



Chapter 12: Indexing and Hashing

Rev. Sep 17, 2008

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 (e.g. $10000 < salary < 40000$)
- Access time
- Insertion time
- Deletion time
- Space overhead





Ordered Indices

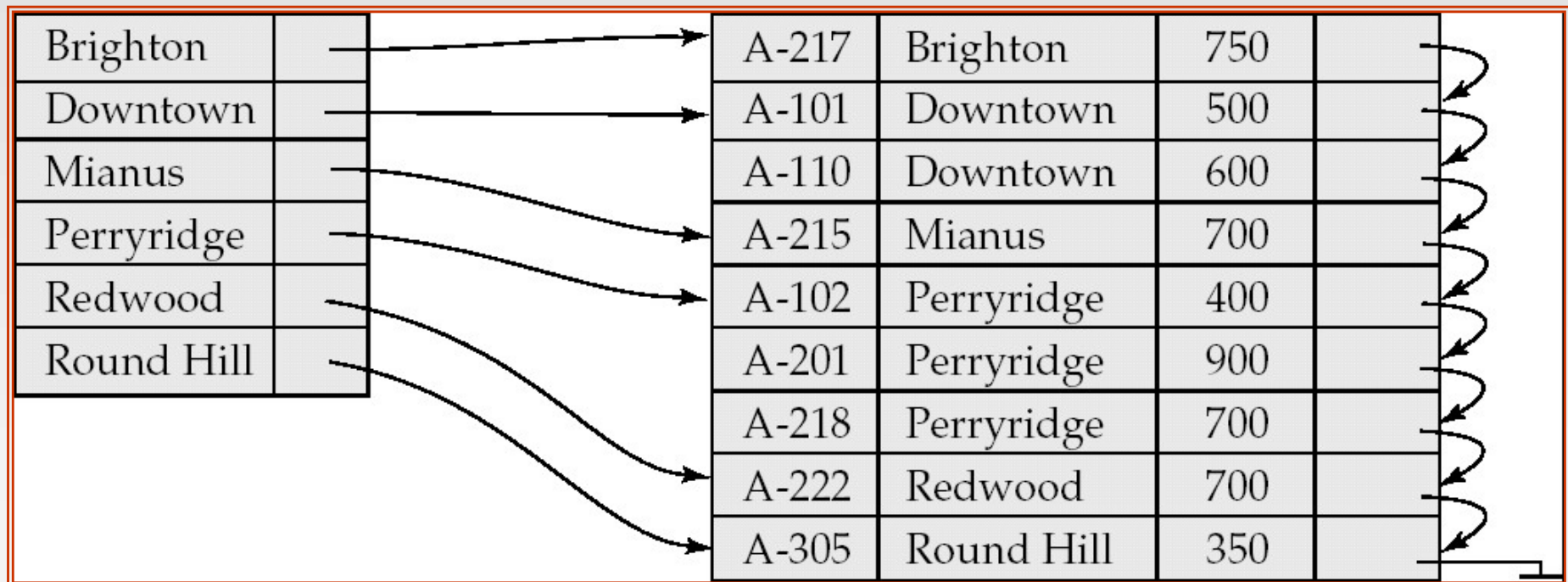
- In an **ordered index**, index entries are stored sorted on the search key value. E.g., author catalog in library.
- **Primary index**: in a sequentially ordered file, the index whose search key specifies the sequential order of the file.
 - Also called **clustering index**
 - The search key of a primary 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.





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

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	

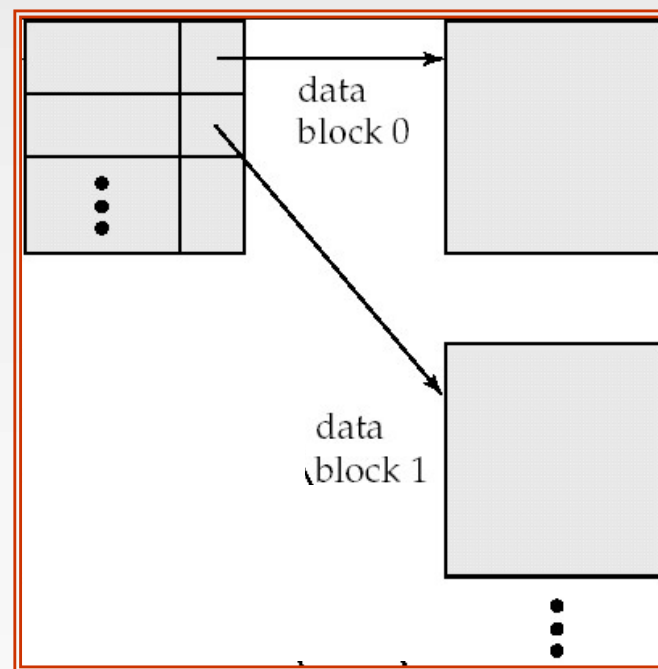
Diagram illustrating a sparse index structure. The index table on the left maps search-key values (Brighton, Mianus, Redwood) to specific record pointers (A-217, A-101, A-110, A-215, A-102, A-201, A-218, A-222, A-305). The main data table on the right shows the records in sequential order of their search-key values (Round Hill, Perryridge, Downtown, Mianus, Brighton). Arrows indicate the mapping from the index to the data records. A vertical line with arrows on the right side of the data table indicates sequential access.





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.





