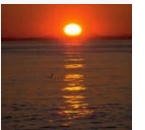




Chapter 12: Indexing and Hashing

Database System Concepts, 5th Ed.

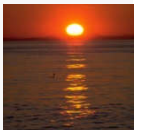
©Silberschatz, Korth and Sudarshan
See www.db-book.com for conditions on re-use





Chapter 12: 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
- Multiple-Key Access



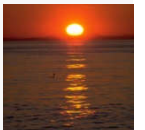


Basic Concepts

- 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
------------	---------

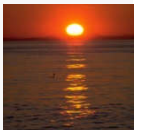
- Index files are typically much smaller than the original file
- Two basic kinds of indices:
 - **Ordered indices:** search keys are stored in sorted order
 - **Hash indices:** search keys are distributed uniformly across “buckets” using a “hash function”.





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

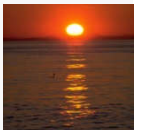




Ordered Indices

Indexing techniques evaluated on basis of:

- 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.

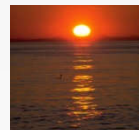




Dense Index Files

- **Dense index** — Index record appears for every search-key value in the file.

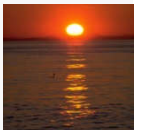
Brighton		→	A-217	Brighton	750	
Downtown		→	A-101	Downtown	500	
Mianus		→	A-110	Downtown	600	
Perryridge		→	A-215	Mianus	700	
Redwood		→	A-102	Perryridge	400	
Round Hill		→	A-201	Perryridge	900	
		→	A-218	Perryridge	700	
		→	A-222	Redwood	700	
		→	A-305	Round Hill	350	





Sparse Index Files

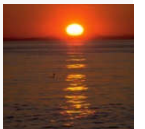
- **Sparse Index:** contains index records for only some search-key values.
 - Applicable when records are sequentially ordered on search-key
- To locate a record with search-key value K we:
 - Find index record with largest search-key value $< K$
 - Search file sequentially starting at the record to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





Example of Sparse Index Files

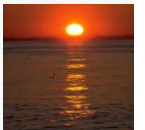
Brighton		A-217	Brighton	750	
Mianus		A-101	Downtown	500	
Redwood		A-110	Downtown	600	
		A-215	Mianus	700	
		A-102	Perryridge	400	
		A-201	Perryridge	900	
		A-218	Perryridge	700	
		A-222	Redwood	700	
		A-305	Round Hill	350	





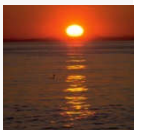
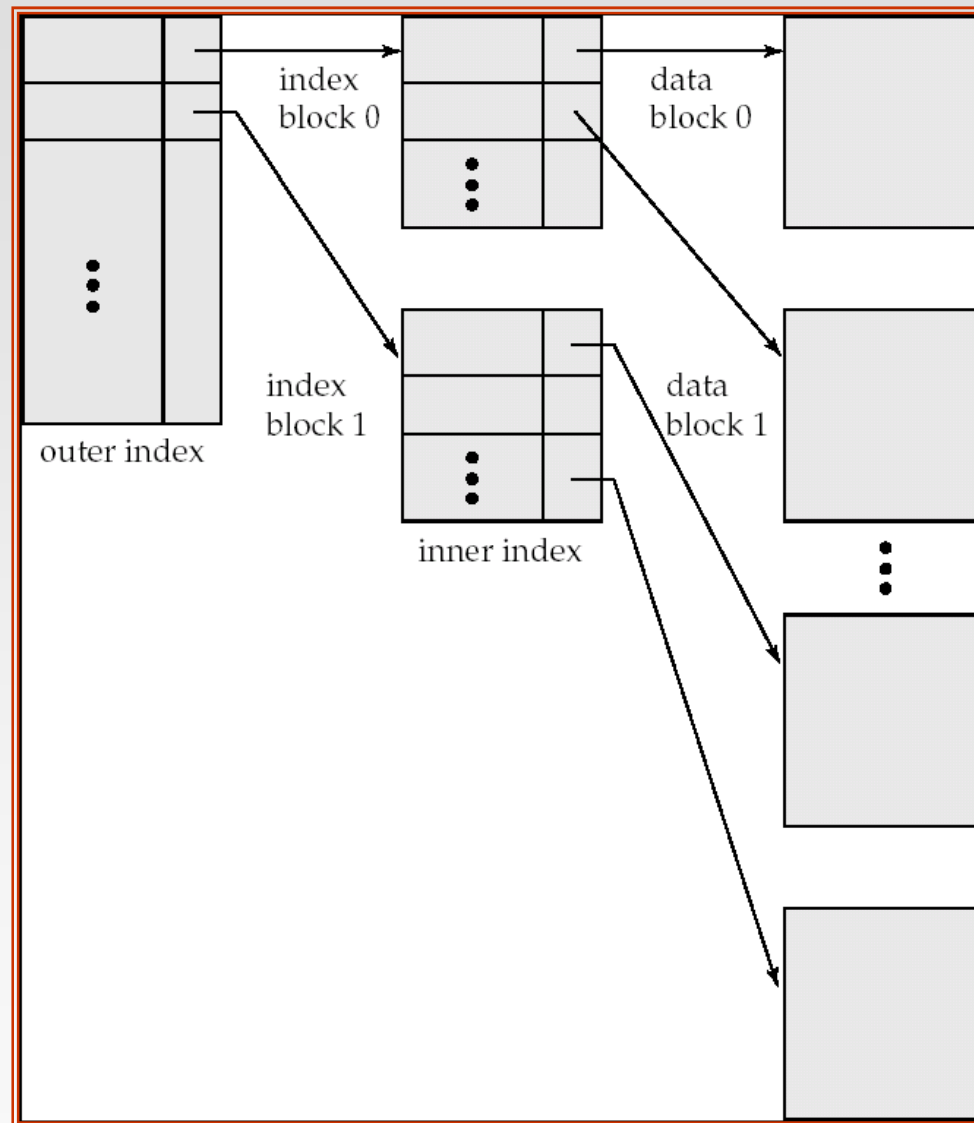
Multilevel Index

- If primary index does not fit in memory, access becomes expensive.
- To reduce number of disk accesses to index records, 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.





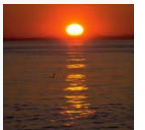
Multilevel Index (Cont.)





Index Update: Deletion

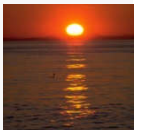
- If deleted record was the only record in the file with its particular search-key value, the search-key is deleted from the index also.
- Single-level index deletion:
 - Dense indices – deletion of search-key is similar to file record deletion.
 - Sparse indices –
 - ▶ if an entry for the search key 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

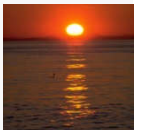
- Single-level index insertion:
 - Perform a lookup using the search-key value appearing in the record to be inserted.
 - Dense indices – if the search-key value does not appear in the index, insert it.
 - Sparse indices – if index stores an entry for each block of the file, no change needs to be made 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 (as well as deletion) algorithms are simple extensions of the single-level algorithms





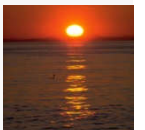
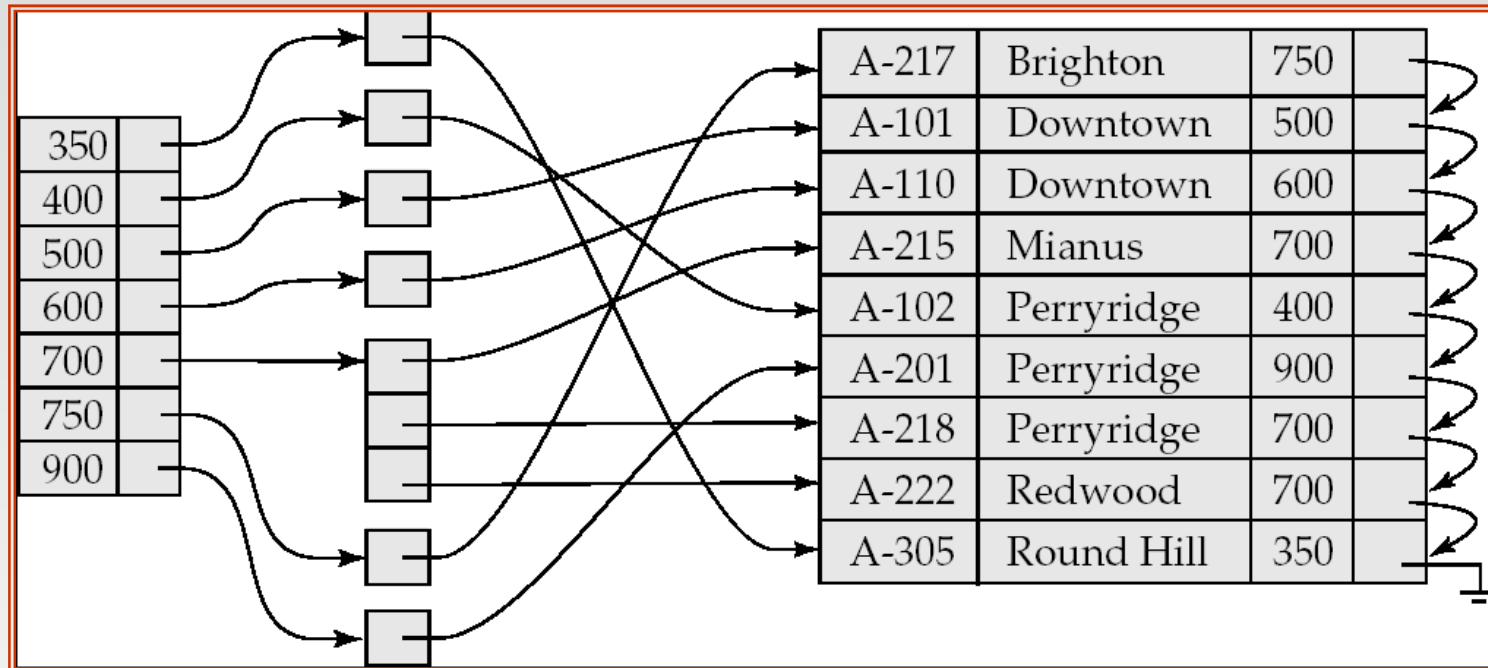
Secondary Indices

- 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 *account* relation stored sequentially by account number, we may want to find all accounts in a particular branch
 - Example 2: as above, but where we want to find all accounts with a specified balance or range of balances
- We can have a secondary index with an index record for each search-key value
 - index record points to a bucket that contains pointers to all the actual records with that particular search-key value.





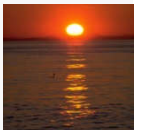
Secondary Index on *balance* field of *account*





Primary and Secondary Indices

- Secondary indices have to be dense.
- Indices offer substantial benefits when searching for records.
- When a file is modified, every index on the file must be updated, Updating indices imposes overhead on database modification.
- 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

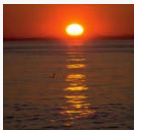




B⁺-Tree Index Files

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.

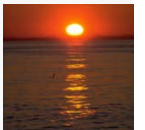




B⁺-Tree Index Files (Cont.)

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.



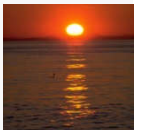


B⁺-Tree Node Structure

- Typical node



- K_i are the search-key values
- The search-keys in a node are ordered
$$K_1 < K_2 < K_3 < \dots < K_{n-1}$$
- Leaf node: P_i ($i=1, \dots, n-1$) are pointers to records or buckets of records.
- Non-leaf nodes: P_i are pointers to children.

