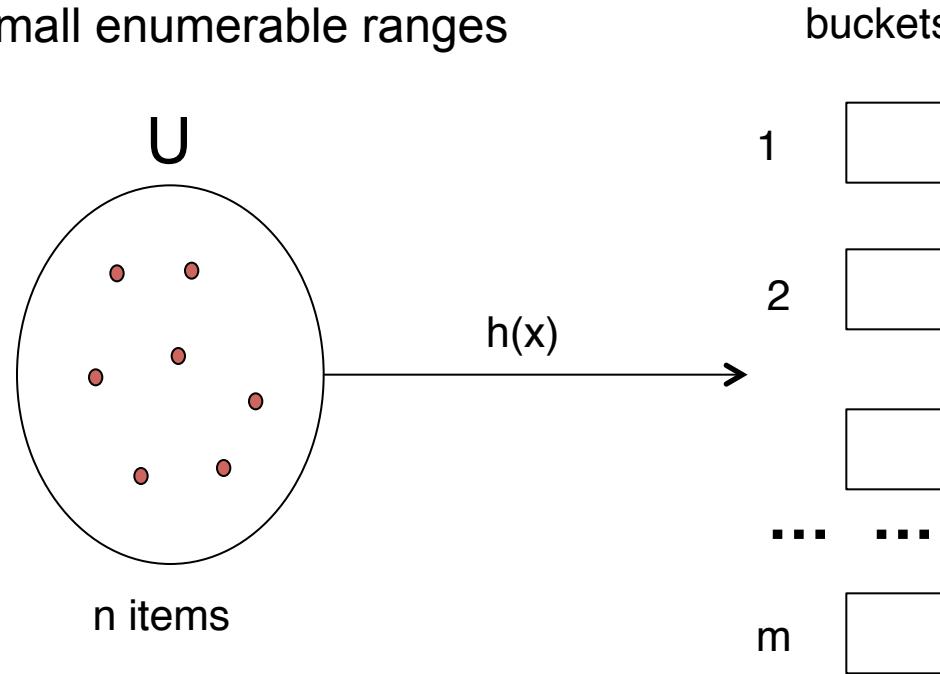
Ordered Indices

- 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.



Hash Indices

- **Hash Indices**: search keys are distributed uniformly across buckets using a hash function
- X = key value
- $h: U \rightarrow [1,2,3,\dots,m]$
- No range queries supported
- Exception: small enumerable ranges





Dense Index Files

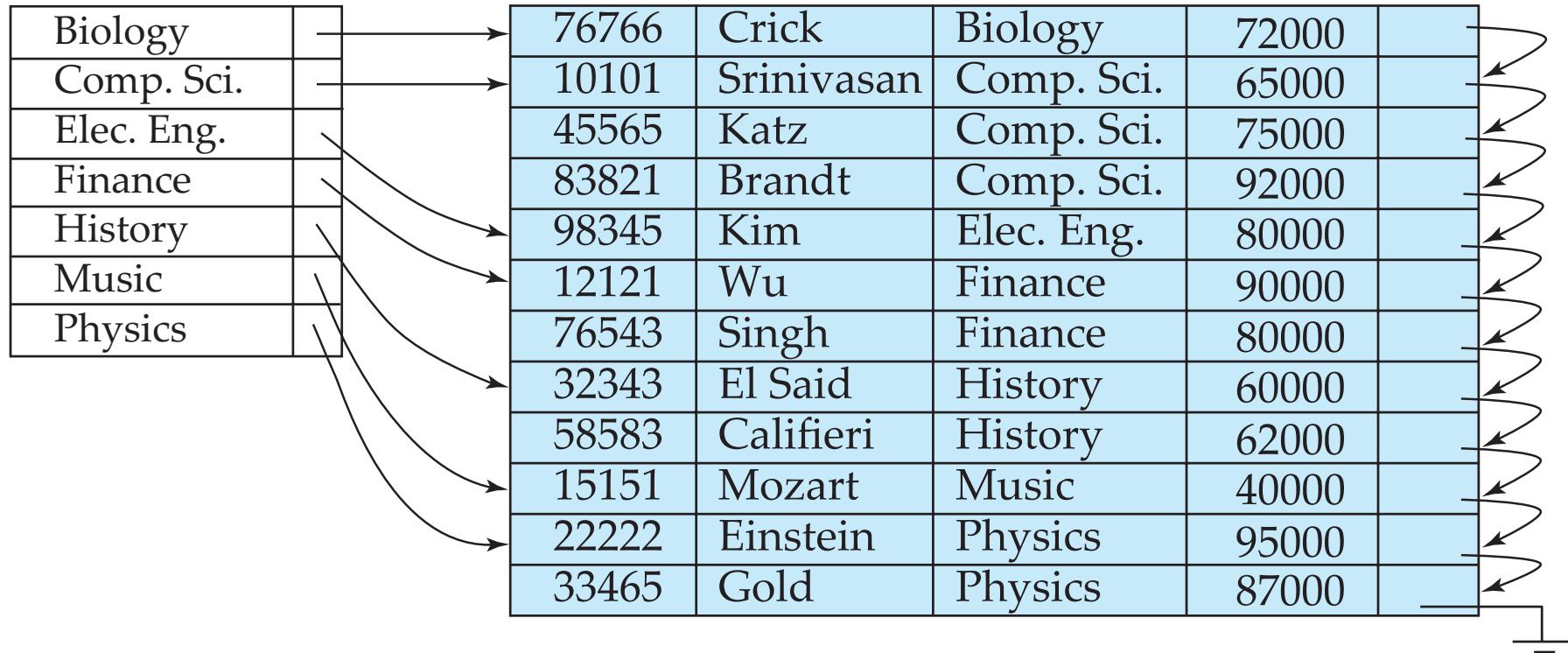
- **Dense index** — Index record appears for every search-key value in the file.
- E.g. index on *ID* attribute of *instructor* relation

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.)

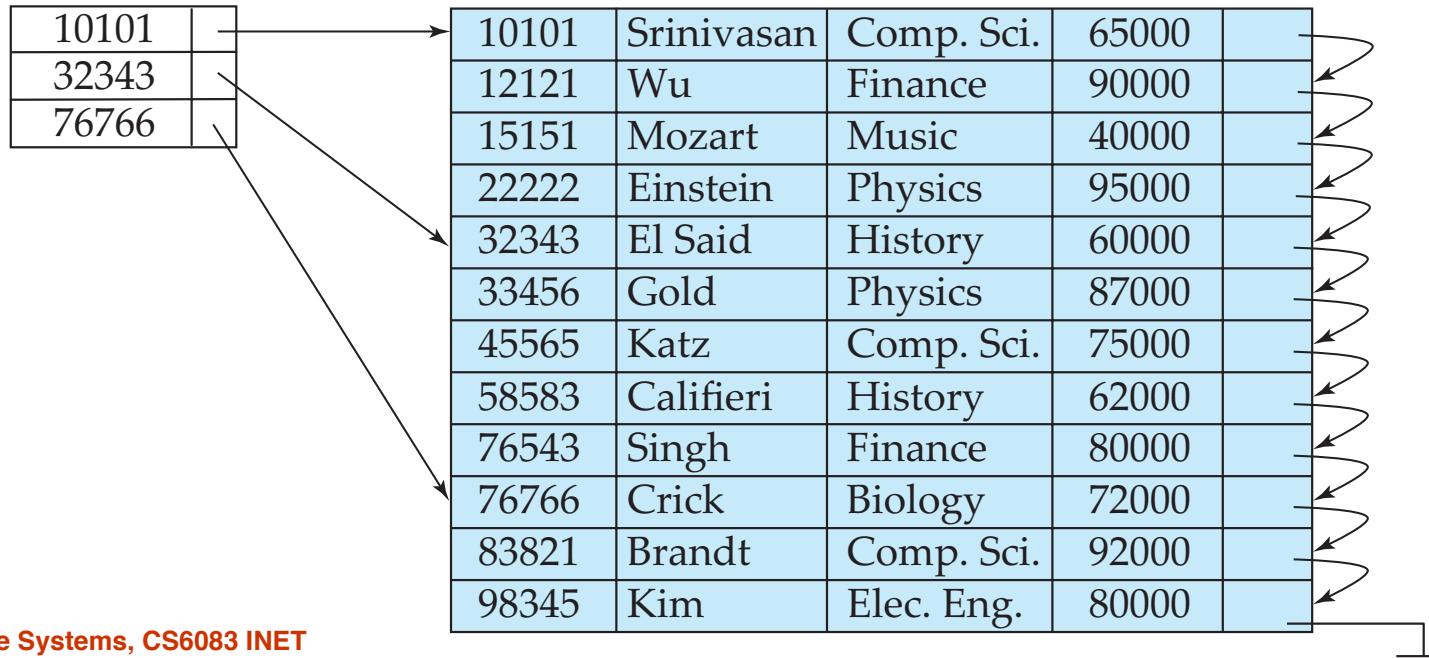
- Dense index on *dept_name*, with *instructor* file sorted on *dept_name*





Sparse Index Files

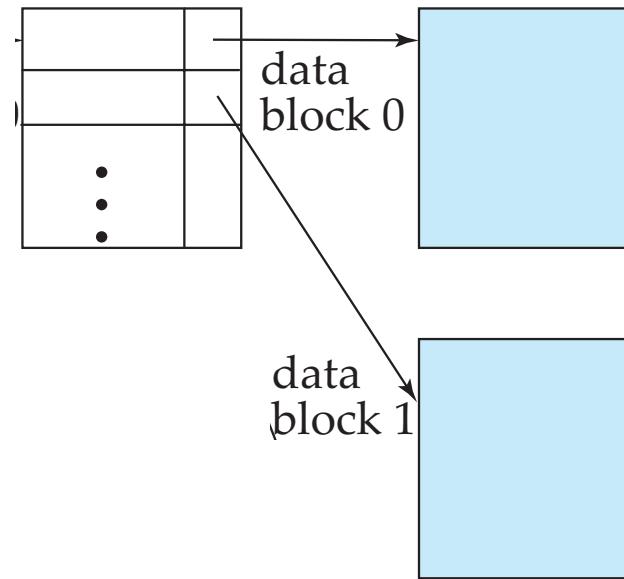
- **Sparse Index:** contains index records for only some search-key values.
 - Applicable when records are sequentially ordered on search-key
- To locate a record with search-key value K we:
 - Find index record with largest search-key value $< K$
 - Search file sequentially starting at the record to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