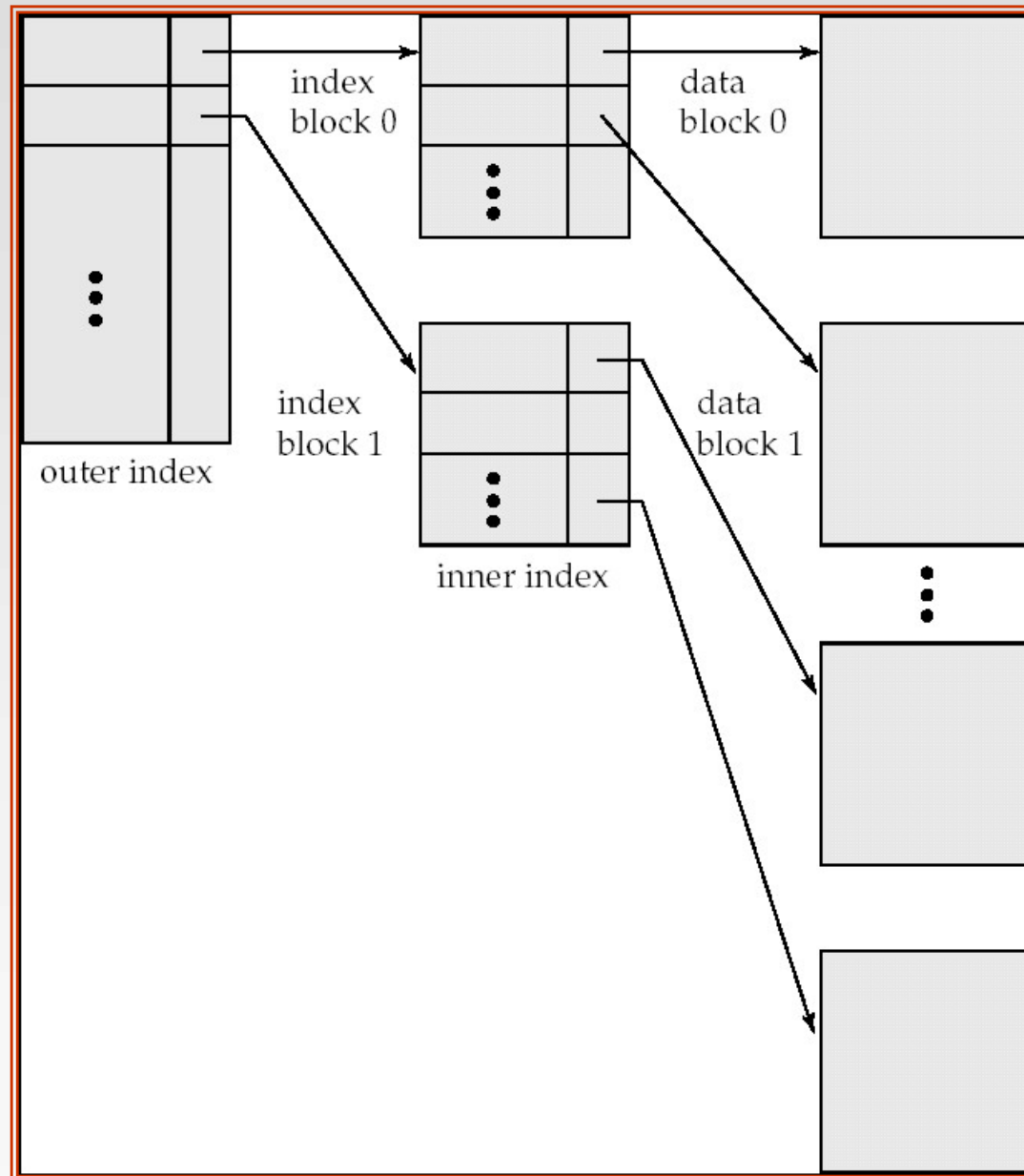
Multilevel Index

- 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.





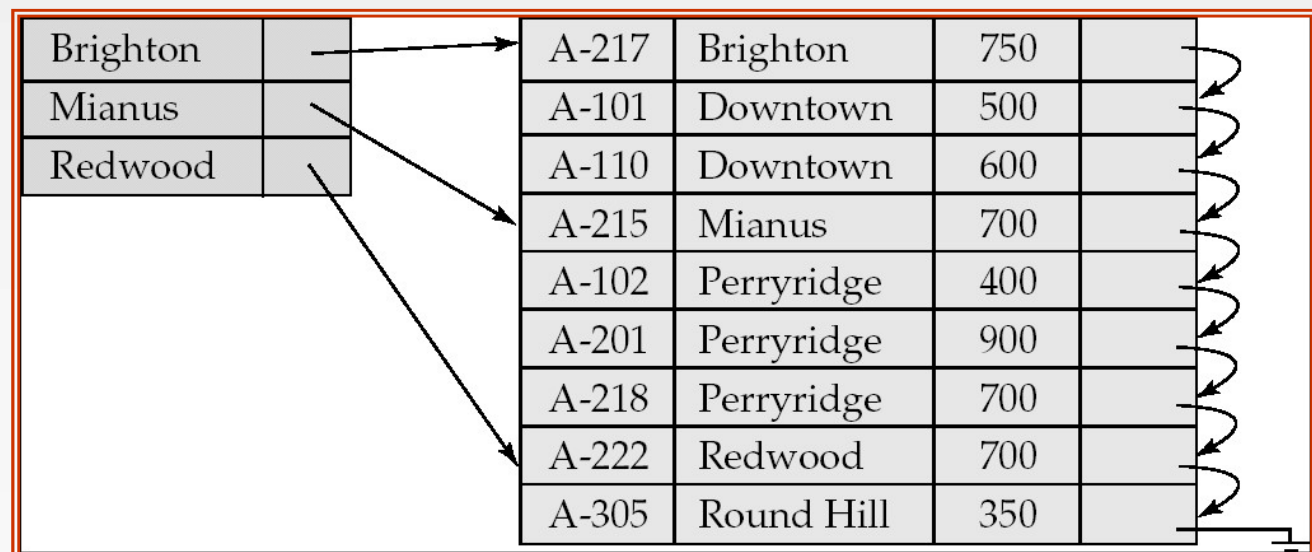
Multilevel Index (Cont.)