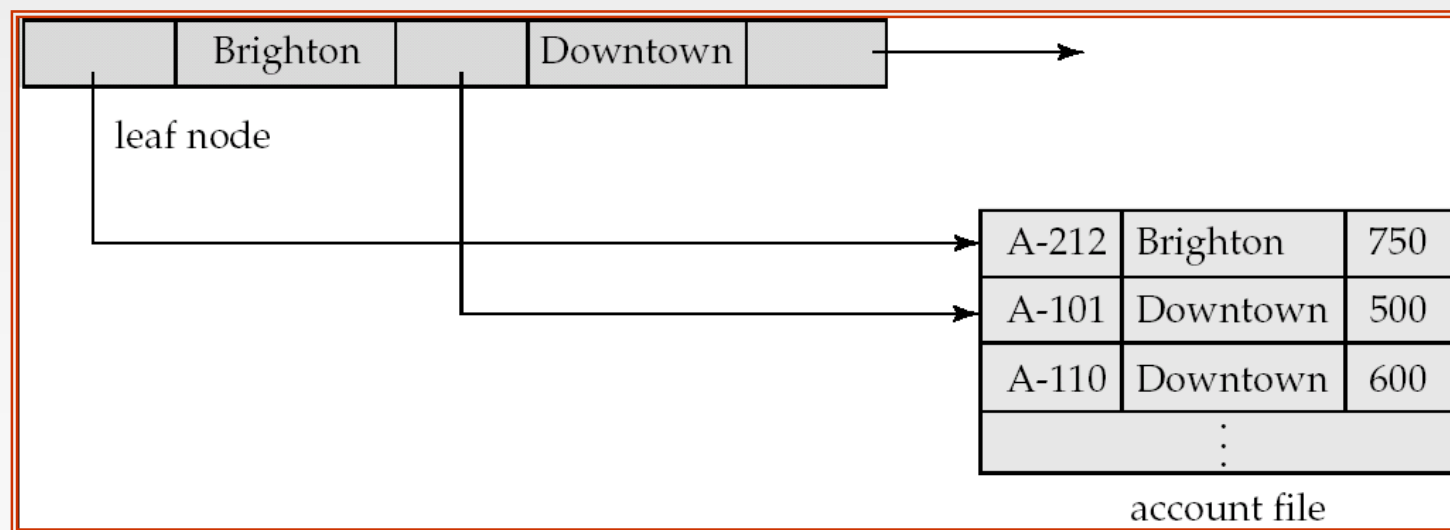




Leaf Nodes in B⁺-Trees

Properties of a leaf node:

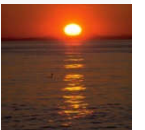
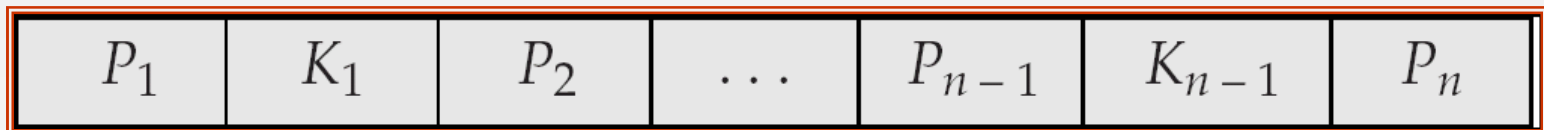
- For $i = 1, 2, \dots, n-1$, pointer P_i either points to a file record with search-key value K_i , or to a bucket of pointers to file records, each record having search-key value K_i . Only need bucket structure if search-key does not form a primary key.
- If L_i, L_j are leaf nodes and $i < j$, L_i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





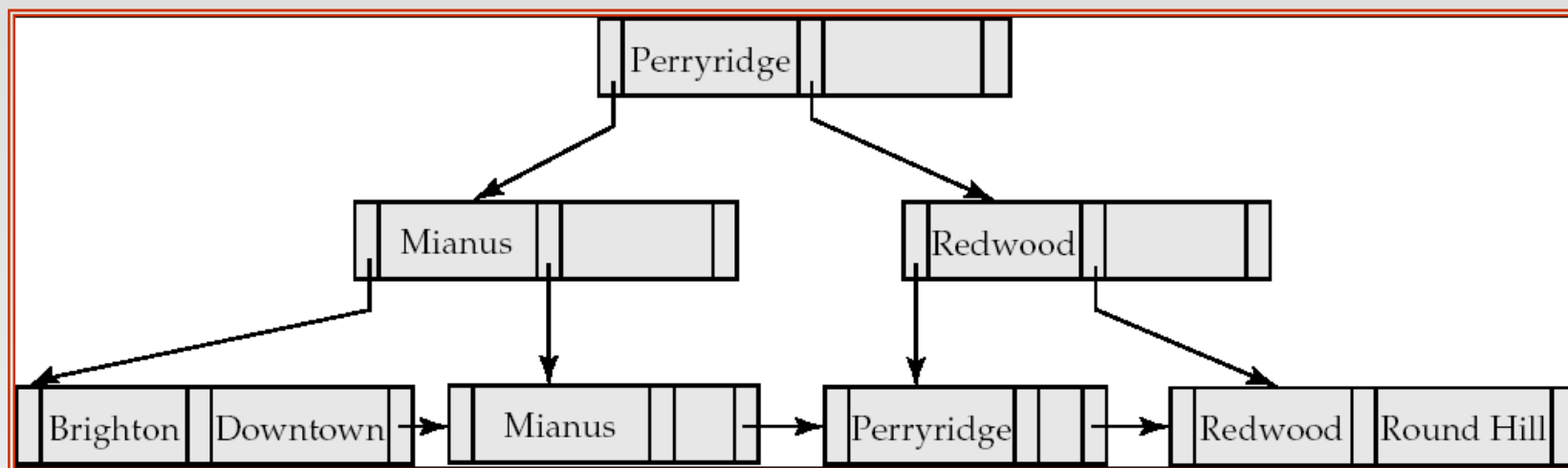
Non-Leaf Nodes in B⁺-Trees

- 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
 - Pointer P_n points to the part of the subtree that contains those key values which are greater than or equal to K_{n-1} ,





Example of a B⁺-tree

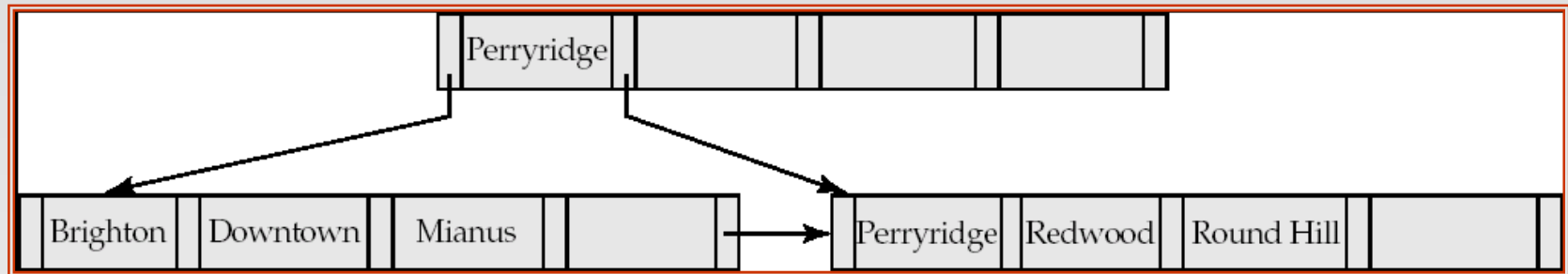


B⁺-tree for *account* file ($n = 3$)



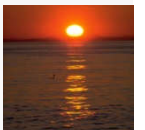


Example of B⁺-tree



B⁺-tree for *account* file ($n = 5$)

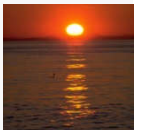
- ❑ Leaf nodes must have between 2 and 4 values ($\lceil (n-1)/2 \rceil$ and $n-1$, with $n = 5$).
- ❑ Non-leaf nodes other than root must have between 3 and 5 children ($\lceil n/2 \rceil$ and n with $n = 5$).
- ❑ Root must have at least 2 children.





Observations about B⁺-trees

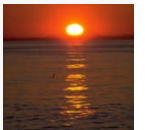
- Since the inter-node connections are done by pointers, “logically” close blocks need not be “physically” close.
- The non-leaf levels of the B⁺-tree form a hierarchy of sparse indices.
- The B⁺-tree contains a relatively small number of levels (logarithmic in the size of the main file), thus searches 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 .
 1. Start with the root node
 1. Examine the node for the smallest search-key value $> k$.
 2. If such a value exists, assume it is K_j . Then follow P_j to the child node
 3. Otherwise $k \geq K_{m-1}$, where there are m pointers in the node. Then follow P_m to the child node.
 2. If the node reached by following the pointer above is not a leaf node, repeat step 1 on the node
 3. Else we have reached a leaf node.
 1. If for some i , key $K_i = k$ follow pointer P_i to the desired record or bucket.
 2. Else no record with search-key value k exists.





procedure find(*value V*)

set C = root node

while C is not a leaf node **do**

begin

Let K_i = smallest search-key value, if any, greater than V

if there is no such value **then begin** // $V > K_{n-1}$

Let m = the number of pointers in the node

set C = node pointed to by P_m

end

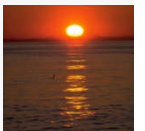
else set C = the node pointed to by P_i

end

if there is a key value K_i in C such that $K_i = V$ //leaf level

then pointer P_i directs us to the desired record or bucket

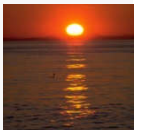
else no record with key value V exists





Queries on B⁺-Trees (Cont.)

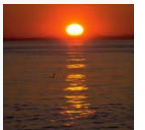
- In processing a query, a path is traversed in the tree from the root to some leaf node.
- If there are K search-key values in the file, the path is no longer than $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$.
- A node is generally the same size as a disk block, typically 4 kilobytes, and n is typically 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 20 milliseconds!





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 a pointer is inserted into the bucket.
- If the search-key value is not there, then add the record to the main file and create a bucket if necessary. Then:
 - If there is room in the leaf node, insert (key-value, pointer) pair in the leaf node
 - Otherwise, split the node (along with the new (key-value, pointer) entry) as discussed in the next slide.



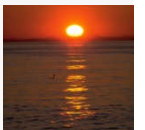


Updates on B⁺-Trees: Insertion (Cont.)

- Splitting a 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.

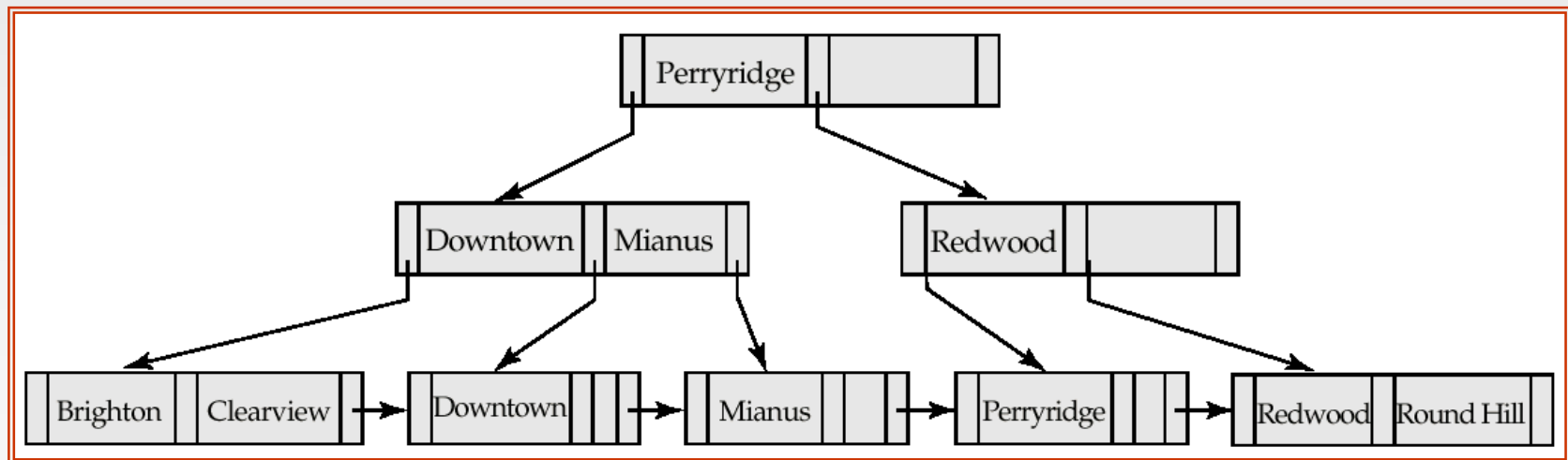
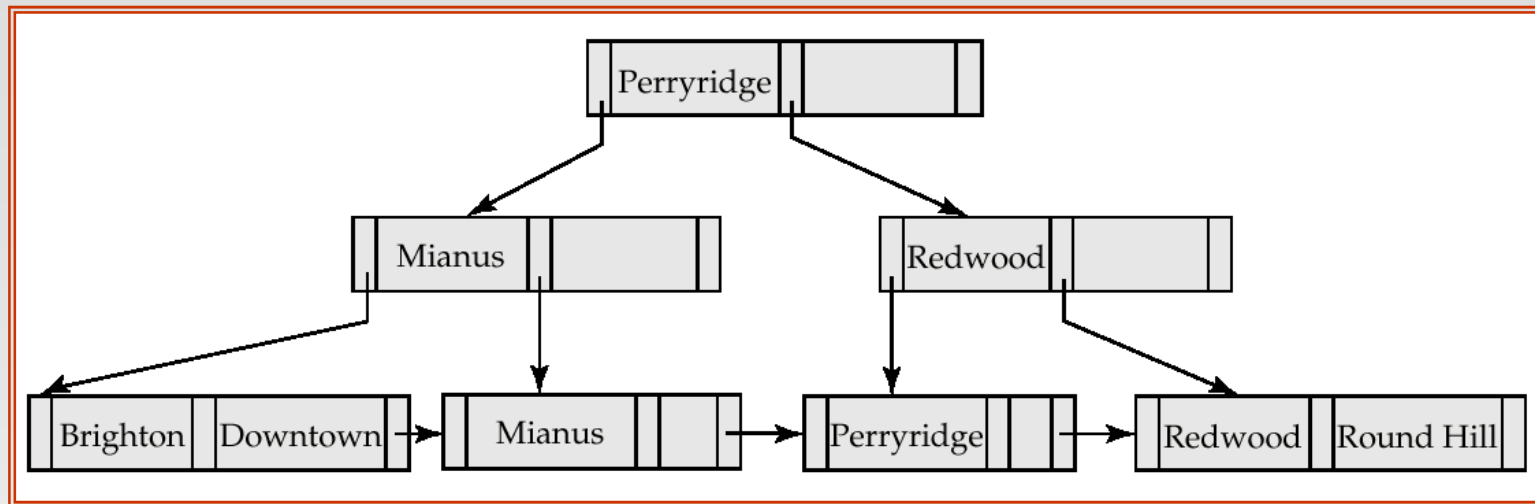


Result of splitting node containing Brighton and Downtown on inserting Clearview

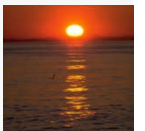




Updates on B⁺-Trees: Insertion (Cont.)