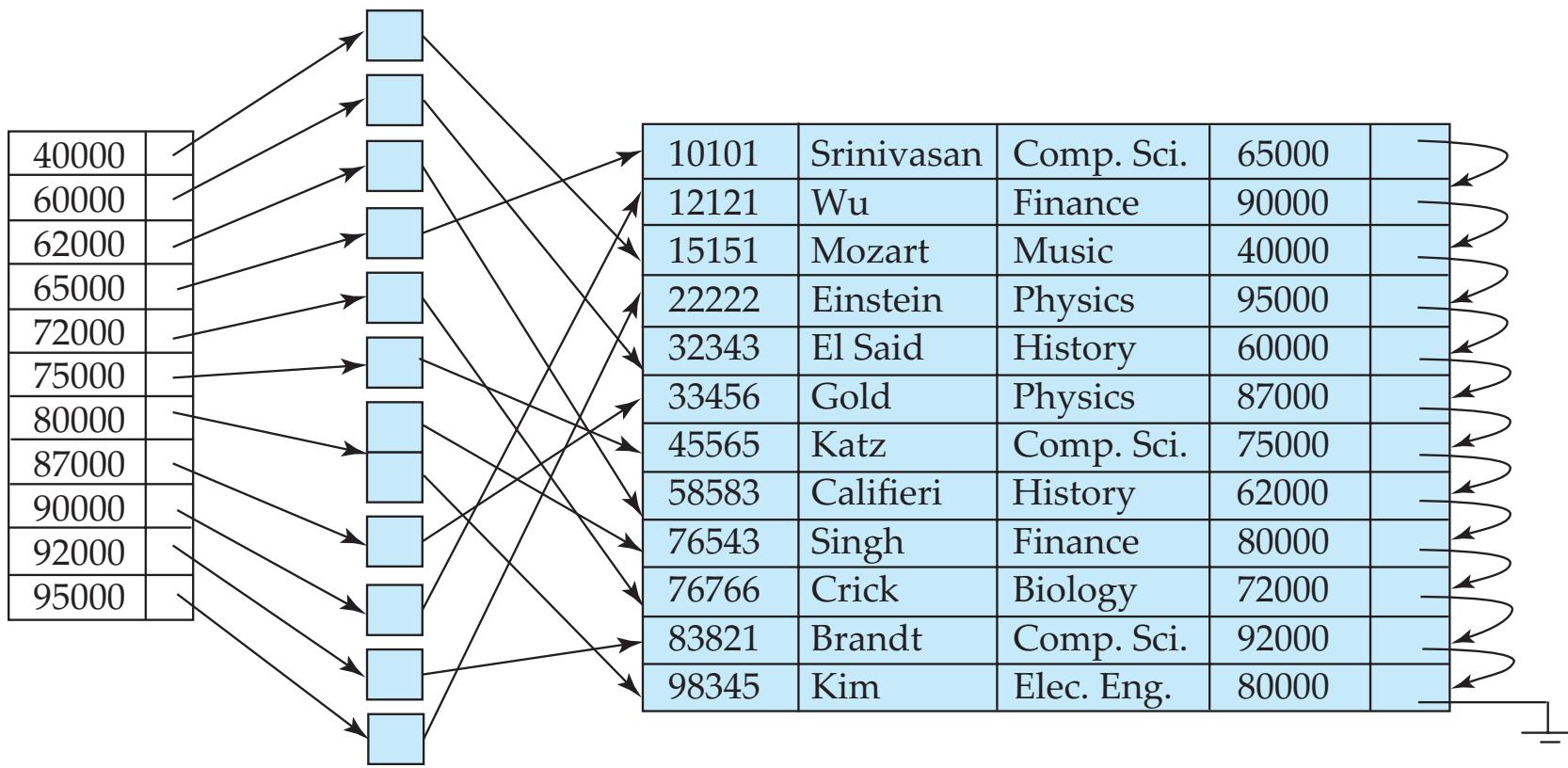
Sparse Index Files (Cont.)

- 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.





Secondary Indices Example

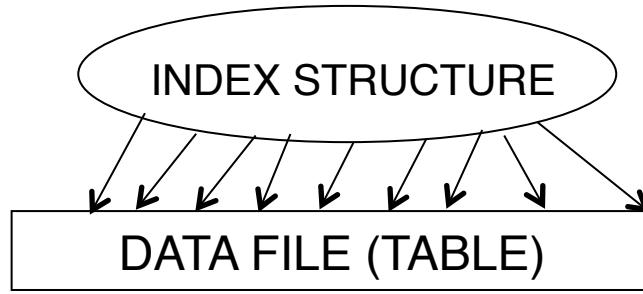


- Index record points to a bucket that contains pointers to all the actual records with that particular search-key value.
- Secondary indices have to be dense

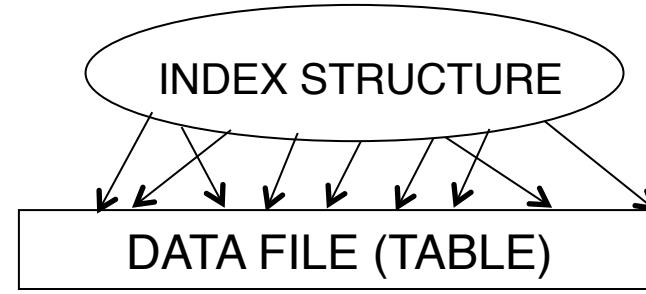


Primary and Secondary Indices

- 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



Primary Index (sorted by key)



Secondary Index (not sorted by key)

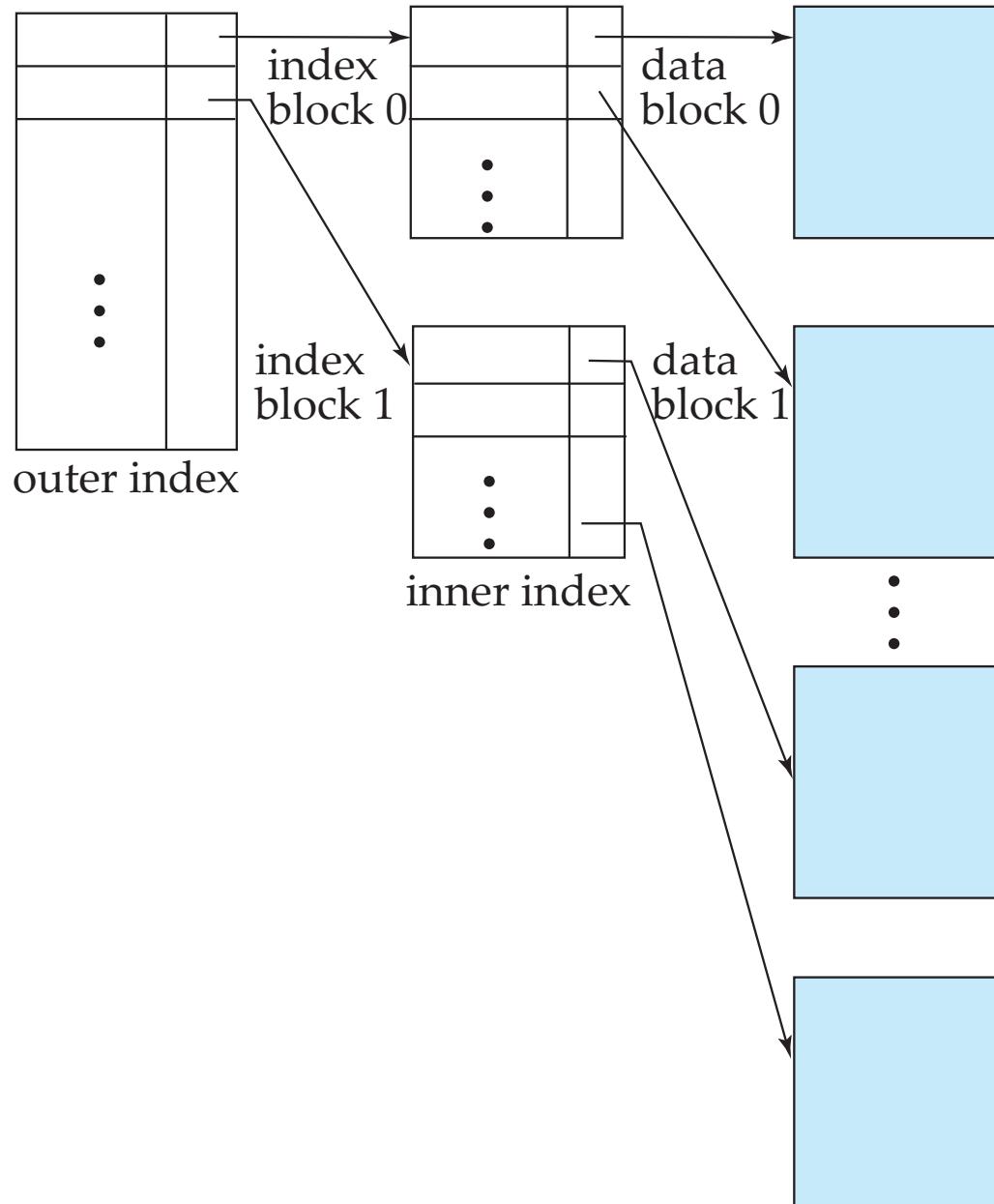


Multilevel Index

- 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.)