Index Update: Record 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: similar to file record deletion.
 - ❑ **Sparse indices** –
 - ▶ if deleted key value exists in the index, the value is replaced by 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.





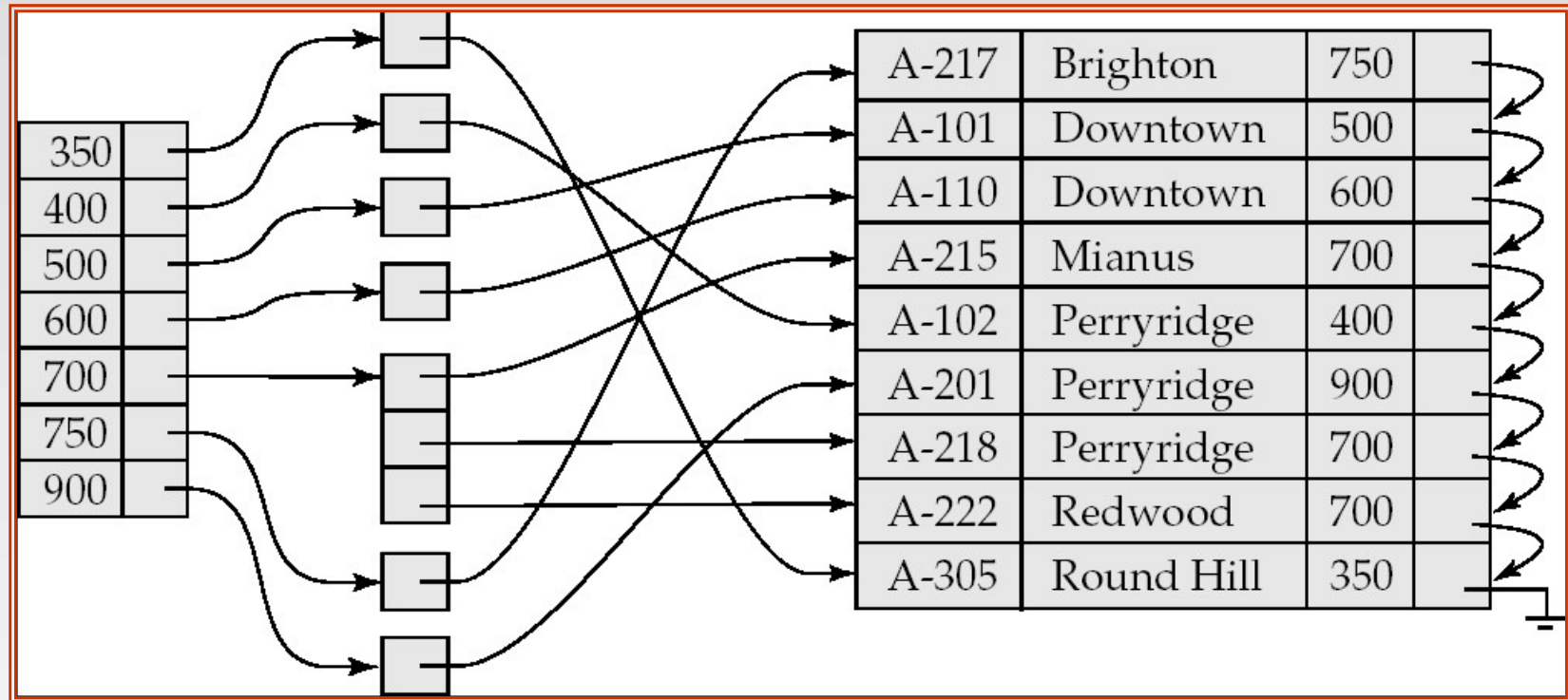
Index Update: Record Insertion

- Single-level index insertion:
 - Perform a lookup using the key value from inserted record
 - **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 Example



Secondary index on *balance* field of *account*

- ❑ 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 micro seconds, 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