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





procedure insert(value K , pointer P)

find the leaf node L that should contain key value K

if (L has less than $n - 1$ key values) **then**

insert_in_leaf(L, K, P) // L is in memory

else begin /* L has $n - 1$ key values already, split it */

 Create node L'

 Copy $L.P_1 \dots L.K_{n-1}$ to a block of memory T that can
 hold n (pointer, key-value) pairs

insert_in_leaf(T, K, P) // T is in memory

 Set $L'.P_1 = L.P_n$; Set $L.P_n = L'$ // leaf blocks are chained

 Erase $L.P_1$ through $L.K_{n-1}$ from L

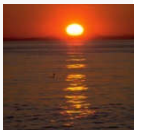
 Copy $T.P_1$ through $T.K_{\lceil n/2 \rceil}$ from T into L starting at $L.P_1$

 Insert $T.P_{\lceil n/2 \rceil + 1}$ through $T.K_n$ from T into L' before $L'.P_1$

 Let K' be the smallest key-value in L'

insert_in_parent(L, K', L')

end





procedure insert_in_leaf(*node L, value K, pointer P*)

if K is less than $L.K_1$

then insert P, K into L just before $L.P_1$

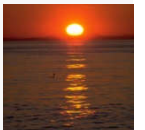
else

begin

 Let K_i be the least value in L that is less than K

 insert P, K into L just after $L.K_i$

end





procedure insert_in_parent(*node N, value K', node N'*) // *N* was split

if *N* is the root of the tree //the last insert was done in the root

then begin

create a new node *R* containing *N, K', N'*

// *N* and *N'* are pointers

make *R* the root of the tree

return

end

Let $P = \text{parent}(N)$

if (*P* has less than *n* pointers)

then insert (*K', N'*) in *P* just after *N*

else begin /* Split *P* */

Copy *P* to a block of memory *T* that can hold *P* and (*K', N'*)

Insert (*K', N'*) into *T* just after *N* // in *T* there are $n+1$ pointers

Erase all entries from *P*; Create node *P'*

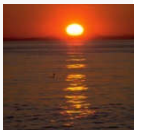
Copy $T.P_1, T.K_1, \dots, T.K_{\lceil n/2 \rceil - 1}, T.P_{\lceil n/2 \rceil}$ into *P*

Let $K'' = T.K_{\lceil n/2 \rceil}$

Copy $T.P_{\lceil n/2 \rceil + 1}, T.K_{\lceil n/2 \rceil + 1}, \dots, T.K_n, T.P_{n+1}$ into *P'*

insert_in_parent(*P, K'', P'*)

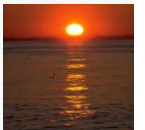
end





Updates on B⁺-Trees: Deletion

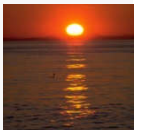
- Find the record to be deleted, and remove it from the main file and from the bucket (if present)
- Remove (search-key value, pointer) from the leaf node if there is no bucket or if the bucket has become empty
- If the node has too few entries due to the removal, and the entries in the node and a sibling 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.





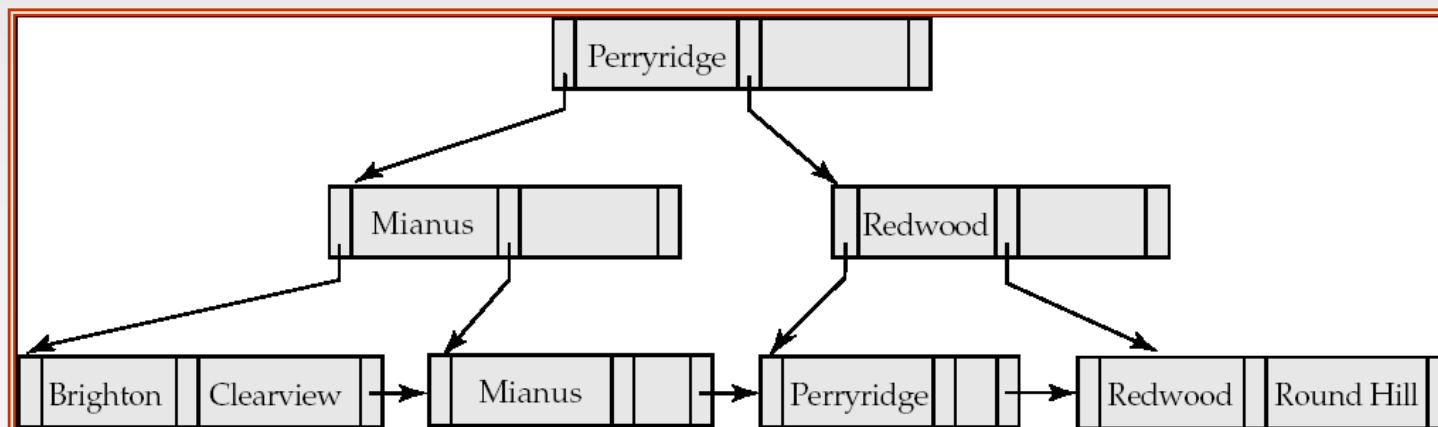
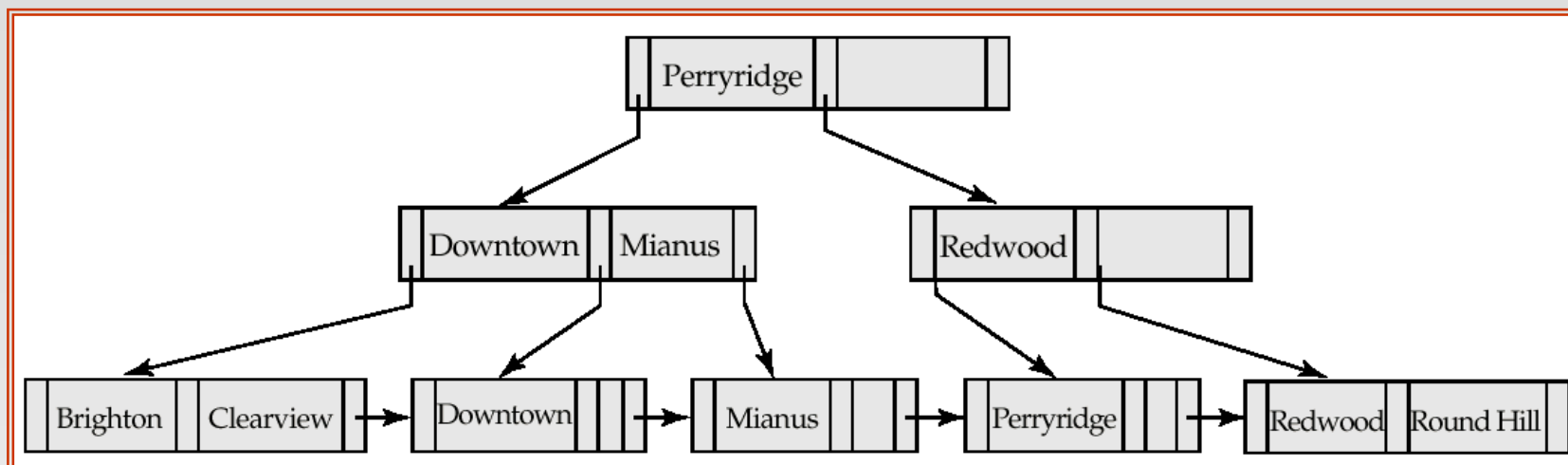
Updates on B⁺-Trees: Deletion

- Otherwise, if the node has too few entries due to the removal, and the entries in the node and a sibling don't 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.



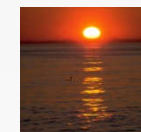


Examples of B⁺-Tree Deletion



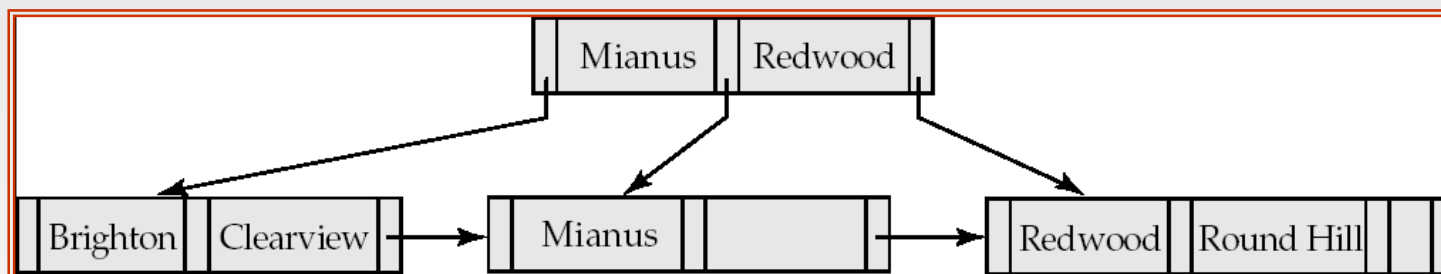
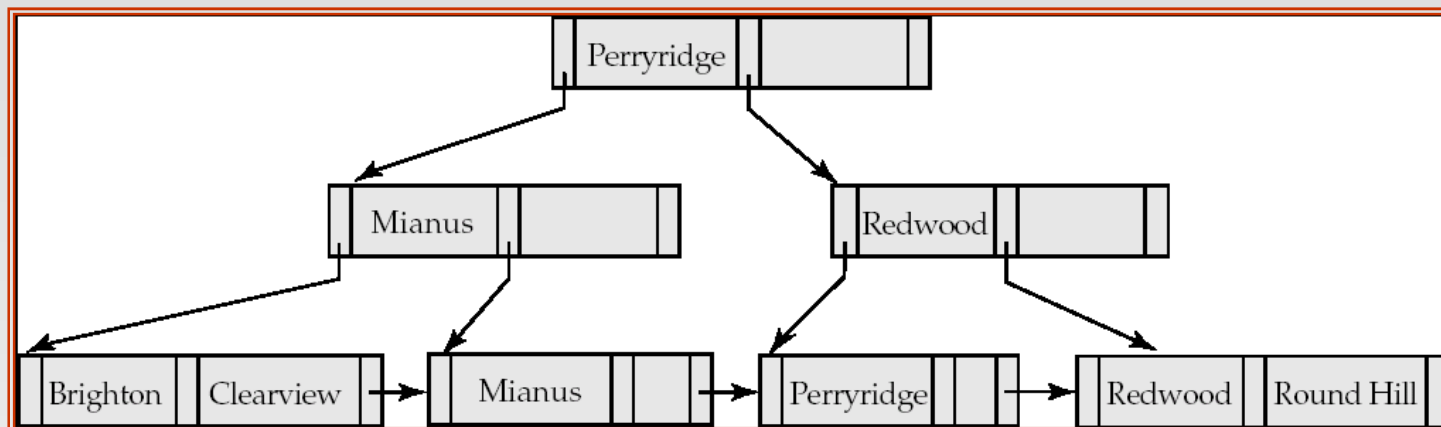
Before and after deleting "Downtown"

- 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.)



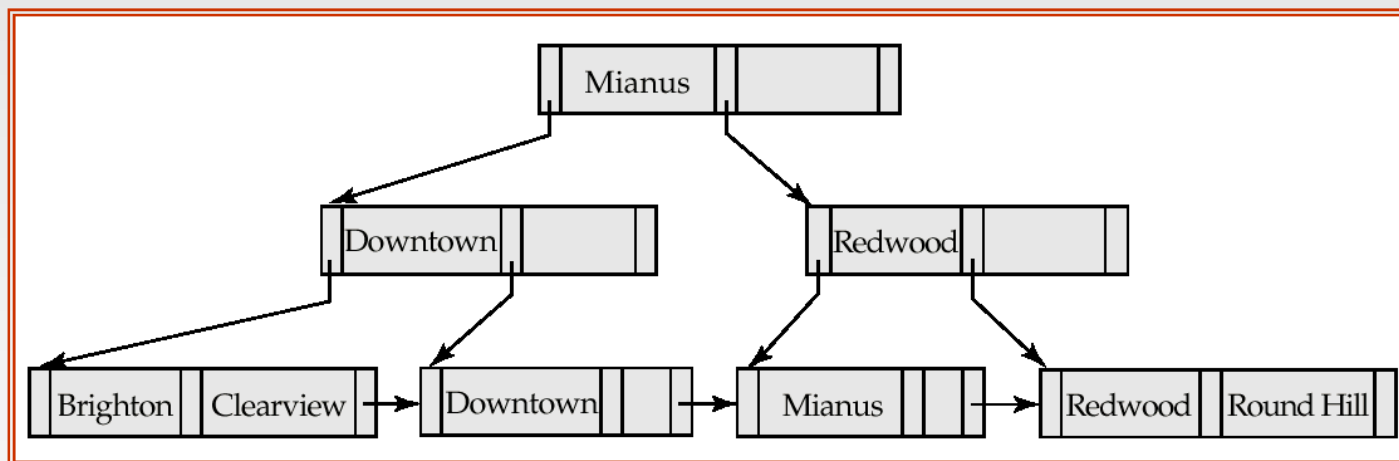
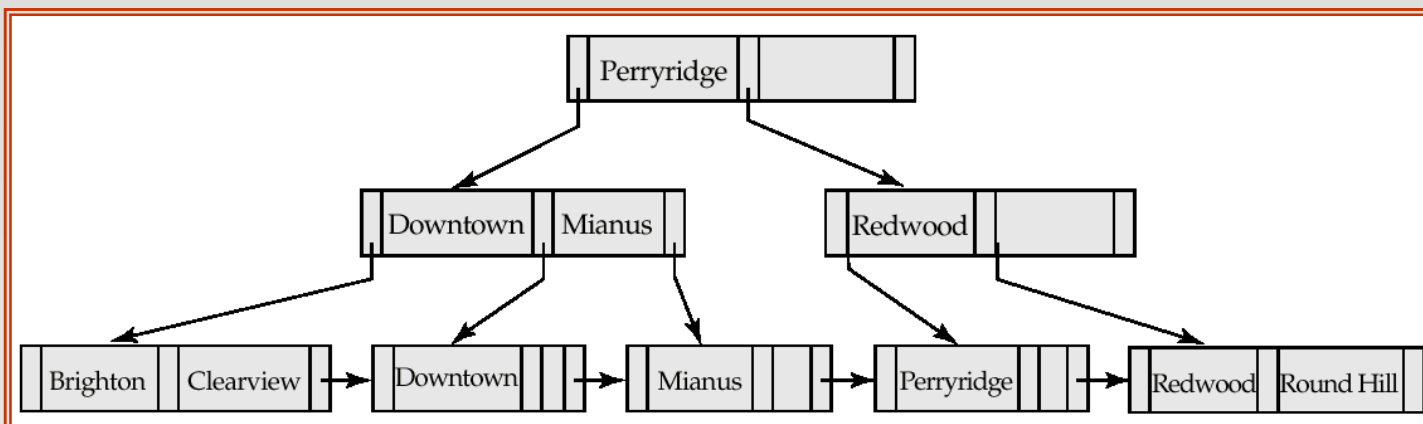
Deletion of “Perryridge” from result of previous example

- ❑ Node with “Perryridge” becomes underfull (actually empty, in this special case) and merged with its sibling.
- ❑ As a result “Perryridge” node’s parent became underfull, and was merged with its sibling (and an entry was deleted from their parent)
- ❑ Root node then had only one child, and was deleted and its child became the new root node



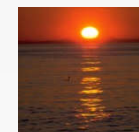


Example of B⁺-tree Deletion (Cont.)