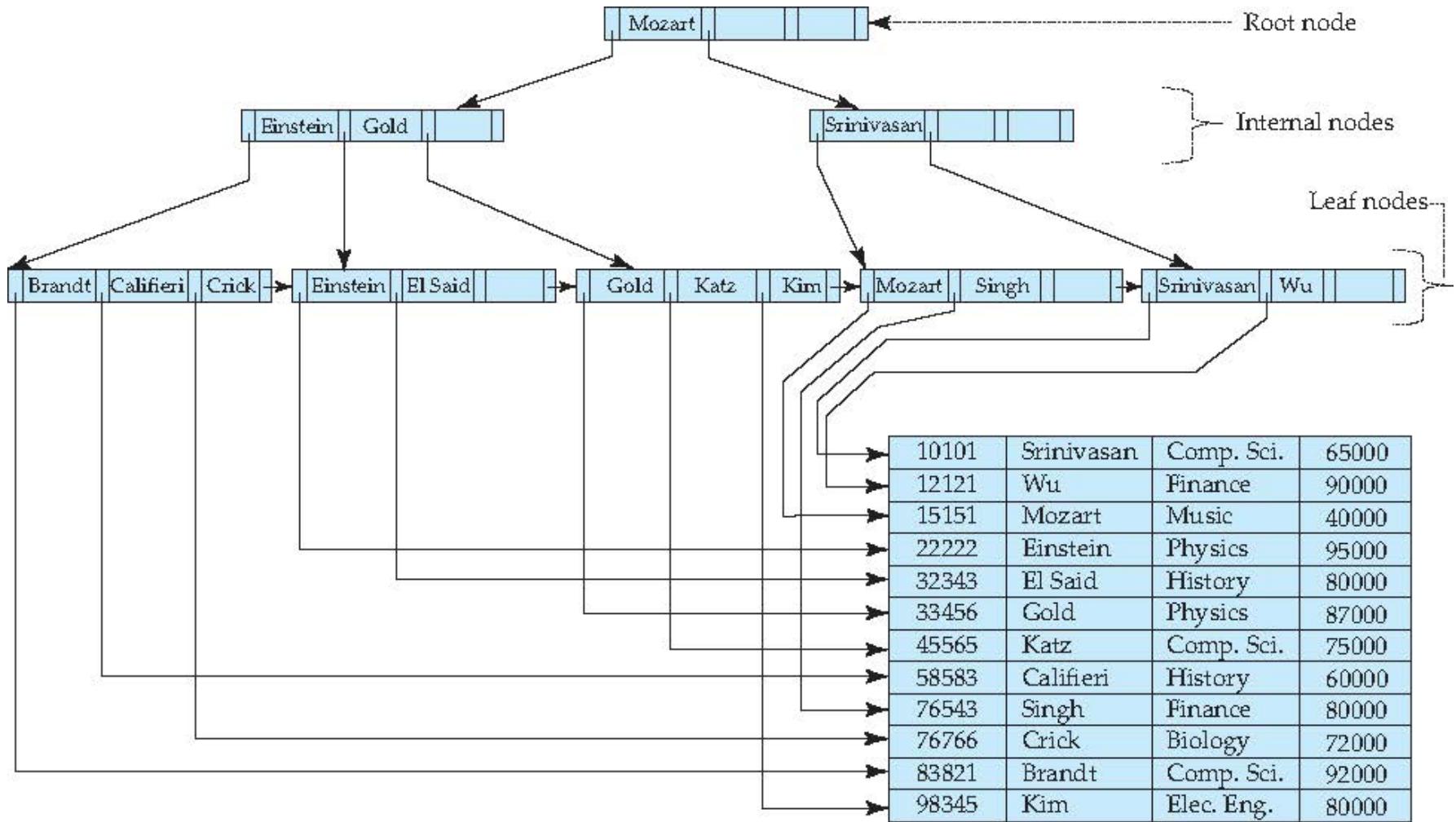
B+-Tree Index Files

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+-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	\dots	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

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

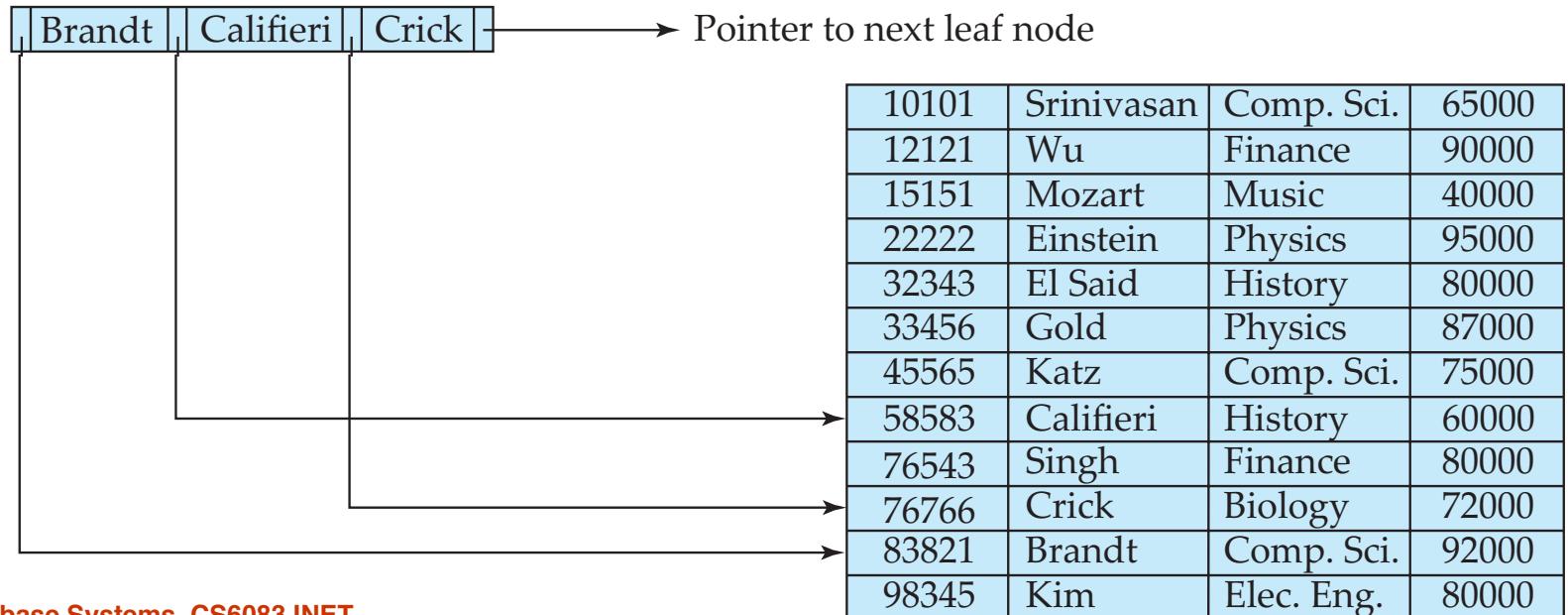


Leaf Nodes in B+-Trees

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





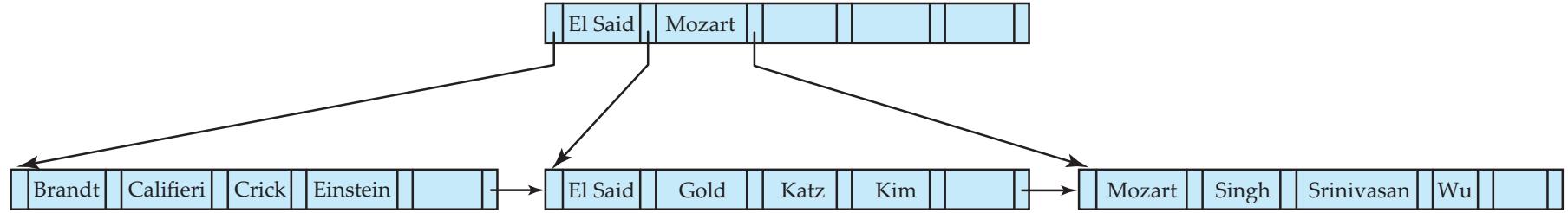
Non-Leaf Nodes in B+-Trees

- Non leaf nodes form a multi-level sparse index on the leaf nodes. For a non-leaf node with n 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	\dots	P_{n-1}	K_{n-1}	P_n
-------	-------	-------	---------	-----------	-----------	-------



Example of B⁺-tree



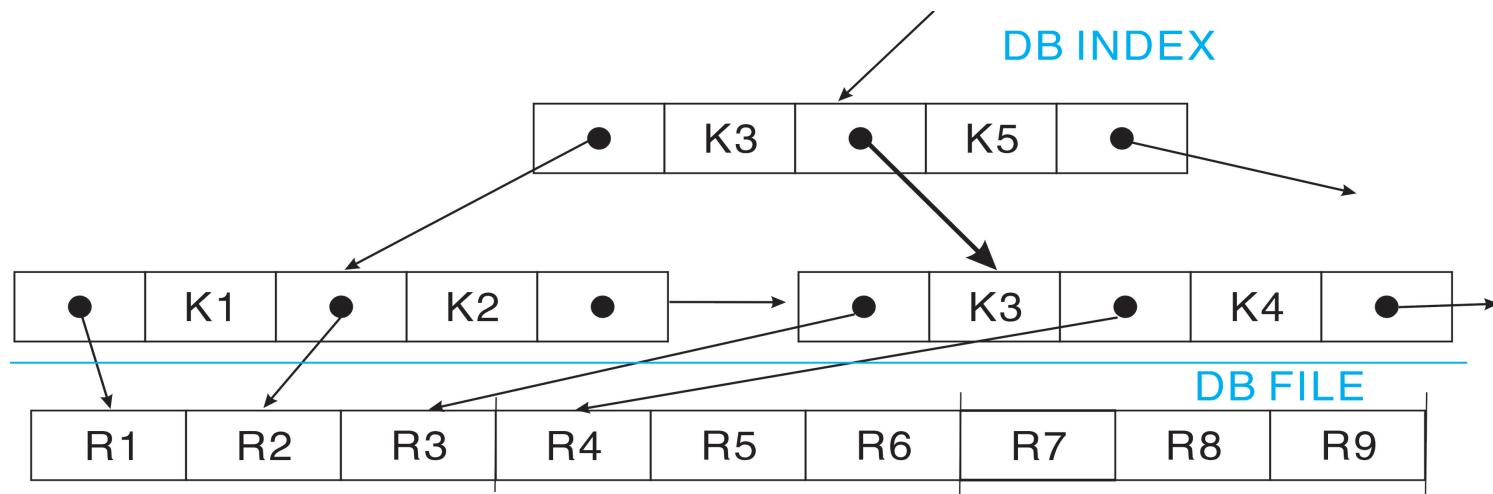
B⁺-tree for *instructor* file ($n = 6$)

- 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.