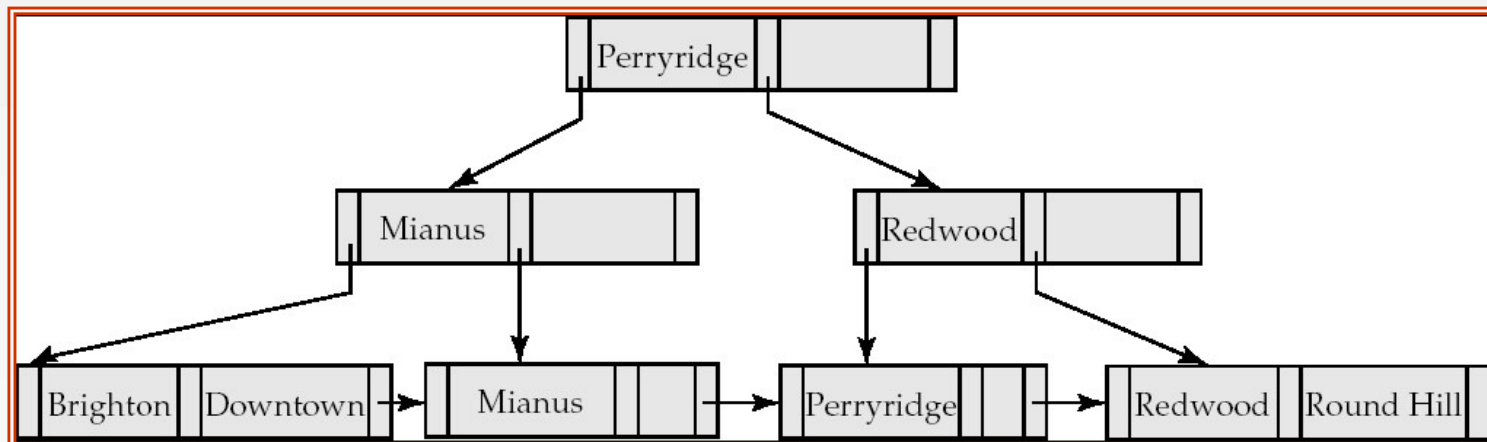




B⁺-Tree Index Files (Cont.)

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
- 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

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

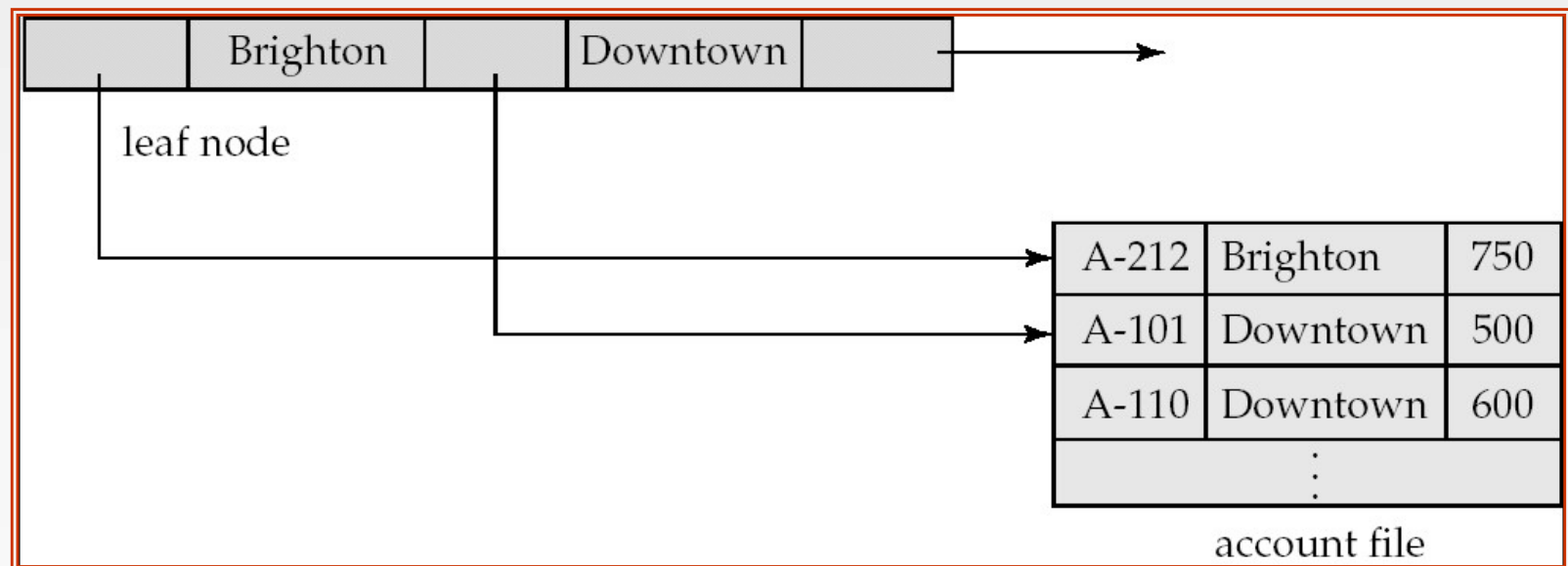




Leaf Nodes in B⁺-Trees

Properties of a leaf node:

- For $i = 1, 2, \dots, n-1$, pointer P_i 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