Before and after deletion of “Perryridge” from earlier example

- Parent of leaf containing Perryridge became underfull, and borrowed a pointer from its left sibling
- Search-key value in the parent’s parent changes as a result





procedure delete(*value K, pointer P*)

 find the leaf node L that contains (K, P)

 delete_entry(L, K, P)

procedure delete_entry(*node N, value K, pointer P*)

 delete from $N (K, P)$

if (N is the root **and** N has only one remaining child)

then make the child of N the new root of the tree and delete N

else if (N has too few values/pointers) **then begin** //if N has enough values nothing to do

 Let N' be the previous or next child of $parent(N)$

 Let K' be the value between pointers N and N' in $parent(N)$

if (entries in N and N' can fit in a single node)

then begin /* Coalesce nodes */

if (N is a predecessor of N') **then** swap_variables(N, N')

if (N is not a leaf)

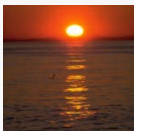
then append K' and all pointers and values in N to N'

else begin append all (K_i, P_i) pairs in N to N' ; set $N'.P_n = N.P_n$;

end (K' will not be inserted in leaf node)

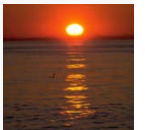
 delete_entry($parent(N), K', N$); delete node N

end





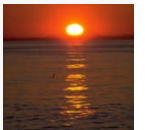
```
else begin /*redistribution: borrow an entry from  $N'$  */  
  if ( $N'$  is a predecessor of  $N$ ) then begin  
    if ( $N$  is a nonleaf node) then begin  
      let  $m$  be such that  $N'.P_m$  is the last pointer in  $N'$   
      remove ( $N'.K_{m-1}$  ,  $N'.P_m$  ) from  $N'$   
      insert ( $N'.P_m$  ,  $K'$  ) as the first pointer and value in  $N$   
      by shifting other pointers and values right  
      replace  $K'$  in  $parent(N)$  by  $N'.K_{m-1}$   
    end  
    else begin  
      let  $m$  be such that ( $N'.P_m$  ,  $N'.K_m$ ) is the last pointer/value  
      pair in  $N'$   
      remove ( $N'.P_m$  ,  $N'.K_m$ ) from  $N'$   
      insert ( $N'.P_m$  ,  $N'.K_m$ ) as the first pointer and value in  $N$ ,  
      by shifting other pointers and values right  
      replace  $K'$  in  $parent(N)$  by  $N'.K_m$   
    end  
  end  
else .... symmetric to the then case....  
  
end  
end
```





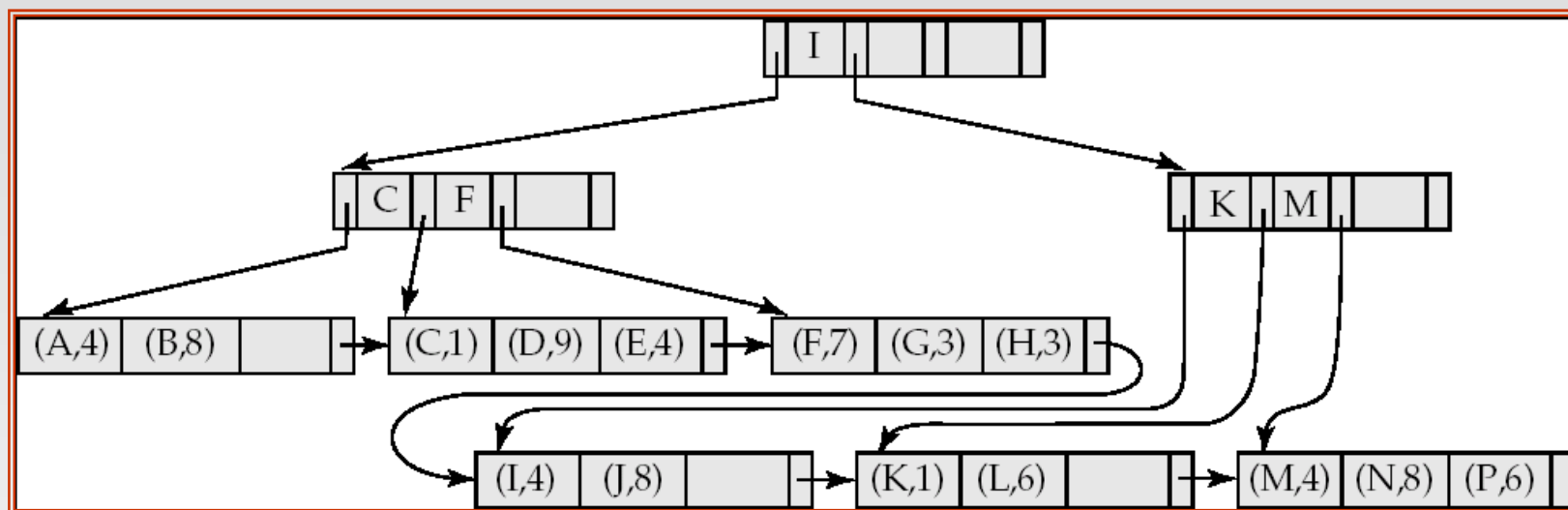
B⁺-Tree File Organization

- 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 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 non-leaf 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.





B⁺-Tree File Organization (Cont.)



Example of B⁺-tree File Organization

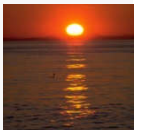
- ❑ Good space utilization important since records use more space than pointers.
- ❑ To improve space utilization, involve more sibling nodes in redistribution during splits and merges
 - ❑ Involving 2 siblings in redistribution (to avoid split / merge where possible) results in each node having at least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





Indexing Strings

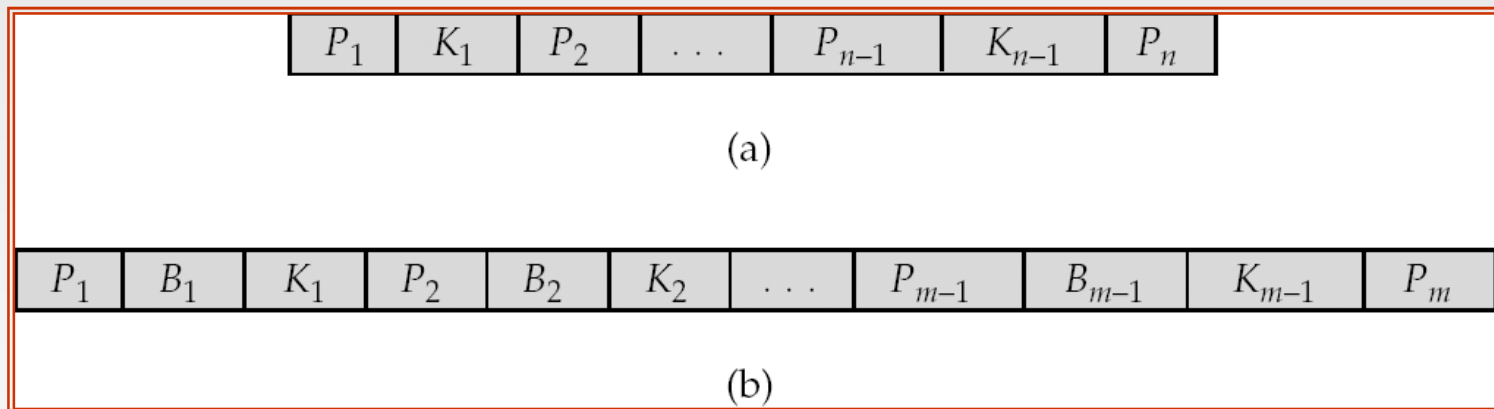
- Variable length strings as keys
 - Variable fanout
 - Use space utilization as criterion for splitting, not number of pointers
- Prefix compression
 - Key values at internal nodes can be prefixes of full key
 - ▶ Keep enough characters to distinguish entries in the subtrees separated by the key value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



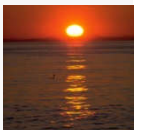


B-Tree Index Files

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

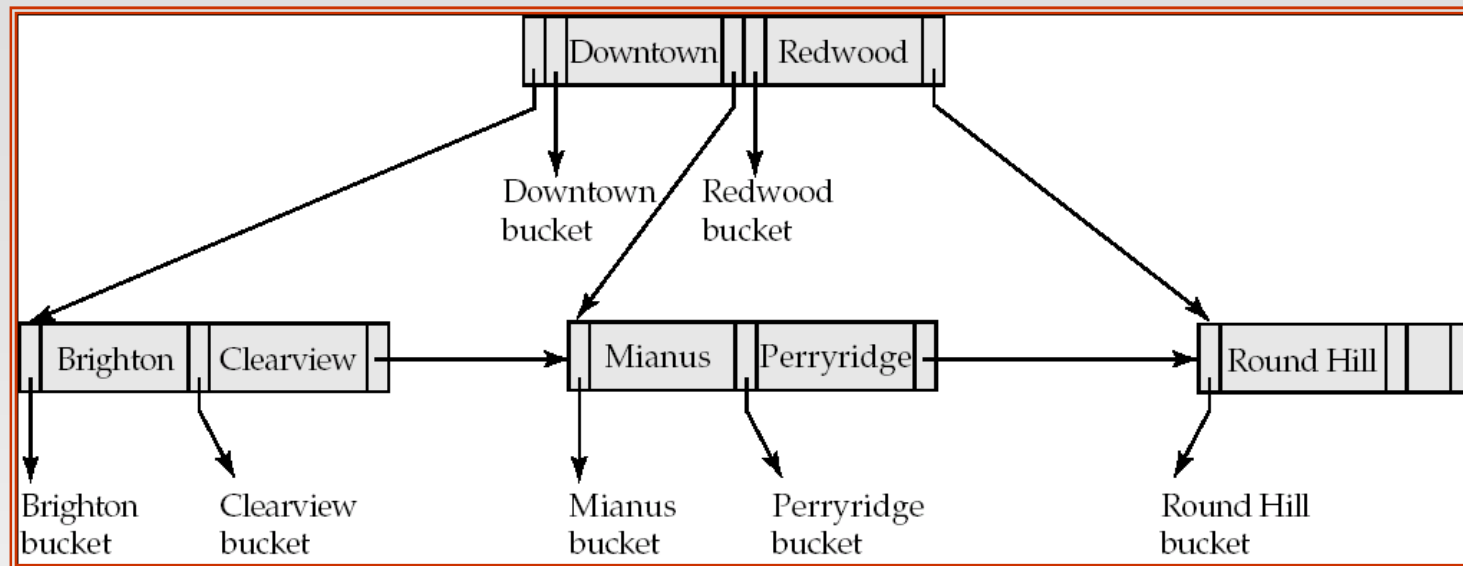


- Non-leaf node – pointers B_i are the bucket or file record pointers.