Another Example

- 4 KB node size
- 8 – byte integer keys
- 8 – byte pointers
- Choose $n = 250$ so that $250*8 + 249*8 = 3992$ bytes
 ≈ 4 KB (including some overhead)





Insertion & Deletion

■ Insert:

- If space left, done
- If not:
 - ▶ If neighbor has space, rebalance with him
 - ▶ Else, split and insert new node in next level

■ Delete

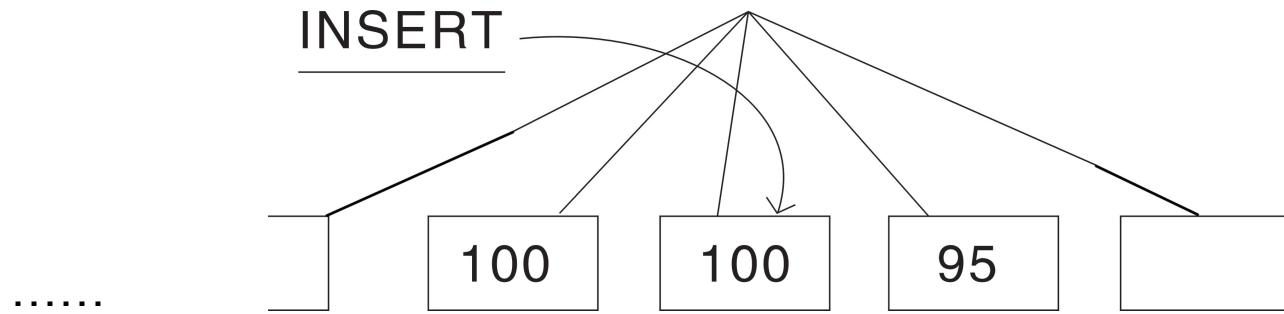
- If still $\geq \lceil n/2 \rceil$ children, done
- Else, try to balance with neighbors
- If impossible, merge with neighbor, and push delete up by one level



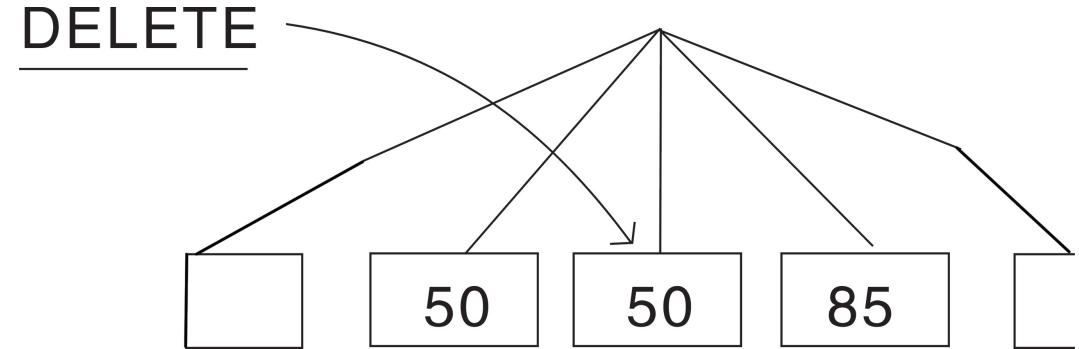
Insertion & Deletion (Cont.)

- Note: numbers are # of index entries in leaves

- $n = 101$



- Could balance with right neighbor (98 each afterwards) or not
- Delete: try to balance



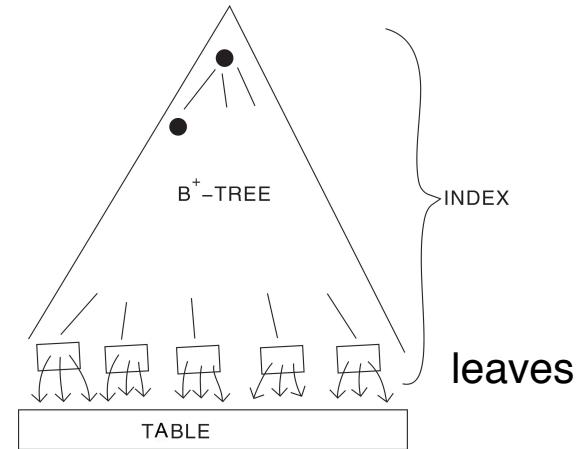
- Never go beyond immediate neighbor!



File vs. Index Organization

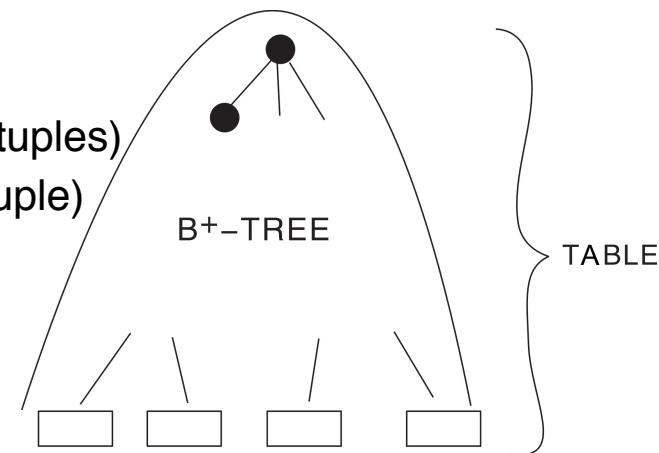
■ Example of B⁺ - Tree (Clustered) Index Organization

- Index and relation are separate



■ Example of B⁺ - Tree File Organization

- Index and relation are in the same file
- Tuples are contained in leaves (instead of pointers to tuples)
- One less indirection (one fewer disk access to find a tuple)
- But more leaf nodes (height may change)
- Can still add secondary indexes
- “The file is the Index”





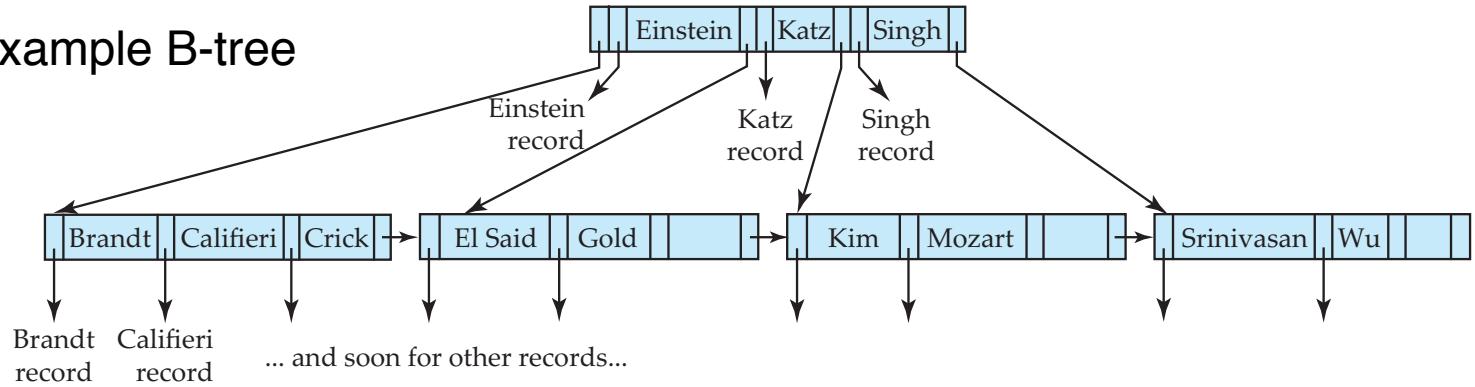
Bulk Loading and Bottom-Up Build

- Inserting entries one-at-a-time into a B⁺-tree requires ≥ 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 in earlier lecture)
 - 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
 - ▶ Build leaf nodes by scanning entries in sorted order, then next level, etc.
 - 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 B-tree; an additional pointer field for each search key in a nonleaf node must be included to store a pointer to record or buckets of records with same key.
- Example B-tree



- Advantages: less nodes than B+-trees, likely to find search-key before reaching leaf
- Disadvantages: more complicated to insert and delete, implementation is harder.