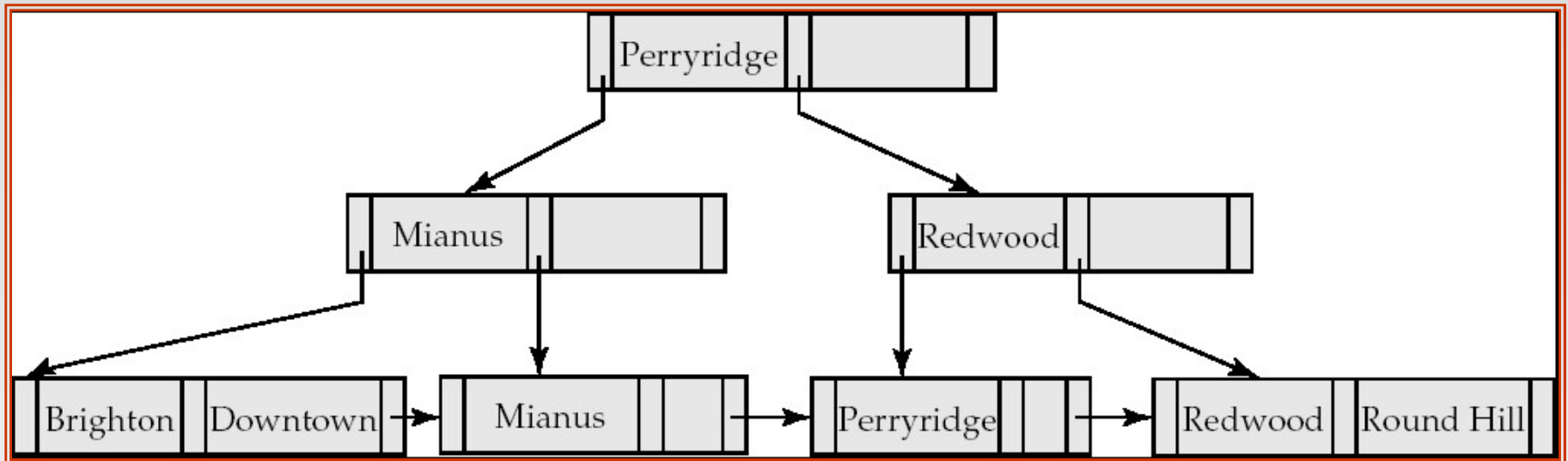
Non-Leaf Nodes in B⁺-Trees

- 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
 - All the search-keys in the subtree to which P_n points have values greater than or equal to K_{n-1}





Example of a B⁺-tree

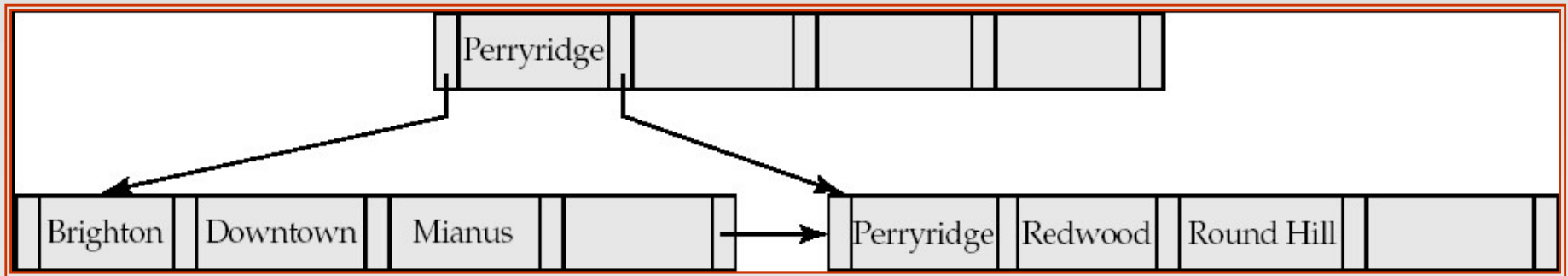


B⁺-tree for *account* file ($n = 3$)





Example of B⁺-tree



B⁺-tree for *account* file ($n = 5$)

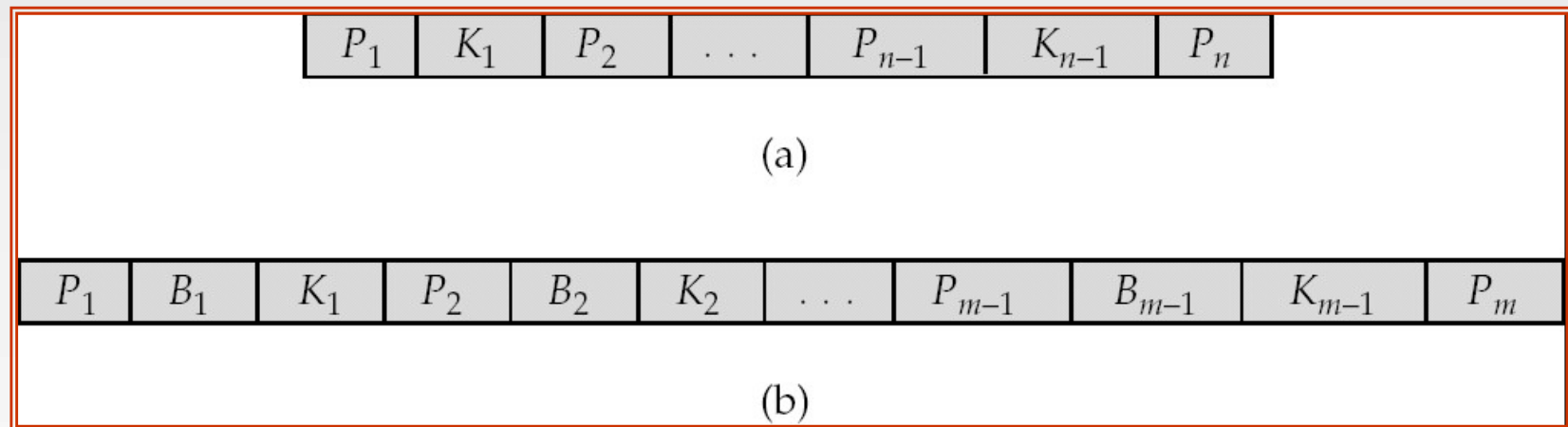
- ❑ Leaf nodes must have between 2 and 4 values ($\lceil (n-1)/2 \rceil$ and $n-1$, with $n = 5$).
- ❑ Non-leaf nodes other than root must have between 3 and 5 children ($\lceil n/2 \rceil$ and n with $n = 5$).
- ❑ Root must have at least 2 children.