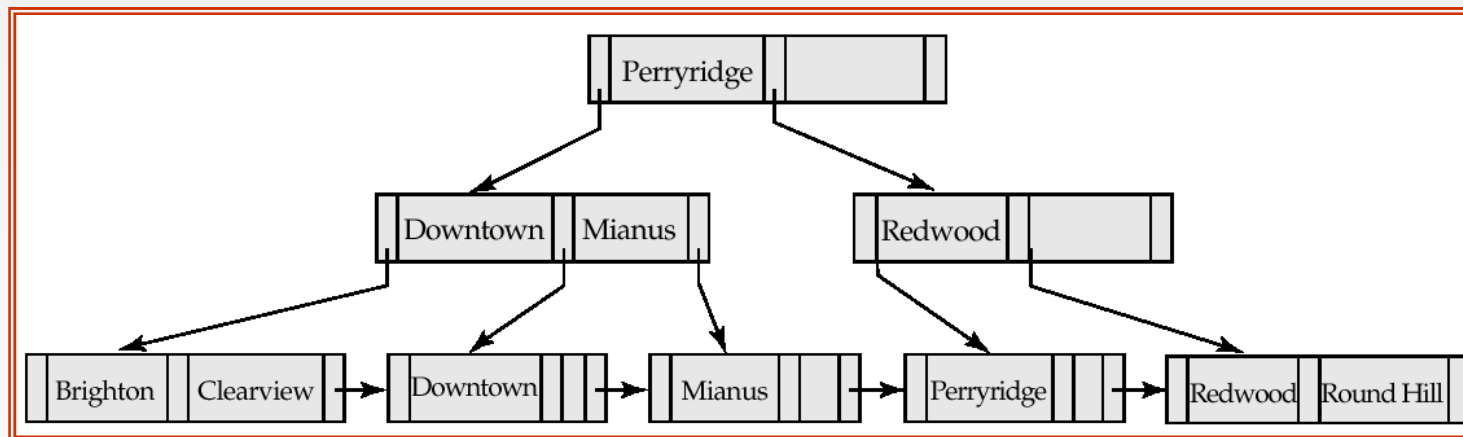




B-Tree Index File Example



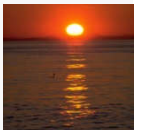
B-tree (above) and B+-tree (below) on same data





B-Tree Index Files (Cont.)

- Advantages of B-Tree indices:
 - May use less tree nodes than a corresponding B⁺-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.





Multiple-Key Access

- Use multiple indices for certain types of queries.

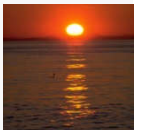
- Example:

select *account_number*

from *account*

where *branch_name* = "Perryridge" **and** *balance* = 1000

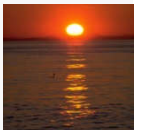
- Possible strategies for processing query using indices on single attributes:
 1. Use index on *branch_name* to find accounts with balances of \$1000; test *branch_name* = "Perryridge".
 2. Use index on *balance* to find accounts with balances of \$1000; test *branch_name* = "Perryridge".
 3. Use *branch_name* index to find pointers to all records pertaining to the Perryridge branch. Similarly use index on *balance*. 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. (*branch_name*, *balance*)
- Lexicographic ordering: $(a_1, a_2) < (b_1, b_2)$ if either
 - $a_1 < b_1$, or
 - $a_1 = b_1$ and $a_2 < b_2$

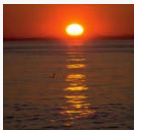




Indices on Multiple Attributes

Suppose we have an index on combined search-key
(*branch_name*, *balance*).

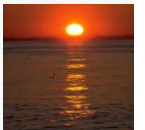
- With the **where** clause
 where *branch_name* = "Perryridge" **and** *balance* = 1000
the index on (*branch_name*, *balance*) can be used to fetch only records that satisfy both conditions.
 - Using separate indices is less efficient — we may fetch many records (or pointers) that satisfy only one of the conditions.
- Can also efficiently handle
 where *branch_name* = "Perryridge" **and** *balance* < 1000
- But cannot efficiently handle
 where *branch_name* < "Perryridge" **and** *balance* = 1000
 - May fetch many records that satisfy the first but not the second condition





Non-Unique Search Keys

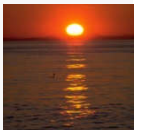
- Alternatives:
 - Buckets on separate block (bad idea)
 - List of tuple pointers with each key
 - ▶ Extra code to handle long lists
 - ▶ Deletion of a tuple can be expensive
 - ▶ Low space overhead, no extra cost for queries
 - Make search key unique by adding a record-identifier
 - ▶ Extra storage overhead for keys
 - ▶ Simpler code for insertion/deletion
 - ▶ Widely used





Other Issues

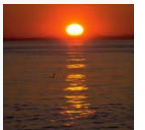
- Covering indices
 - Add extra attributes to index so (some) queries can avoid fetching the actual records
 - ▶ Particularly useful for secondary indices
 - ex. Non-clustering index on *account-number* of the *account* relation store the value of *balance*
 - Can store extra attributes only at leaf
- 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 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





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.

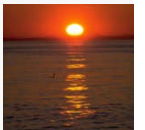




Example of Hash File Organization (Cont.)

Hash file organization of *account* file, using *branch_name* as key
(See figure in next slide.)

- There are 10 buckets,
- the *i-th* letter in the alphabet is represented by the integer *i*.
- The hash function returns the sum of the representations of the characters modulo 10
 - E.g. $h(\text{Perryridge}) = 5$ $h(\text{Round Hill}) = 3$ $h(\text{Brighton}) = 3$





Example of Hash File Organization

Hash file organization of *account* file, using *branch_name* as key
(see previous slide for details).

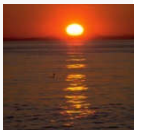
bucket 0			
bucket 1			
bucket 2			
bucket 3	A-217	Brighton	750
	A-305	Round Hill	350
bucket 4	A-222	Redwood	700
bucket 5	A-102	Perryridge	400
	A-201	Perryridge	900
	A-218	Perryridge	700
bucket 6			
bucket 7	A-215	Mianus	700
bucket 8	A-101	Downtown	500
	A-110	Downtown	600
bucket 9			





Hash Functions

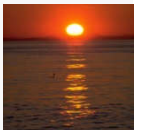
- ❑ Worst hash function maps all search-key values to the same bucket; this makes access time proportional to the number of search-key values in the file.
- ❑ An ideal hash function is **uniform**, i.e., each bucket is assigned the same number of search-key values from the set of *all* possible values.
- ❑ 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

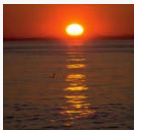
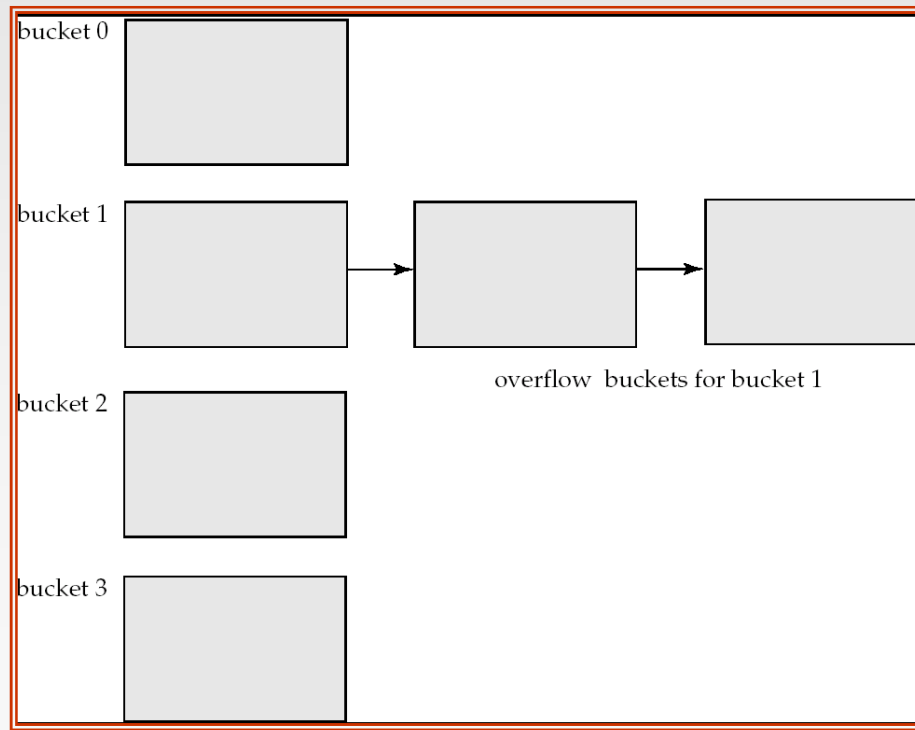
- Bucket overflow can occur because of
 - Insufficient buckets
 - Skew in distribution of records. This can occur due to two reasons:
 - ▶ multiple records have same search-key value
 - ▶ chosen hash function produces non-uniform distribution of key values
- Although the probability of bucket overflow can be reduced, it cannot be eliminated; it is handled by using *overflow buckets*.





Handling of Bucket Overflows (Cont.)

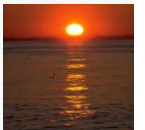
- **Overflow chaining** – the overflow buckets of a given bucket are chained together in a linked list.
- Above scheme is called **closed hashing**.
 - An alternative, called **open hashing**, which does not use overflow buckets, is not suitable for database applications.





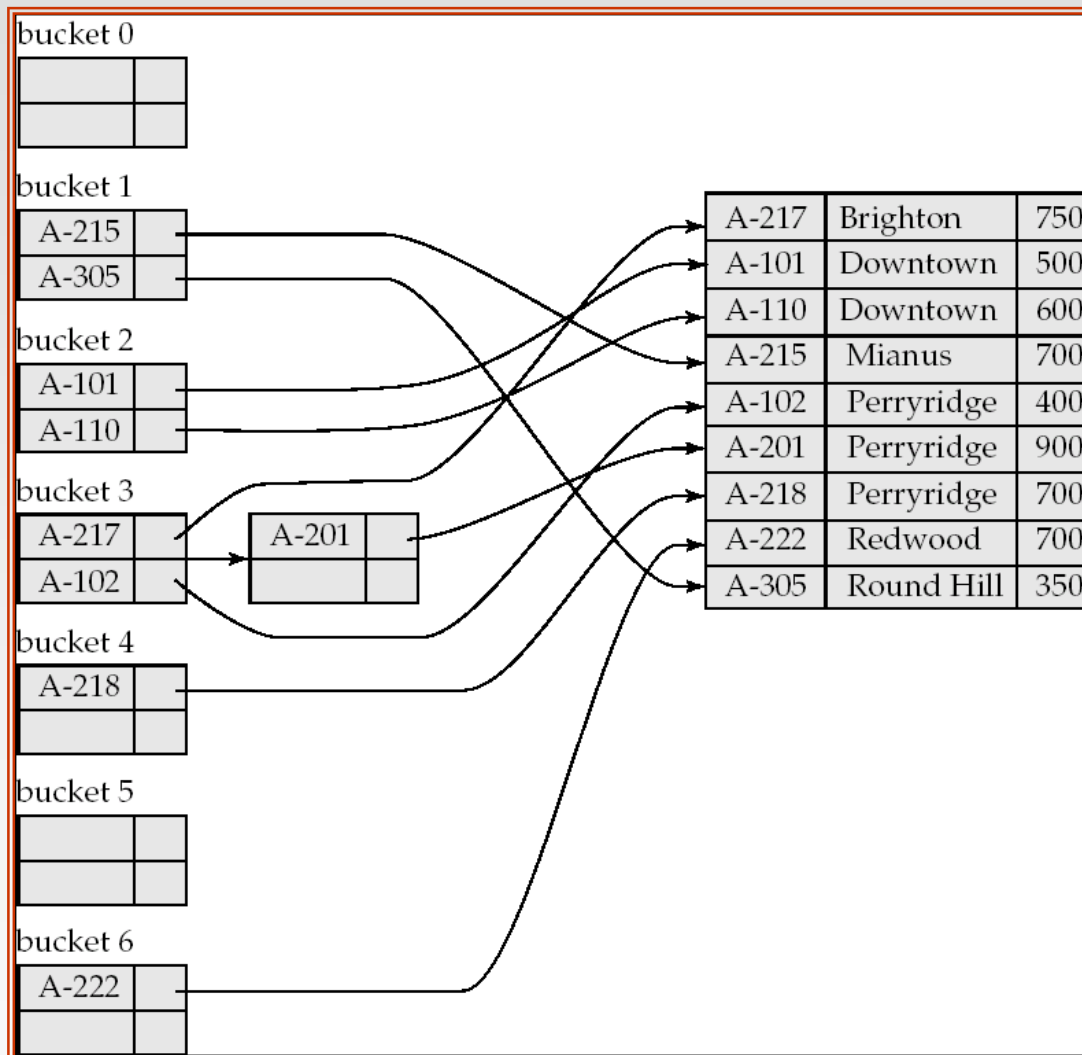
Hash Indices

- Hashing can be used not only for file organization, but also for index-structure creation.
- A **hash index** organizes the search keys, with their associated record pointers, into a hash file structure.
- Strictly speaking, hash indices are always secondary indices
 - if the file itself is organized using hashing, a separate primary hash index on it using the same search-key is unnecessary.
 - However, we use the term hash index to refer to both secondary index structures and hash organized files.





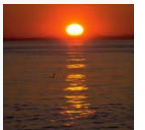
Example of Hash Index





Deficiencies of Static Hashing

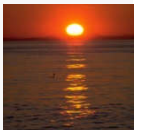
- In static hashing, function h maps search-key values to a fixed set of B of bucket addresses.
 - Databases grow with time. If initial number of buckets is too small, performance will degrade due to too much overflows.
 - If file size at some point in the future is anticipated and number of buckets allocated accordingly, significant amount of space will be wasted initially.
 - If database shrinks, again space will be wasted.
 - One option is periodic re-organization of the file with a new hash function, but it is very expensive.
- These problems can be avoided by using techniques that allow the number of buckets to be modified dynamically.