Dense Non-Clustered B+-Tree Index

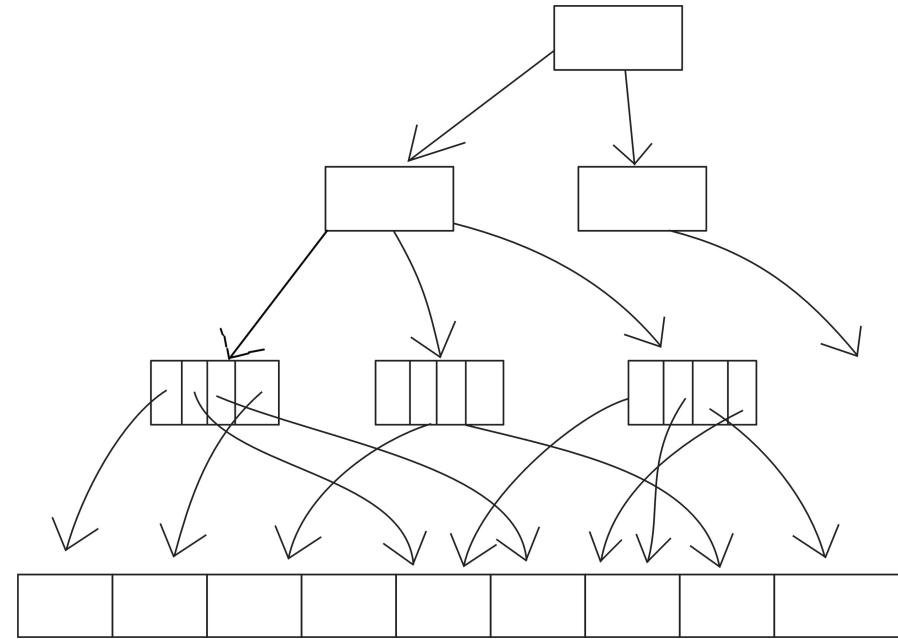
$n = 250$

1 root

16 ~ 64 nodes

4000 ~ 8000 nodes

1000000 records

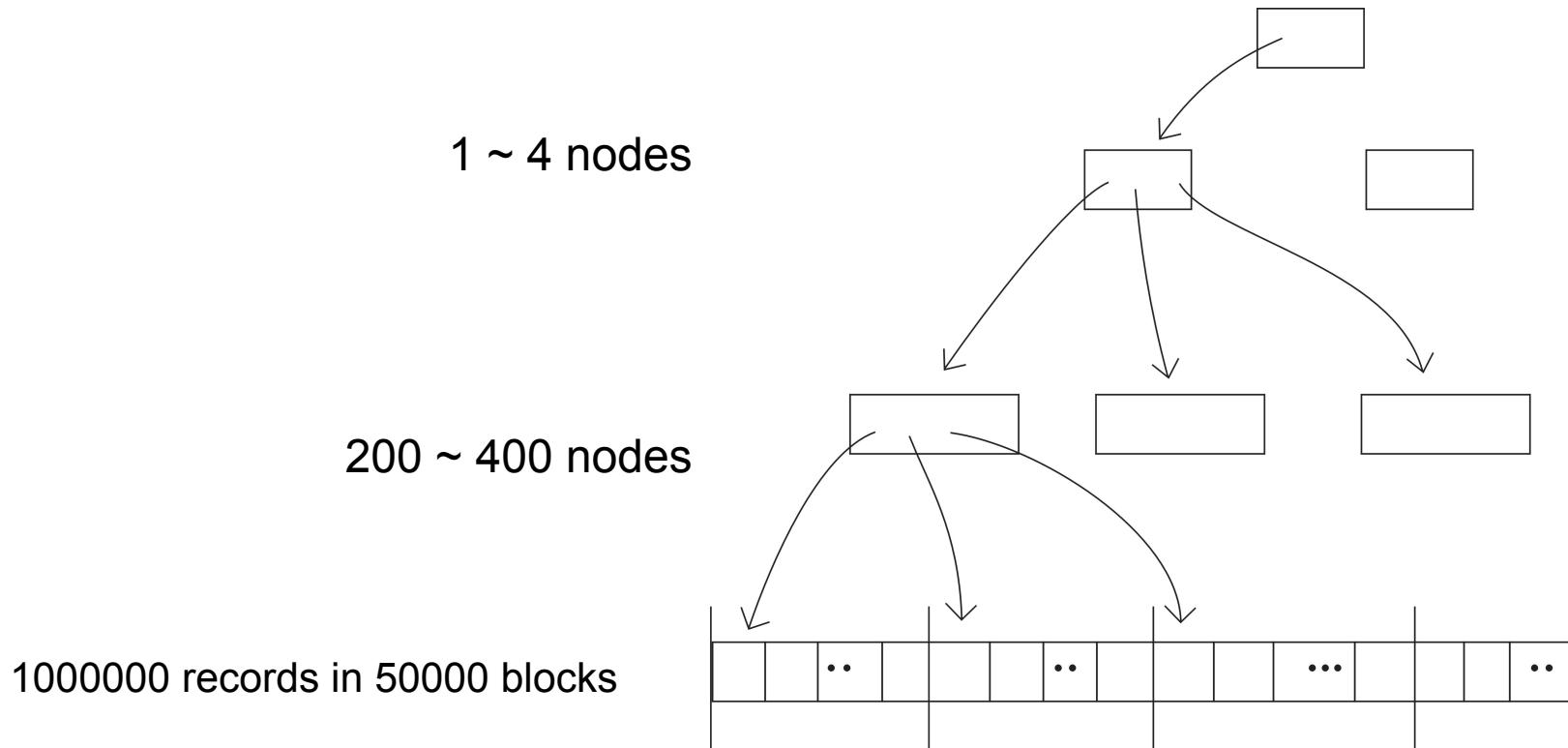


- Leaf level: 1000000 (key, pointer) pairs
- Next level: 4000 ~ 8000 pairs



Sparse, Clustered B⁺-Tree Index

- $n = 250$, 20 records/block





Analysis of B⁺ - Tree Costs

- Assume occupancy factor
- E.G. 80% \rightarrow degree 200 for $n = 250$
- Or upper / lower bounds as shown
- Also occupancy factor for DB file if sparse index
- Then count cost in terms of disk accesses and/or transfer costs



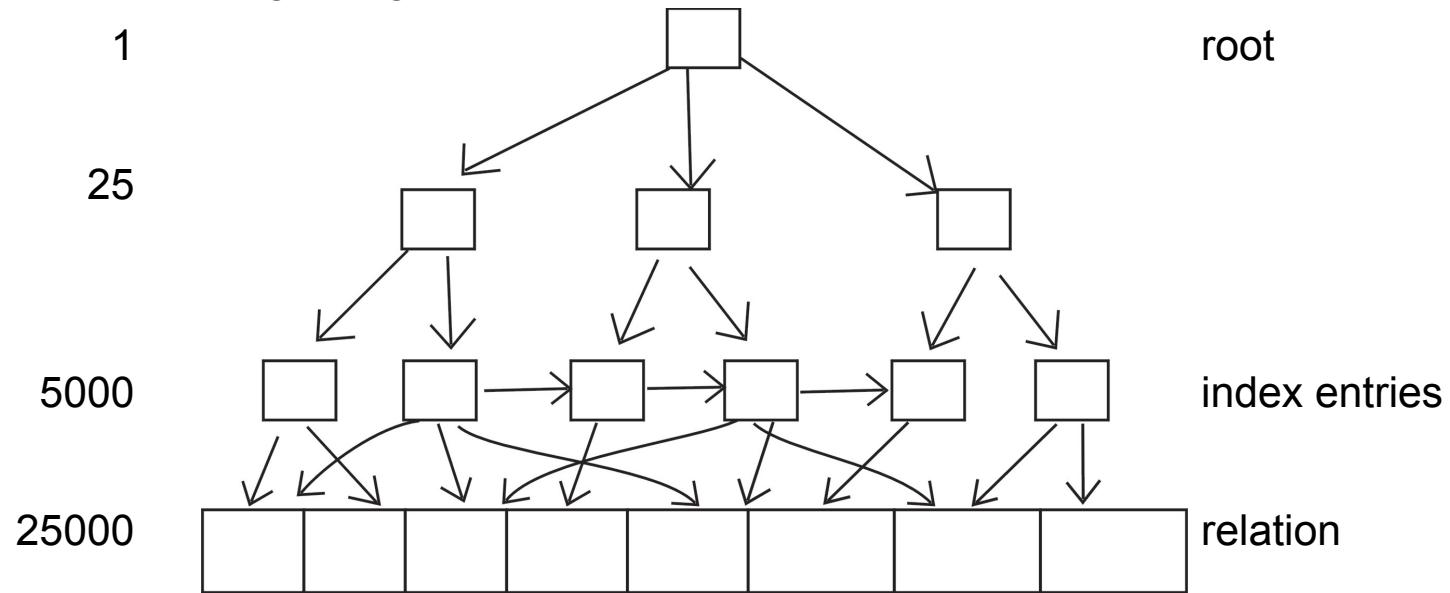
Example

- Disk with $T_s + T_r = 10\text{ms}$ and $tr = 10 \text{ MB/s}$
- Table with 1000000 records of size 100 bytes
- Assume block size 4 KB
 - 25000 blocks with 40 records each (100 % occupancy) and $T_B = 10\text{ms} + 0.4\text{ms} = 10.4\text{ms}$ (time per block)



1. Unclustered B+ Tree

- Key of size = 8 bytes and Pointer = 8 bytes
- $n = 250$ (approximation)
- Assume average degree 200 (80%)

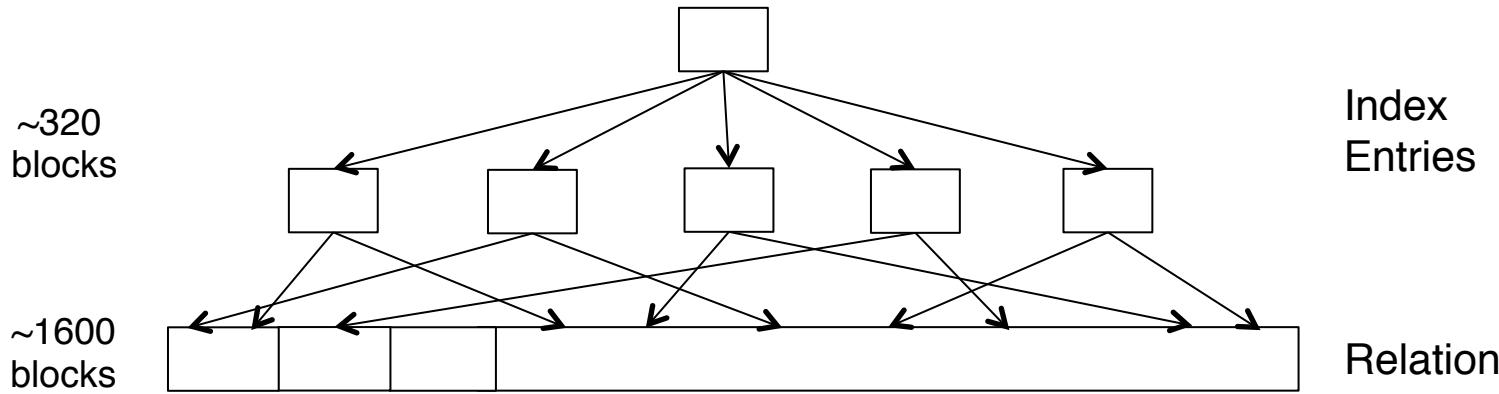


- Find record: $4 * 10.4\text{ms} = 41.6\text{ms}$
- Range Query for 1000 records: $1000 * 10.4\text{ms} + 5 * 10.4\text{ms} + 2 * 10.4\text{ms} \approx 10.5\text{s}$