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

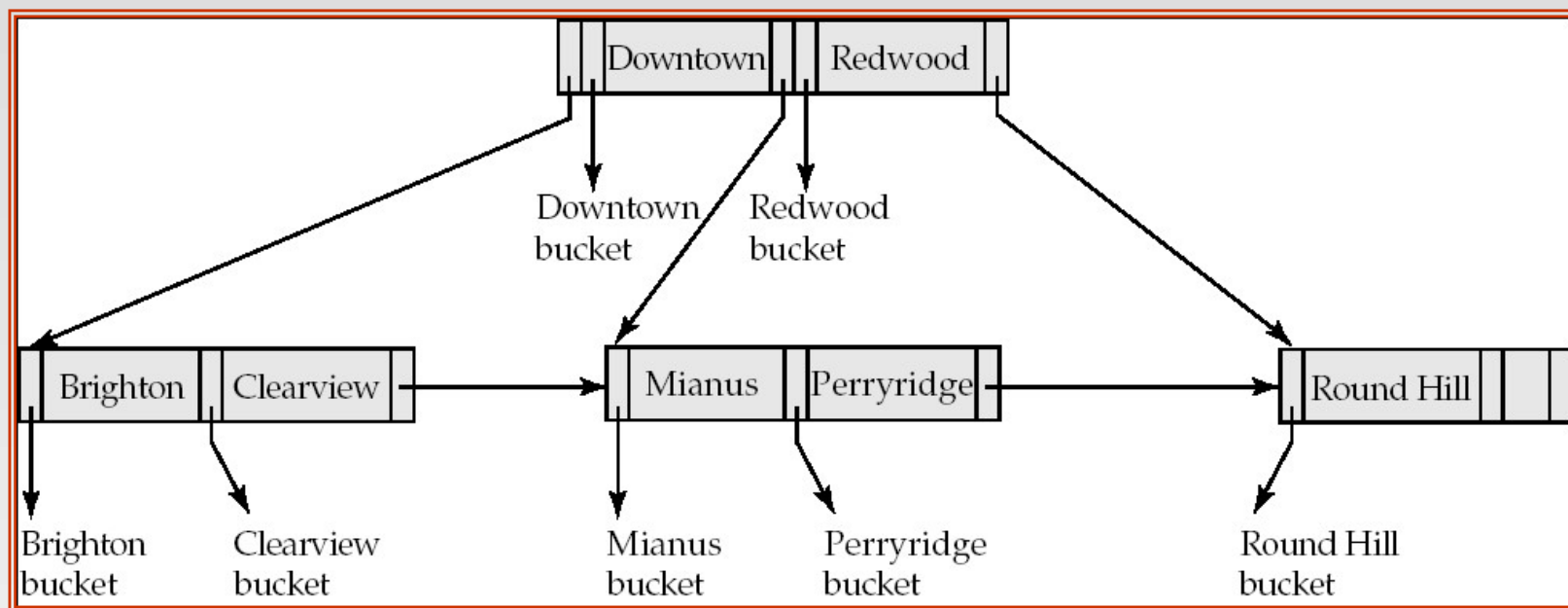


- Nonleaf node – pointers B_i are the bucket or file record pointers.

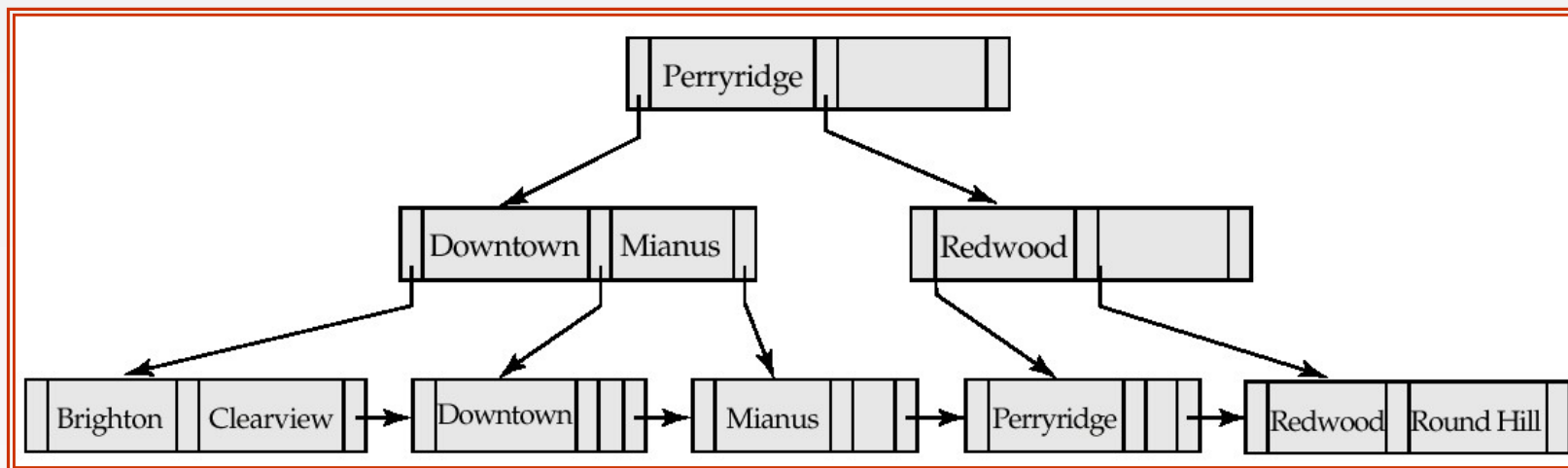




B-Tree Index File Example



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





B-Tree Index Files (Cont.)