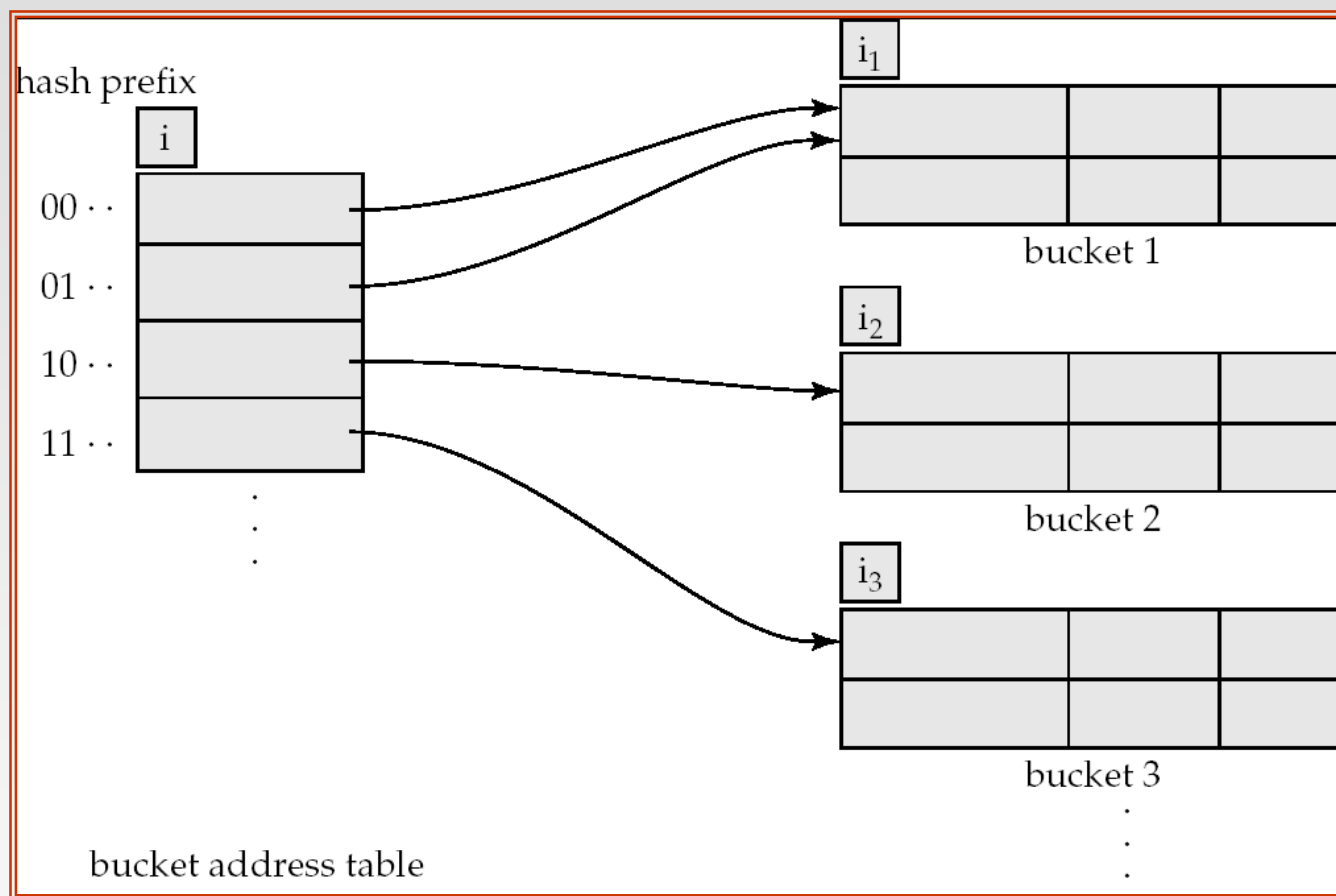
Dynamic Hashing

- 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.
 - 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



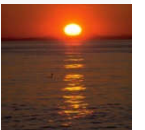
In this structure, $i_2 = i_3 = i$, whereas $i_1 = i - 1$ (see next slide for details)





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)

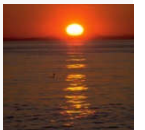




Updates in Extendable Hash Structure

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

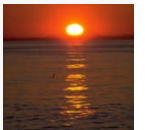
- If $i > i_j$ (more than one pointer to bucket j)
 - allocate a new bucket z , and set i_j and i_z to the old i_j ± 1 .
 - make the second half of the bucket address table entries pointing to j to point to z
 - remove and reinsert each record in bucket j .
 - recompute new bucket for K_j and insert record in the bucket (further splitting is required if the bucket is still full)
- If $i = i_j$ (only one pointer to bucket j)
 - 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.





Updates in Extendable Hash Structure (Cont.)

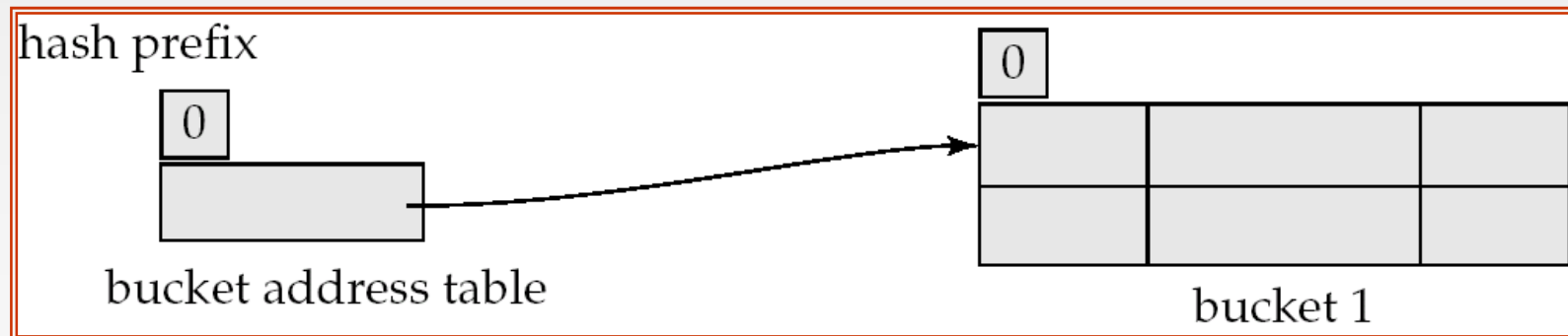
- When inserting a value, if the bucket is full after several splits (that is, i reaches some limit b) create an overflow bucket instead of splitting bucket entry table further.
- To delete a key value,
 - locate it in its bucket and remove it.
 - The bucket itself can be removed if it becomes empty (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



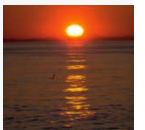


Use of Extendable Hash Structure: Example

<i>branch_name</i>	$h(\text{branch_name})$
Brighton	0010 1101 1111 1011 0010 1100 0011 0000
Downtown	1010 0011 1010 0000 1100 0110 1001 1111
Mianus	1100 0111 1110 1101 1011 1111 0011 1010
Perryridge	1111 0001 0010 0100 1001 0011 0110 1101
Redwood	0011 0101 1010 0110 1100 1001 1110 1011
Round Hill	1101 1000 0011 1111 1001 1100 0000 0001



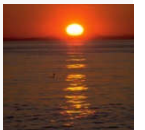
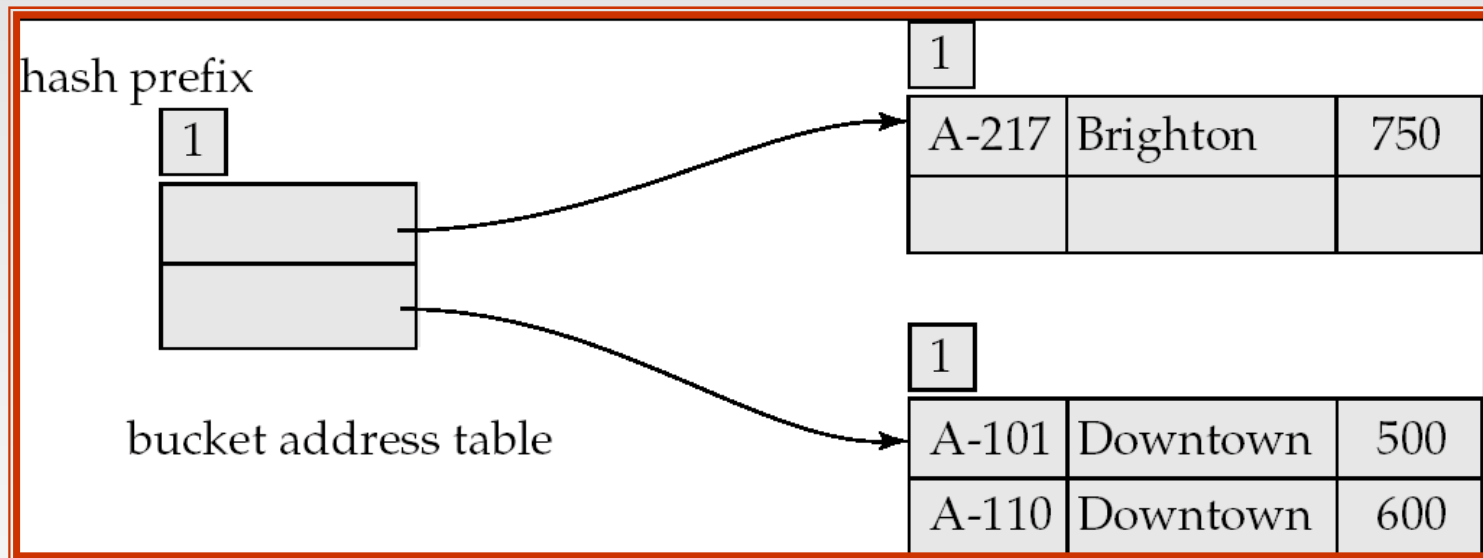
Initial Hash structure, bucket size = 2





Example (Cont.)

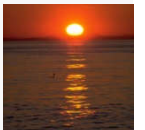
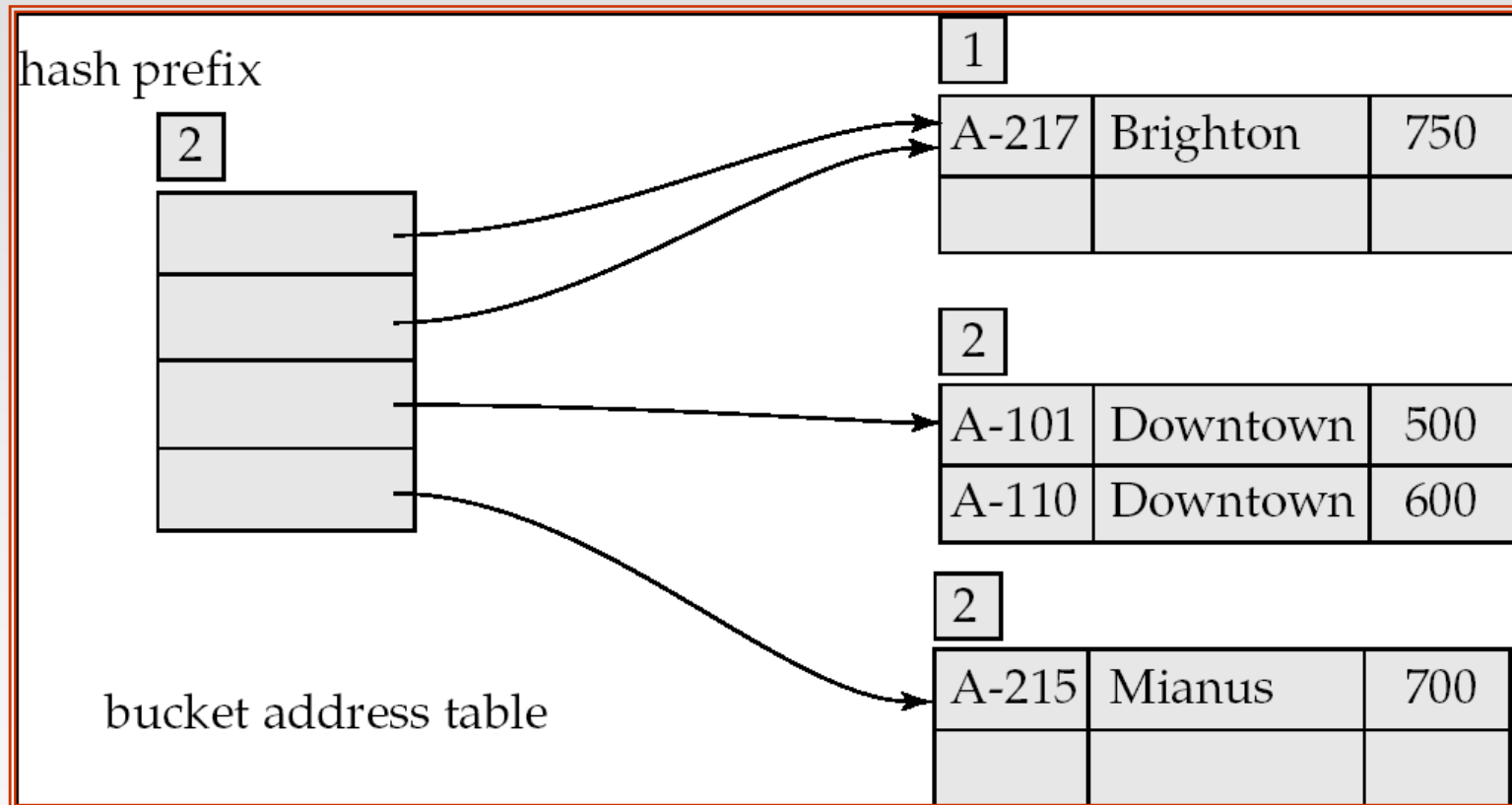
- Hash structure after insertion of one Brighton and two Downtown records