2. Same with block size 64 KB

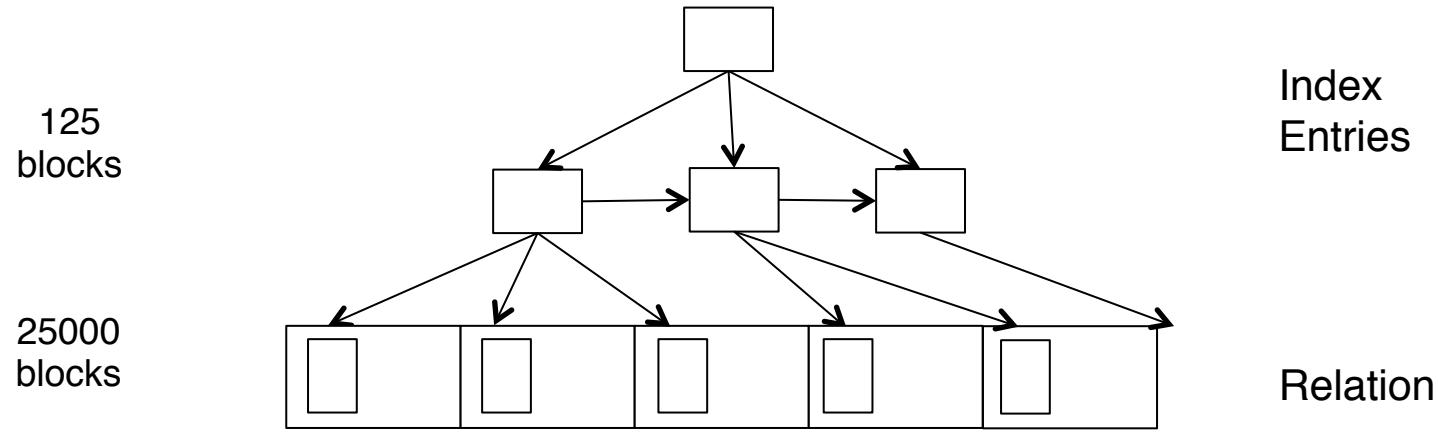
- ~ 1600 blocks with 640 records each
- $n = 4000$ and $T_B = 16.4\text{ms}$
- Assume average degree 3200



- One item: $3 * 16.4\text{ms} = 49.2\text{ms}$
- Retrieve 1000: $1000 * 16.4\text{ms} + 2 * 16.4\text{ms} = 16.5\text{s}$ ($n = 4000$ not realistic)



3. Clustered B+-Tree, Sparse, 4KB blocks



- One record $\rightarrow 3 * 10.4\text{ms} = 31.2\text{ms}$
- 1000 consecutive records $\rightarrow (25+2) * 10.4 = 280\text{ms}$



4. Clustered, sparse, 64K Blocks

- 64 KB/page → 16.4ms/page
- 1 record → 49.2ms (still 3 levels)
- 1000 records → $4 * 16.4 = 65.6\text{ms}$



5. Scanning a relation

- 4KB page, 25000 pages of data
- $25000 * 10.4\text{ms} \approx 280\text{s}$ (4KB block model)
- -----
- 64KB page, 1600 page of data
- $1600 * 16.4\text{ms} \approx 26\text{s}$ (64KB block model)
- -----
- Scan model
- $10\text{ms} + 10 \text{ sec} = 10.01 \text{ sec}$ (To scan 100MB file)



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 is 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



Static Hashing

- A **bucket** is a unit of storage containing one or more records (a bucket is typically a disk block).
- In a **hash file organization** we obtain the bucket of a record directly from its search-key value using a **hash function**.
- Hash function h is a function from the set of all search-key values K to the set of all bucket addresses B .
- Hash function is used to locate records for access, insertion as well as deletion.
- Records with different search-key values may be mapped to the same bucket; thus entire bucket has to be searched sequentially to locate a record.



Example of Hash File Organization

Hash file organization of *instructor* file, using *dept_name* as key

- There are 10 buckets,
- The binary representation of the i th character is assumed to be the integer i .
- The hash function returns the sum of the binary representations of the characters modulo 10

$$h(\text{Music}) = 1$$

$$h(\text{History}) = 2$$

$$h(\text{Physics}) = 3$$

$$h(\text{Elec. Eng.}) = 3$$

bucket 0			

bucket 1			
15151	Mozart	Music	40000

bucket 2			
32343	El Said	History	80000
58583	Califieri	History	60000

bucket 3			
22222	Einstein	Physics	95000
33456	Gold	Physics	87000
98345	Kim	Elec. Eng.	80000

bucket 4			
12121	Wu	Finance	90000
76543	Singh	Finance	80000

bucket 5			
76766	Crick	Biology	72000

bucket 6			
10101	Srinivasan	Comp. Sci.	65000
45565	Katz	Comp. Sci.	75000
83821	Brandt	Comp. Sci.	92000

bucket 7			



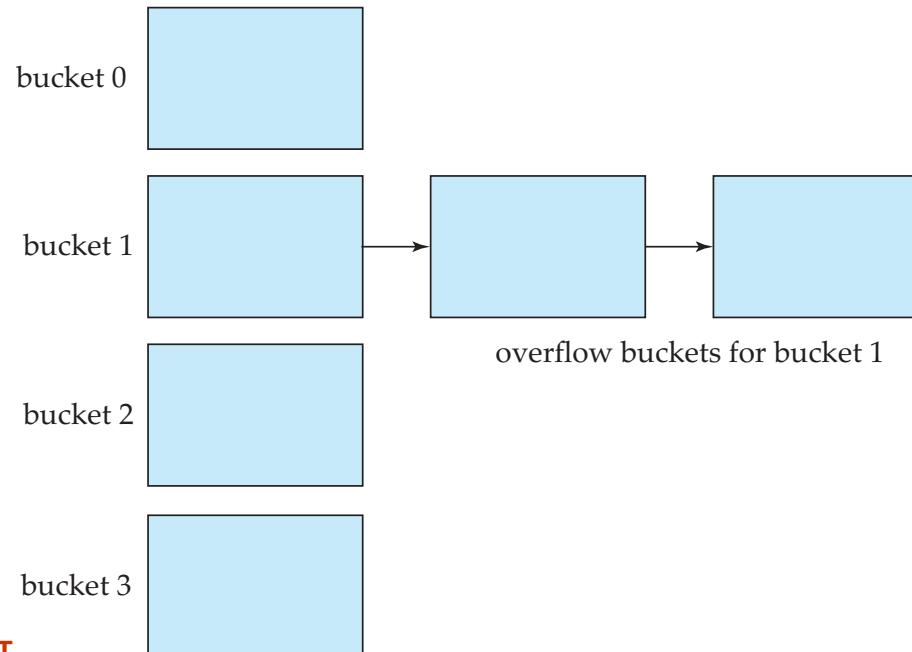
Hash Functions

- Worst hash function maps all search-key values to the same bucket; this makes access time proportional to the number of search-key values in the file.
- An ideal hash function is **uniform**, i.e., each bucket is assigned the same number of search-key values from the set of *all* possible values.
- Ideal hash function is **random**, so each bucket will have the same number of records assigned to it irrespective of the *actual distribution* of search-key values in the file.
- Typical hash functions perform computation on the internal binary representation of the search-key.
 - For example, for a string search-key, the binary representations of all the characters in the string could be added and the sum modulo the number of buckets could be returned.