- Advantages of B-Tree indices:
 - May use less tree nodes than a corresponding B⁺-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 branch name Perryridge; test *balance* = 1000
 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*).

- For
 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





Hashing

Database System Concepts, 5th Ed.

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





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

- 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
 - 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. .





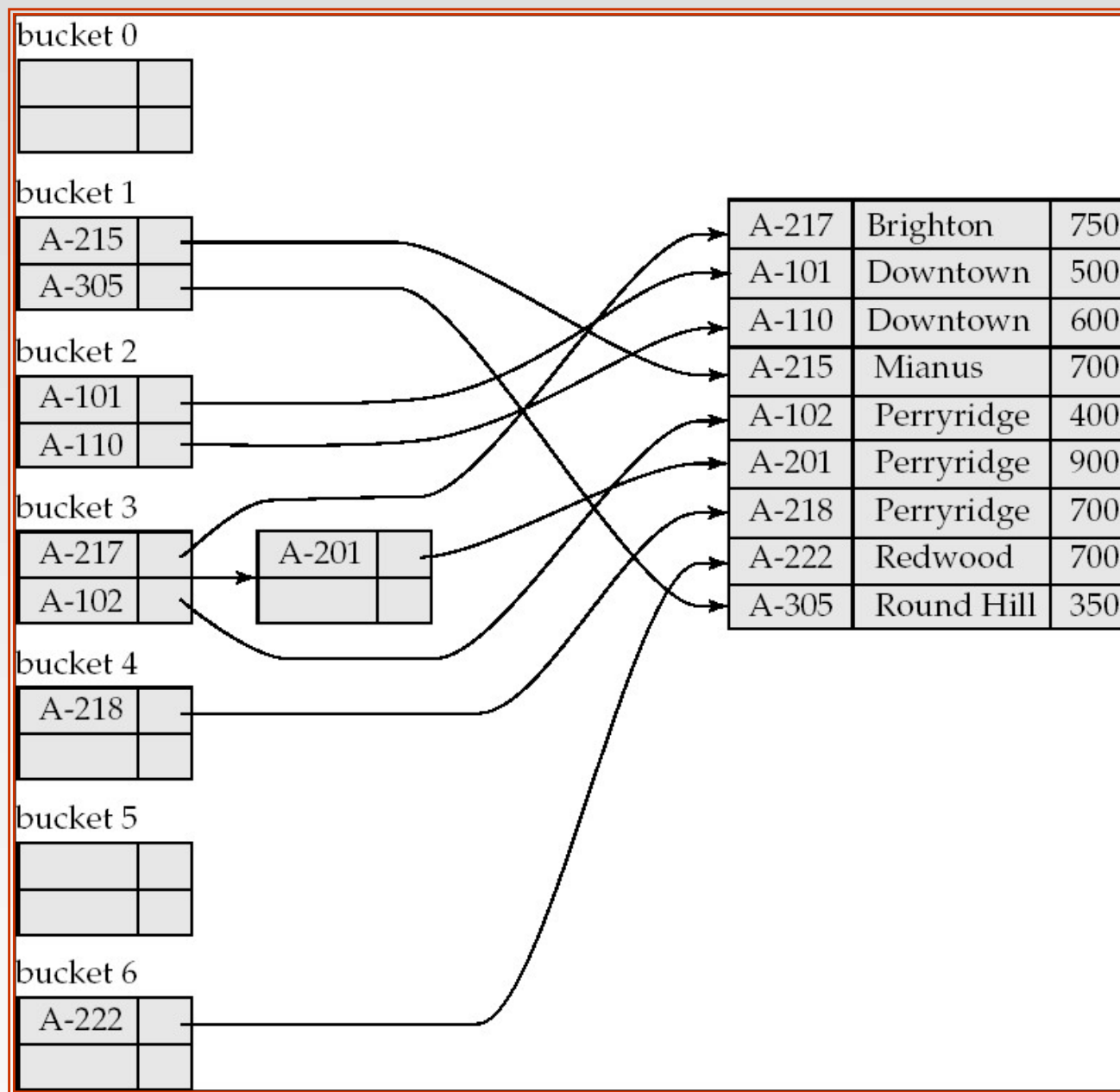
Hash Indices

- 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 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.