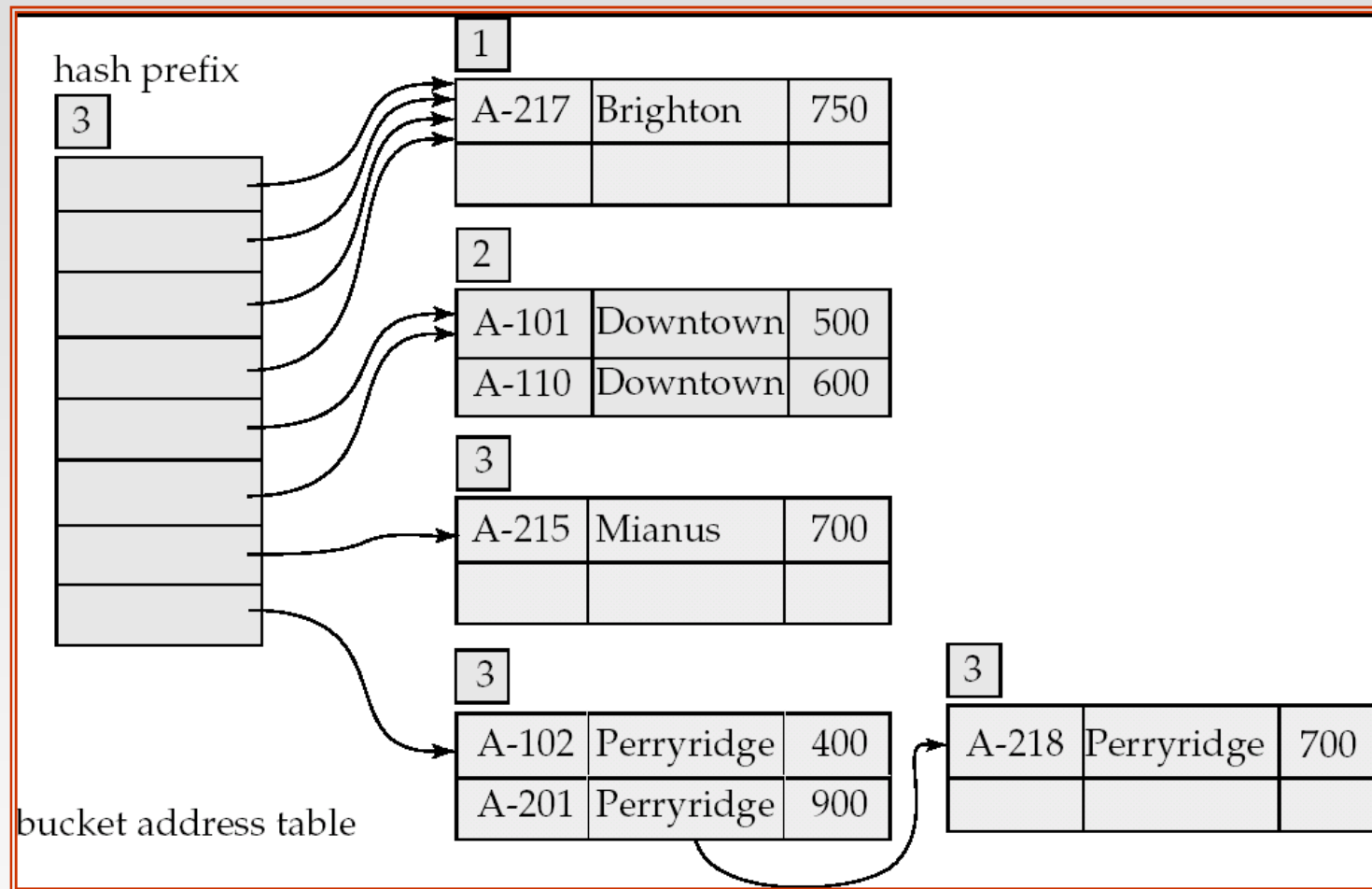
Example (Cont.)

Hash structure after insertion of Mianus record

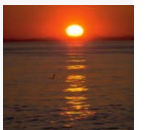




Example (Cont.)



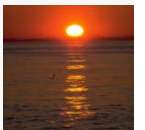
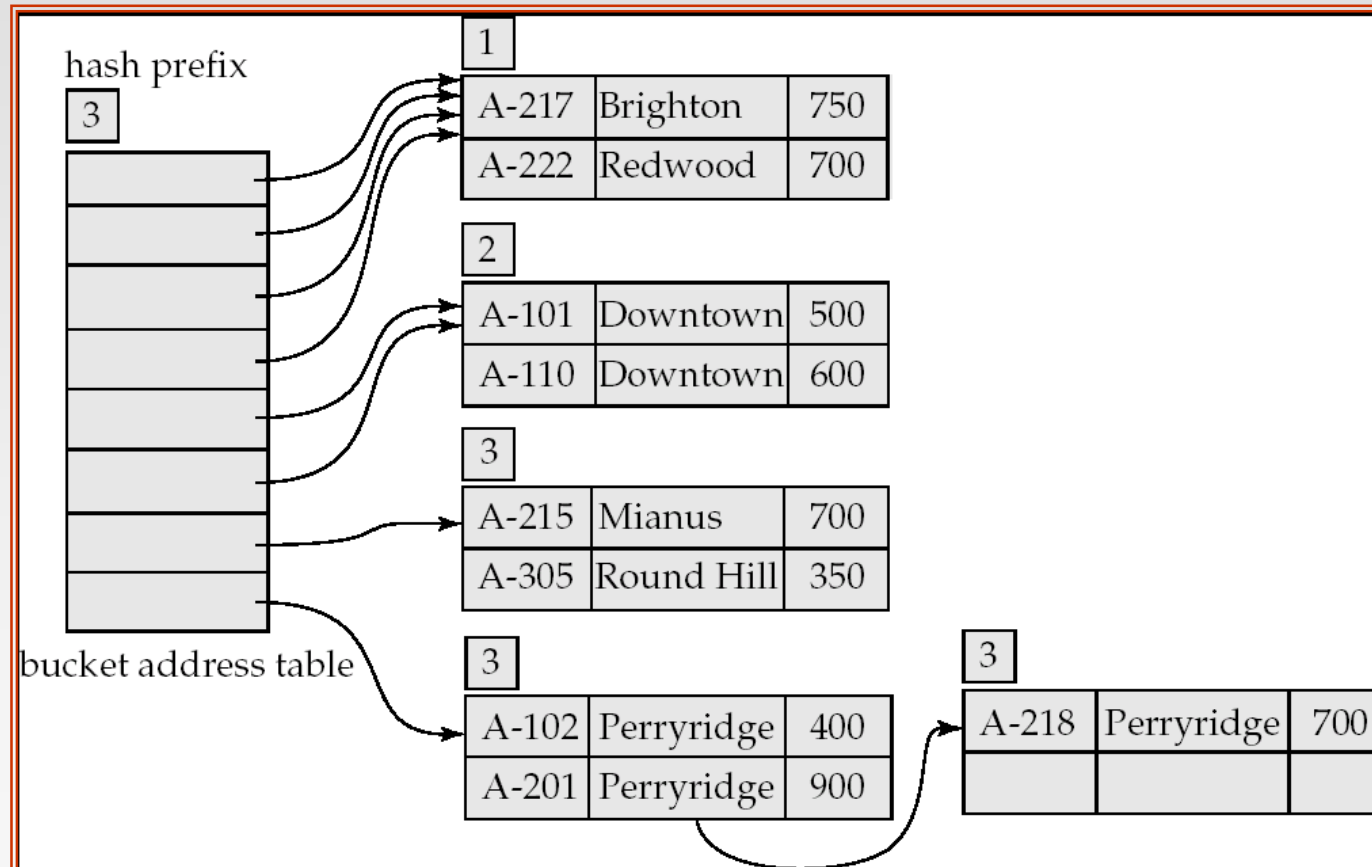
Hash structure after insertion of three Perryridge records





Example (Cont.)

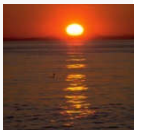
- Hash structure after insertion of Redwood and Round Hill records





Extendable Hashing vs. Other Schemes

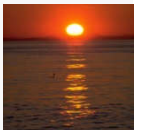
- Benefits of extendable hashing:
 - Hash performance does not degrade with growth of file
 - Minimal space overhead
- Disadvantages of extendable hashing
 - Extra level of indirection to find desired record
 - Bucket address table may itself become very big (larger than memory)
 - ▶ Need a tree structure to locate desired record in the structure!
 - Changing size of bucket address table is an expensive operation
- **Linear hashing** is an alternative mechanism which avoids these disadvantages at the possible cost of more bucket overflows





Comparison of Ordered Indexing and Hashing

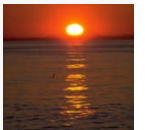
- ❑ Cost of periodic re-organization
- ❑ Relative frequency of insertions and deletions
- ❑ Is it desirable to optimize average access time at the expense of worst-case access time?
- ❑ Expected type of queries:
 - ❑ Hashing is generally better at retrieving records having a specified value of the key.
 - ❑ If range queries are common, ordered indices are to be preferred





Bitmap Indices

- Bitmap indices are a special type of index designed for efficient querying on multiple keys
- Records in a relation are assumed to be numbered sequentially from, say, 0
 - Given a number n it must be easy to retrieve record n
 - ▶ Particularly easy if records are of fixed size
- Applicable on attributes that take on a relatively small number of distinct values
 - E.g. gender, country, state, ...
 - E.g. income-level (income broken up into a small number of levels such as 0-9999, 10000-19999, 20000-50000, 50000- infinity)
- A bitmap is simply an array of bits

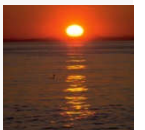




Bitmap Indices (Cont.)

- In its simplest form a bitmap index on an attribute has a bitmap for each value of the attribute
 - Bitmap has as many bits as records
 - In a bitmap for value v , the bit for a record is 1 if the record has the value v for the attribute, and is 0 otherwise

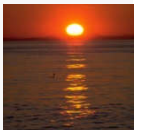
record number	<i>name</i>	<i>gender</i>	<i>address</i>	<i>income_level</i>	Bitmaps for <i>gender</i>		Bitmaps for <i>income_level</i>	
					m	1 0 0 1 0		
					f	0 1 1 0 1		
0	John	m	Perryridge	L1			L1	1 0 1 0 0
1	Diana	f	Brooklyn	L2			L2	0 1 0 0 0
2	Mary	f	Jonestown	L1			L3	0 0 0 0 1
3	Peter	m	Brooklyn	L4			L4	0 0 0 1 0
4	Kathy	f	Perryridge	L3			L5	0 0 0 0 0





Bitmap Indices (Cont.)

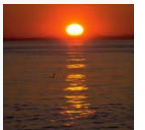
- Bitmap indices are useful for queries on multiple attributes
 - not particularly useful for single attribute queries
- Queries are answered using bitmap operations
 - Intersection (and)
 - Union (or)
 - Complementation (not)
- Each operation takes two bitmaps of the same size and applies the operation on corresponding bits to get the result bitmap
 - E.g. $100110 \text{ AND } 110011 = 100010$
 $100110 \text{ OR } 110011 = 110111$
 $\text{NOT } 100110 = 011001$
 - Males with income level L1: $10010 \text{ AND } 10100 = 10000$
 - ▶ Can then retrieve required tuples.
 - ▶ Counting number of matching tuples is even faster





Bitmap Indices (Cont.)

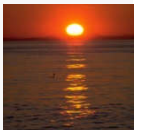
- Bitmap indices generally very small compared with relation size
 - E.g. if record is 100 bytes, space for a single bitmap is 1/800 of space used by relation.
 - ▶ If number of distinct attribute values is 8, bitmap is only 1% of relation size
- Deletion needs to be handled properly
 - **Existence bitmap** to note if there is a valid record at a record location
 - Needed for complementation
 - ▶ $\text{not}(A=v)$: *(NOT bitmap-A-v) AND ExistenceBitmap*
- Should keep bitmaps for all values, even null value
 - To correctly handle SQL null semantics for $\text{NOT}(A=v)$:
 - ▶ intersect above result with *(NOT bitmap-A-Null)*





Efficient Implementation of Bitmap Operations

- Bitmaps are packed into words; a single word and (a basic CPU instruction) computes and of 32 or 64 bits at once
 - E.g. 1-million-bit maps can be anded with just 31,250 instruction
- Counting number of 1s can be done fast by a trick:
 - Use each byte to index into a precomputed array of 256 elements each storing the count of 1s in the binary representation
 - ▶ Can use pairs of bytes to speed up further at a higher memory cost
 - Add up the retrieved counts
- Bitmaps can be used instead of Tuple-ID lists at leaf levels of B⁺-trees, for values that have a large number of matching records
 - Worthwhile if $> 1/64$ of the records have that value, assuming a tuple-id is 64 bits
 - Above technique merges benefits of bitmap and B⁺-tree indices





Index Definition in SQL

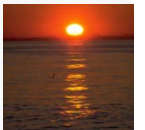
- Create an index

create index <index-name> **on** <relation-name>
(<attribute-list>)

E.g.: **create index** *b-index* **on** *branch(branch_name)*

- Use **create unique index** to indirectly specify and enforce the condition that the search key is a candidate key is a candidate key.
 - Not really required if SQL **unique** integrity constraint is supported
- To drop an index

drop index <index-name>

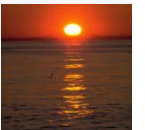




End of Chapter

Database System Concepts, 5th Ed.