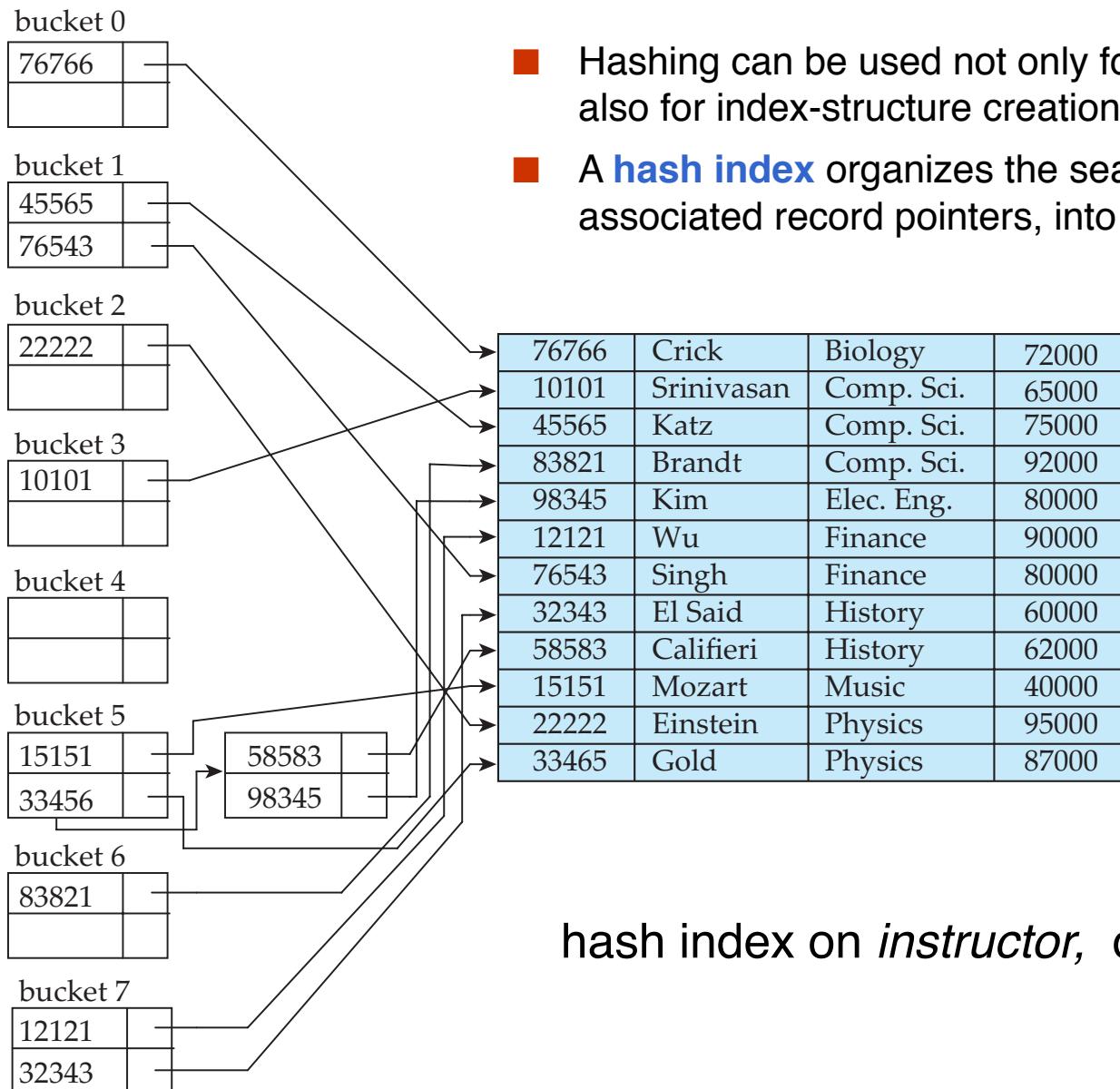
Handling of Bucket Overflows

- Bucket overflow can occur because of
 - Insufficient buckets
 - Skew in distribution of records. This can occur due to two reasons:
 - ▶ multiple records have same search-key value
 - ▶ chosen hash function produces non-uniform distribution of key values
- Although the probability of bucket overflow can be reduced, it cannot be eliminated; it is handled by using ***overflow buckets***.
- **Overflow chaining** – the overflow buckets of a given bucket are chained together in a linked list in a scheme called ***closed hashing***.





Hash Indices





Deficiencies of Static Hashing

- In static hashing, function h maps search-key values to a fixed set of B of bucket addresses. Databases grow or shrink with time.
 - If initial number of buckets is too small, and file grows, performance will degrade due to too much overflows.
 - If space is allocated for anticipated growth, a significant amount of space will be wasted initially (and buckets will be underfull).
 - If database shrinks, again space will be wasted.
- One solution: periodic re-organization of the file with a new hash function
 - Expensive, disrupts normal operations
- Better solution: allow the number of buckets to be modified dynamically.

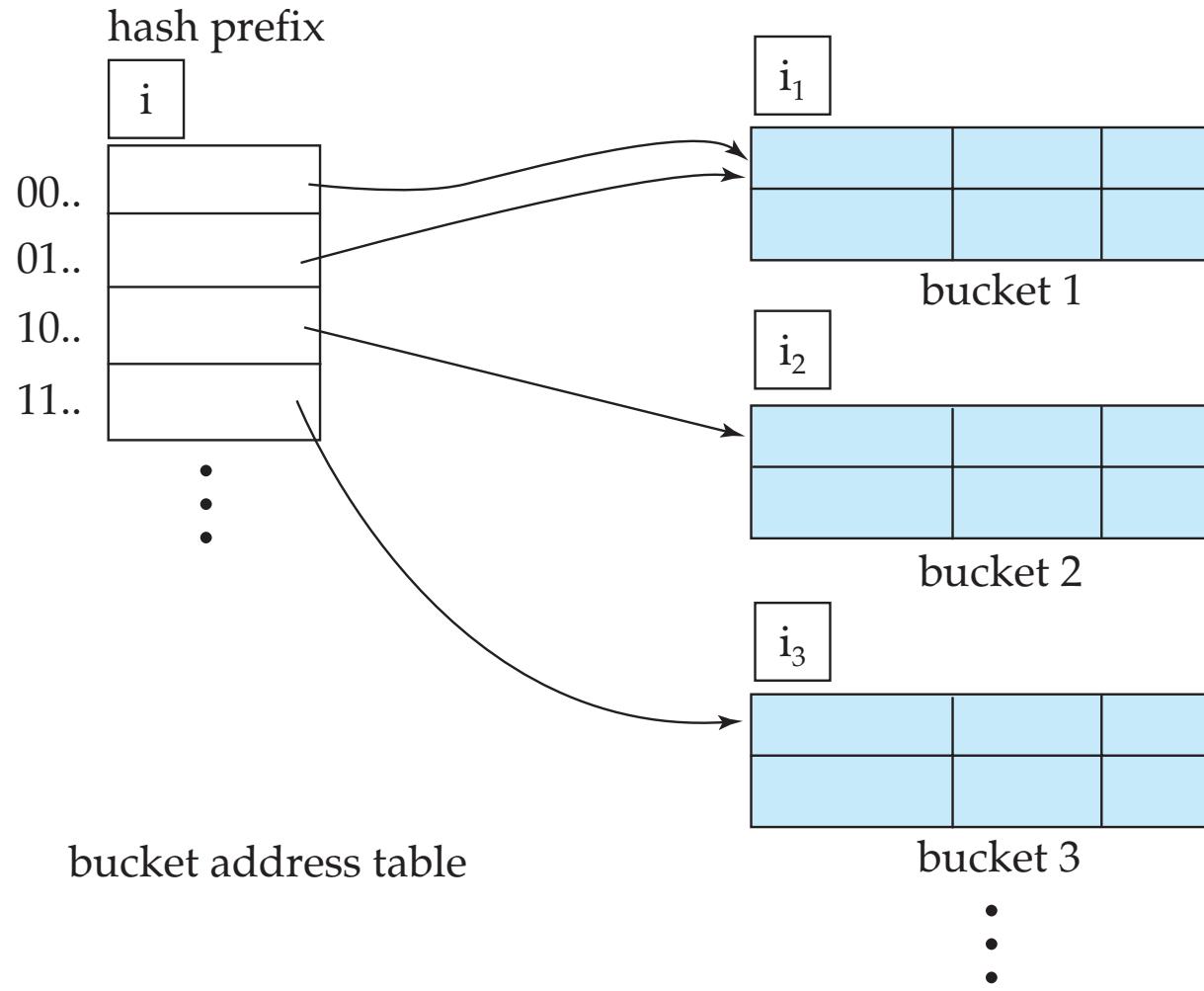


Dynamic Hashing

- Good for database that grows and shrinks in size
- Allows the hash function to be modified dynamically
- **Extendable hashing** – one form of dynamic hashing
 - Hash function generates values over a large range – typically b -bit integers, with $b = 32$.
 - At any time use only a prefix of the hash function to index into a table of bucket addresses.
 - Let the length of the prefix be i bits, $0 \leq i \leq 32$.
 - ▶ Bucket address table size = 2^i . Initially $i = 0$
 - ▶ Value of i grows and shrinks as the size of the database grows and shrinks.
 - Multiple entries in the bucket address table may point to a bucket (why?)
 - Thus, actual number of buckets is $< 2^i$
 - ▶ The number of buckets also changes dynamically due to coalescing and splitting of buckets.



General Extendable Hash Structure



In this structure, $i_2 = i_3 = i$, whereas $i_1 = i - 1$



Use of Extendable Hash Structure

- Each bucket j stores a value i_j
 - All the entries that point to the same bucket have the same values on the first i_j bits.
- To locate the bucket containing search-key K_j :
 1. Compute $h(K_j) = X$
 2. Use the first i high order bits of X as a displacement into bucket address table, and follow the pointer to appropriate bucket
- To insert a record with search-key value K_j
 - follow same procedure as look-up and locate the bucket, say j .
 - If there is room in the bucket j insert record in the bucket.
 - Else the bucket must be split and insertion re-attempted (next slide.)
 - ▶ Overflow buckets used instead in some cases (will see shortly)



Insertion in Extendable Hash Structure (Cont)

To split a bucket j when inserting record with search-key value K_j :

- If $i > i_j$ (more than one pointer to bucket j)
 - allocate a new bucket z , and set $i_j = i_z = (i_j + 1)$
 - Update the second half of the bucket address table entries originally pointing to j , to point to z
 - remove each record in bucket j and reinsert (in j or z)
 - recompute new hash for K_j and insert record in the bucket (further splitting is required if the bucket is still full)
- If $i = i_j$ (only one pointer to bucket j)
 - If i reaches some limit b , or too many splits have happened in this insertion, create an overflow bucket
 - Else
 - ▶ increment i and double the size of the bucket address table.
 - ▶ replace each entry in the table by two entries that point to the same bucket.
 - ▶ recompute new bucket address table entry for K_j
Now $i > i_j$ so use the first case above.