©Silberschatz, Korth and Sudarshan
See www.db-book.com for conditions on re-use





Partitioned Hashing

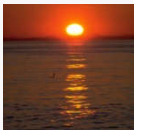
- Hash values are split into segments that depend on each attribute of the search-key.

(A_1, A_2, \dots, A_n) for n attribute search-key

- Example: $n = 2$, for *customer*, search-key being (*customer-street*, *customer-city*)

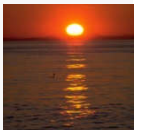
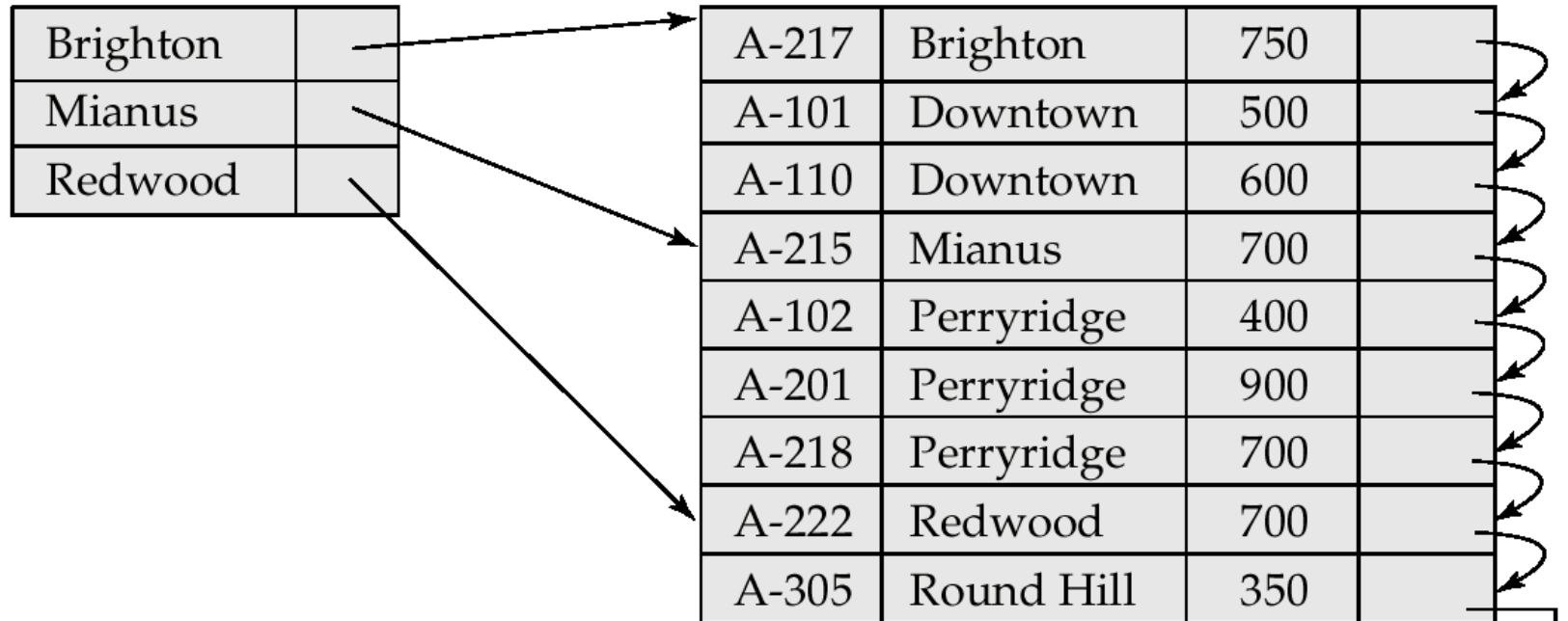
<i>search-key value</i>	<i>hash value</i>
(Main, Harrison)	101 111
(Main, Brooklyn)	101 001
(Park, Palo Alto)	010 010
(Spring, Brooklyn)	001 001
(Alma, Palo Alto)	110 010

- To answer equality query on single attribute, need to look up multiple buckets. Similar in effect to grid files.



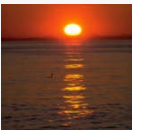


Sequential File For *account* Records





Deletion of “Perryridge” From the B⁺-Tree of Figure 12.12





Sample *account* File

A-217	Brighton	750
A-101	Downtown	500
A-1 10	Downtown	600
A-215	Mianus	700
A-102	Perryridge	400
A-201	Perryridge	900
A-218	Perryridge	700
A-222	Redwood	700
A-305	Round Hill	350

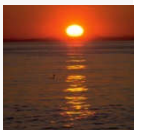




Figure 12.2

Brighton		→	A-217	Brighton	750		↙
Downtown		→	A-101	Downtown	500		↙
Mianus		→	A-110	Downtown	600		↙
Perryridge		→	A-215	Mianus	700		↙
Redwood		→	A-102	Perryridge	400		↙
Round Hill		→	A-201	Perryridge	900		↙
		→	A-218	Perryridge	700		↙
		→	A-222	Redwood	700		↙
		→	A-305	Round Hill	350		↙

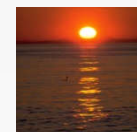




Figure 12.14

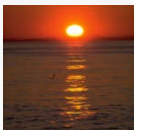
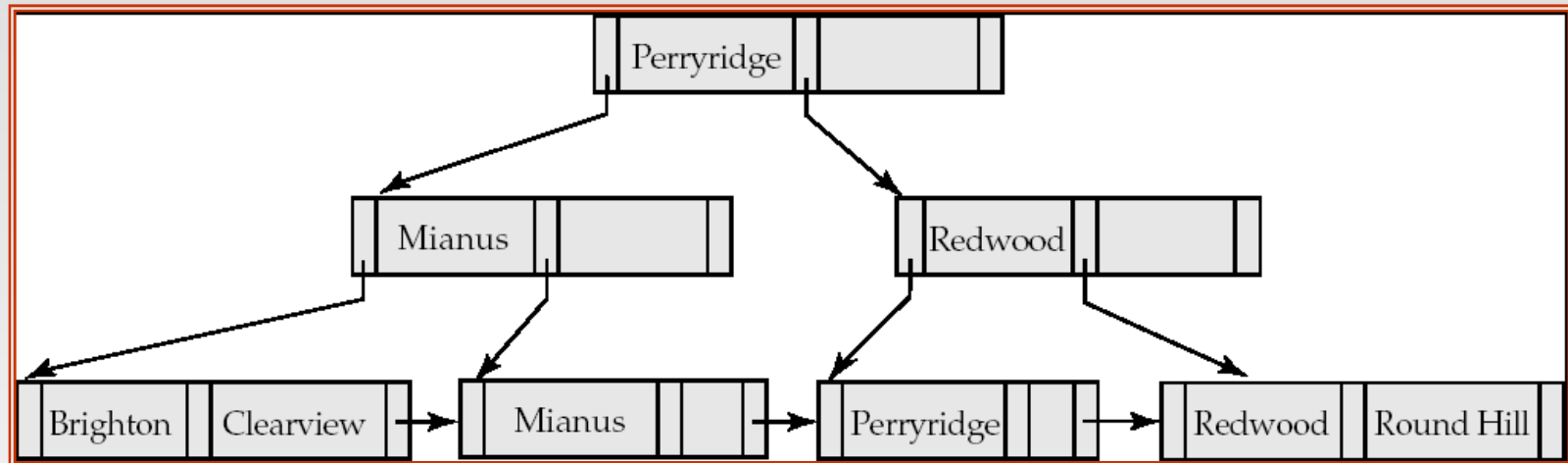




Figure 12.25

A-217	Brighton	750
A-101	Downtown	500
A-1 10	Downtown	600
A-215	Mianus	700
A-102	Perryridge	400
A-201	Perryridge	900
A-218	Perryridge	700
A-222	Redwood	700
A-305	Round Hill	350

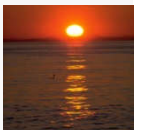
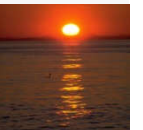




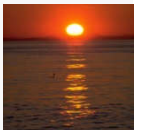
Figure 12.33





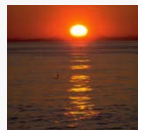
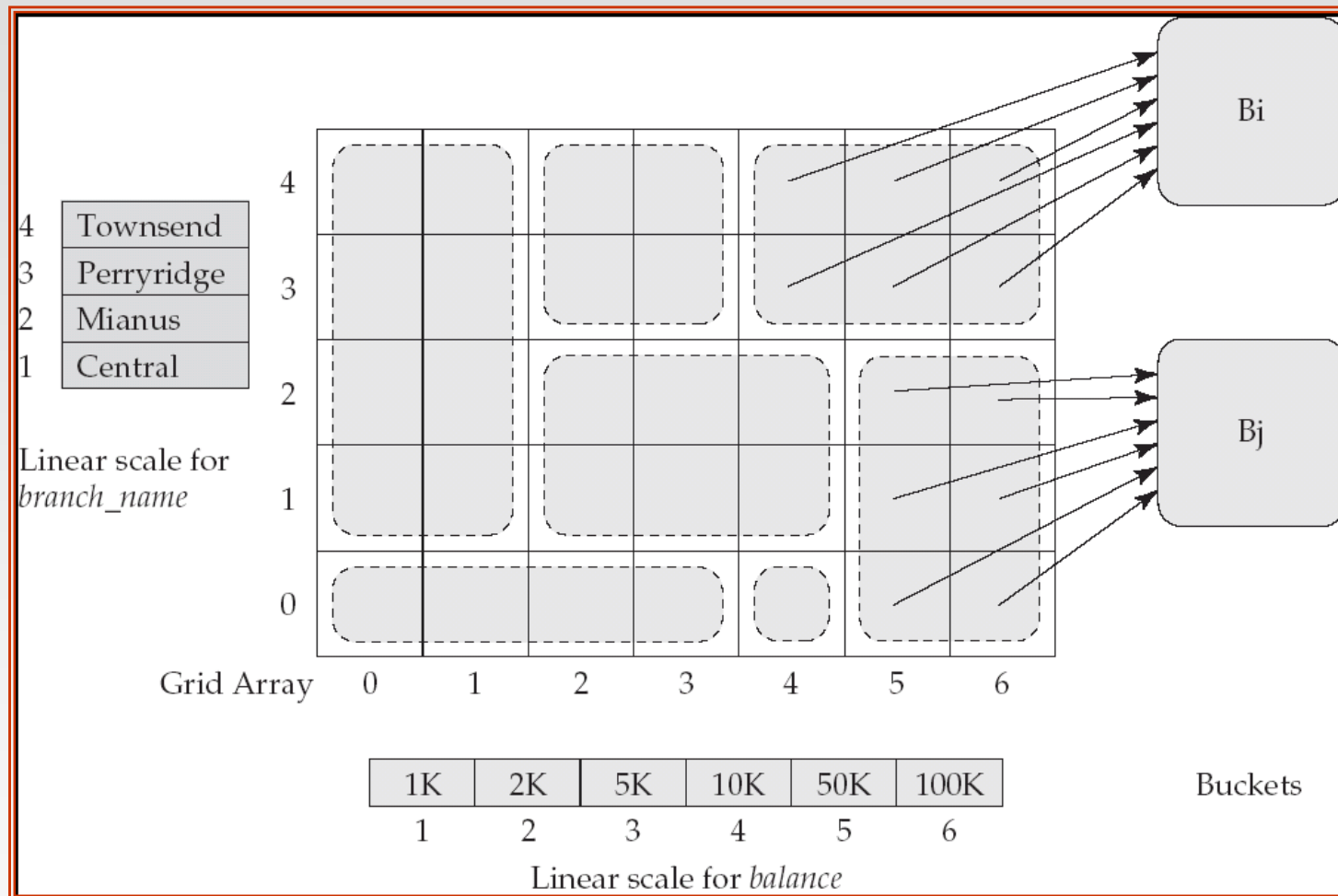
Grid Files

- Structure used to speed the processing of general multiple search-key queries involving one or more comparison operators.
- The **grid file** has a single grid array and one linear scale for each search-key attribute. The grid array has number of dimensions equal to number of search-key attributes.
- Multiple cells of grid array can point to same bucket
- To find the bucket for a search-key value, locate the row and column of its cell using the linear scales and follow pointer





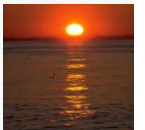
Example Grid File for *account*





Queries on a Grid File

- A grid file on two attributes A and B can handle queries of all following forms with reasonable efficiency
 - $(a_1 \leq A \leq a_2)$
 - $(b_1 \leq B \leq b_2)$
 - $(a_1 \leq A \leq a_2 \wedge b_1 \leq B \leq b_2), .$
- E.g., to answer $(a_1 \leq A \leq a_2 \wedge b_1 \leq B \leq b_2)$, use linear scales to find corresponding candidate grid array cells, and look up all the buckets pointed to from those cells.





Grid Files (Cont.)

- During insertion, if a bucket becomes full, new bucket can be created if more than one cell points to it.
 - Idea similar to extendable hashing, but on multiple dimensions
 - If only one cell points to it, either an overflow bucket must be created or the grid size must be increased
- Linear scales must be chosen to uniformly distribute records across cells.
 - Otherwise there will be too many overflow buckets.
- Periodic re-organization to increase grid size will help.
 - But reorganization can be very expensive.
- Space overhead of grid array can be high.
- R-trees (Chapter 23) are an alternative