Deletion in Extendable Hash Structure

- To delete a key value,
 - locate it in its bucket and remove it.
 - The bucket itself can be removed if it becomes empty (with appropriate updates to the bucket address table).
 - Coalescing of buckets can be done (can coalesce only with a “*buddy*” bucket having same value of i_j and same $i_j - 1$ prefix, if it is present)
 - Decreasing bucket address table size is also possible
 - ▶ Note: decreasing bucket address table size is an expensive operation and should be done only if number of buckets becomes much smaller than the size of the table

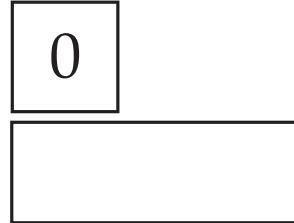


Use of Extendable Hash Structure: Example

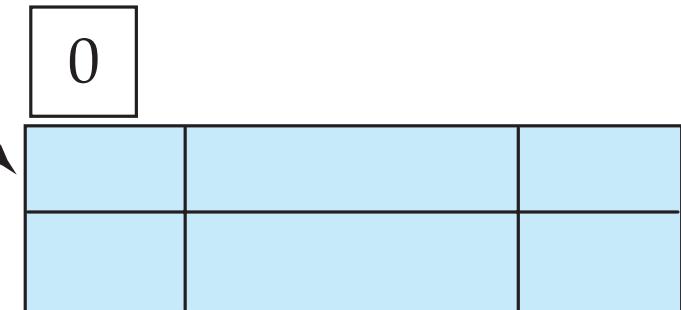
$dept_name$	$h(dept_name)$
Biology	0010 1101 1111 1011 0010 1100 0011 0000
Comp. Sci.	1111 0001 0010 0100 1001 0011 0110 1101
Elec. Eng.	0100 0011 1010 1100 1100 0110 1101 1111
Finance	1010 0011 1010 0000 1100 0110 1001 1111
History	1100 0111 1110 1101 1011 1111 0011 1010
Music	0011 0101 1010 0110 1100 1001 1110 1011
Physics	1001 1000 0011 1111 1001 1100 0000 0001

- Initial Hash structure; bucket size = 2

hash prefix



bucket address table

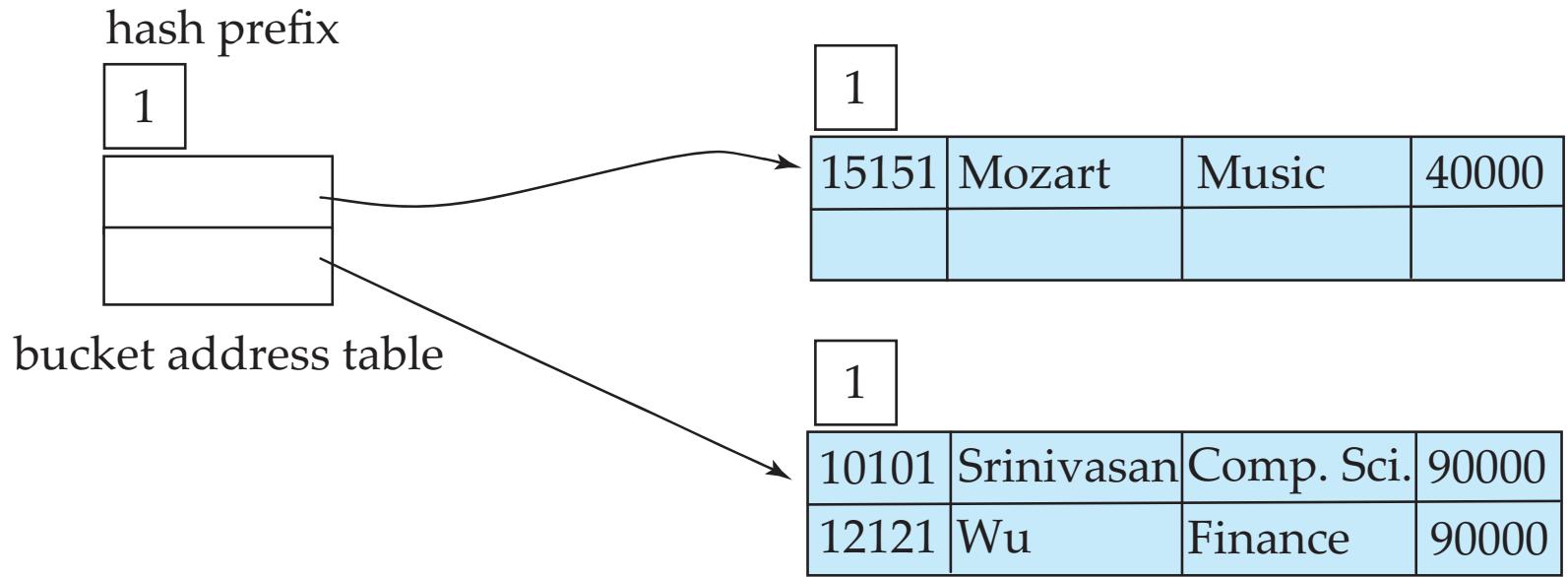


bucket 1



Example (Cont.)

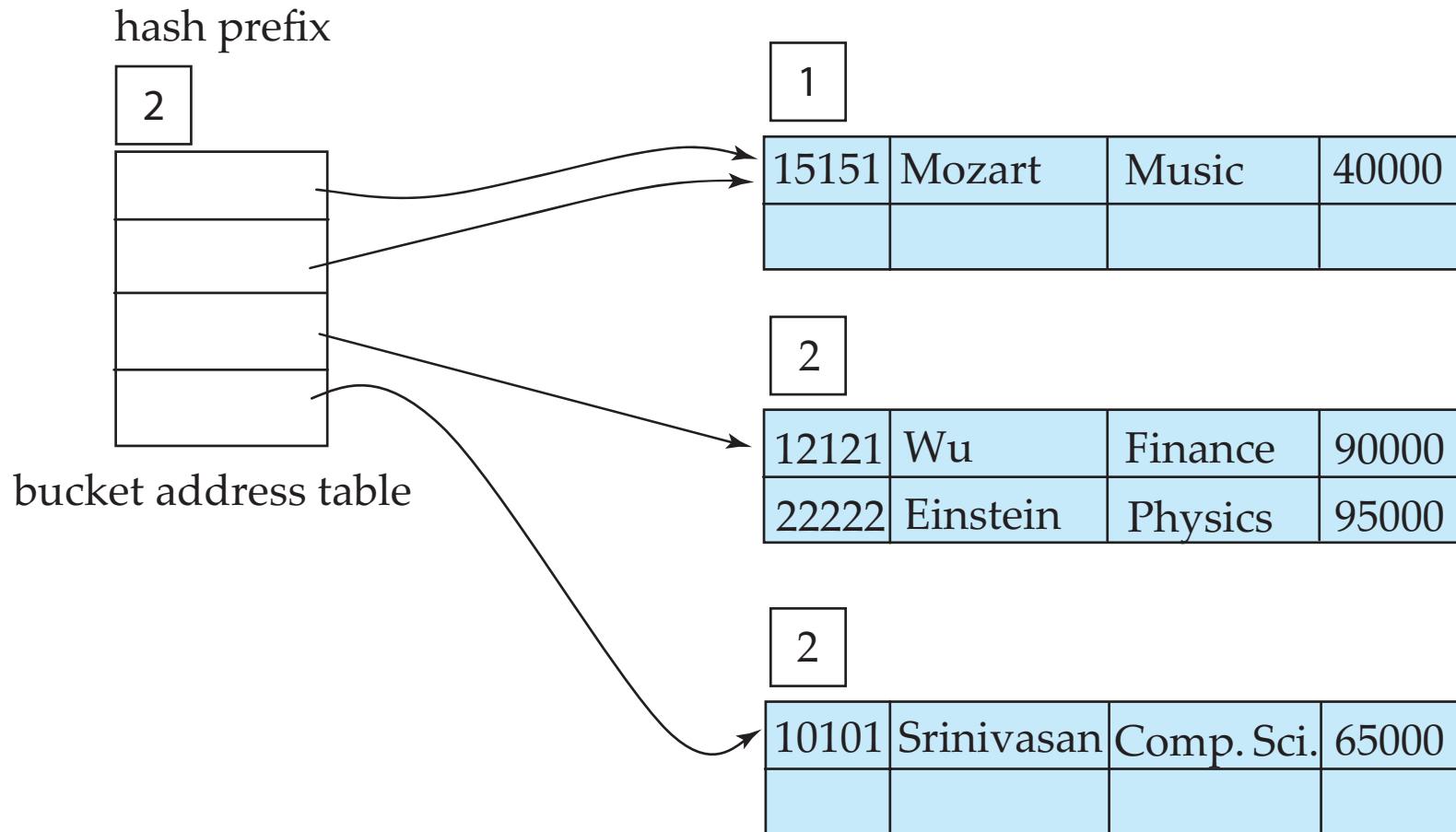
- Hash structure after insertion of “Mozart”, “Srinivasan”, and “Wu” records





Example (Cont.)

■ Hash structure after insertion of Einstein record





Extendable Hashing vs. Other Schemes

- Benefits of extendable hashing:
 - Hash performance does not degrade with growth of file
 - Minimal space overhead
- Disadvantages of extendable hashing
 - Extra level of indirection to find desired record
 - Bucket address table may itself become very big (larger than memory)
 - ▶ Cannot allocate very large contiguous areas on disk either
 - ▶ Solution: B⁺-tree structure to locate desired record in bucket address table
 - Changing size of bucket address table is an expensive operation



Comparison of Ordered Indexing and Hashing

- Cost of periodic re-organization
- Relative frequency of insertions and deletions
- Is it desirable to optimize average access time at the expense of worst-case access time?
- Expected type of queries:
 - Hashing is generally better at retrieving records having a specified value of the key.
 - If range queries are common, ordered indices are to be preferred
- In practice:
 - PostgreSQL supports hash indices, but discourages use due to poor performance
 - Oracle supports static hash organization, but not hash indices
 - SQLServer supports only B+-trees



Index Definition in SQL

- Create an index

```
create index <index-name> on <relation-name>  
<attribute-list>
```

E.g.: **create index b-index on instructor(dept_name)**

- Use **create unique index** to indirectly specify and enforce the condition that the search key is a candidate key.

- Not really required if SQL **unique** integrity constraint is supported

- To drop an index

```
drop index <index-name>
```

- Most database systems allow specification of type of index, and clustering.