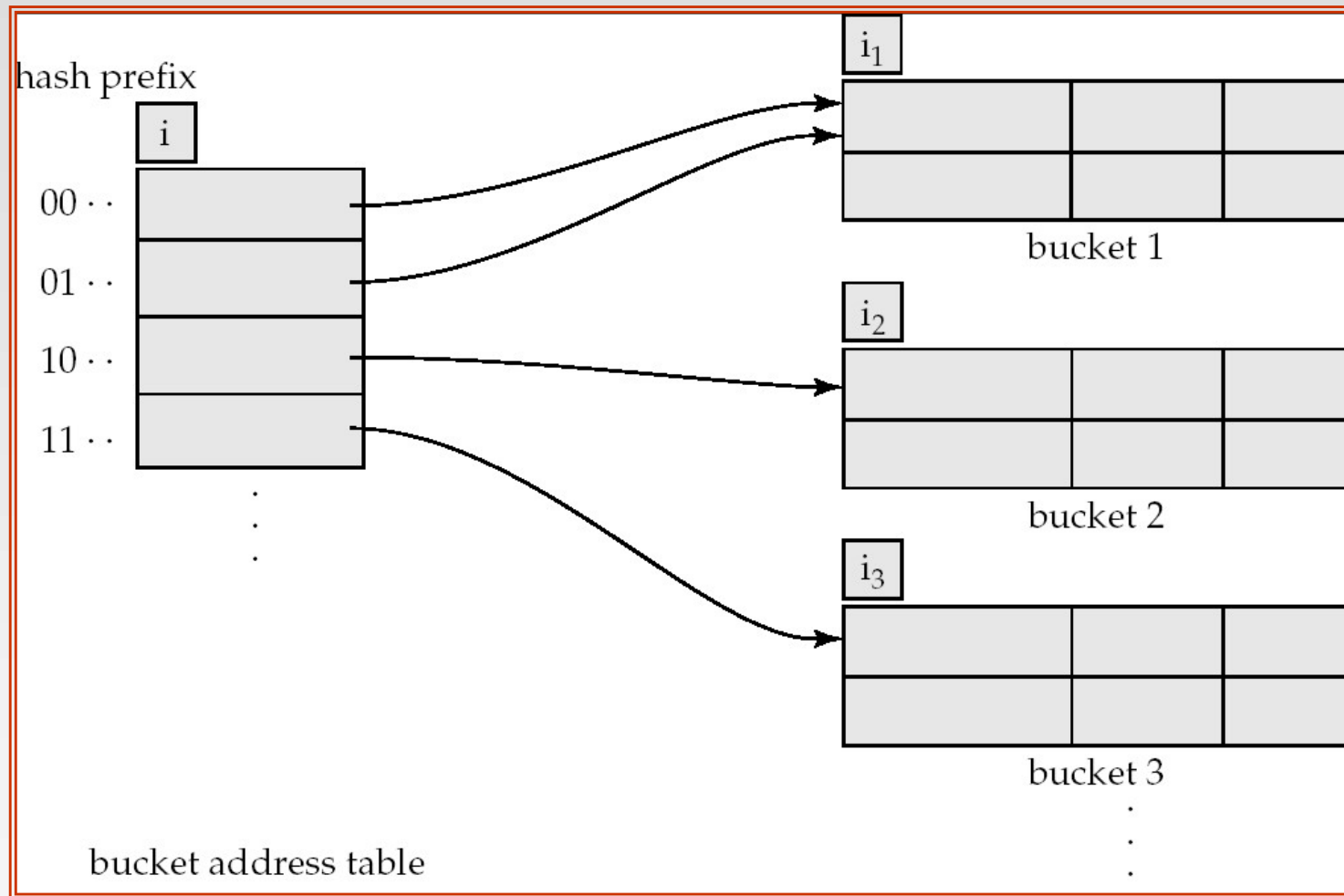
Dynamic Hashing

- 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 $i_1 = i - 1$ (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)





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 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)
 - 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

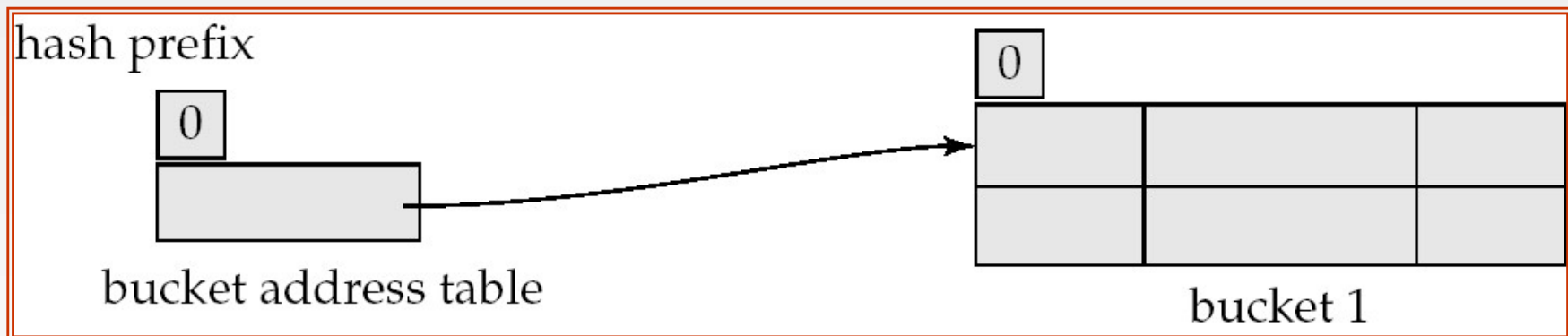
- To delete a key value,
 - locate it in its bucket and remove it.
 - The bucket itself can be removed if it becomes empty (with appropriate updates to the bucket address table).
 - Coalescing of buckets can be done (can coalesce only with a “*buddy*” bucket having same value of i_j and same $i_j - 1$ 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

<i>branch_name</i>	<i>h(branch_name)</i>
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



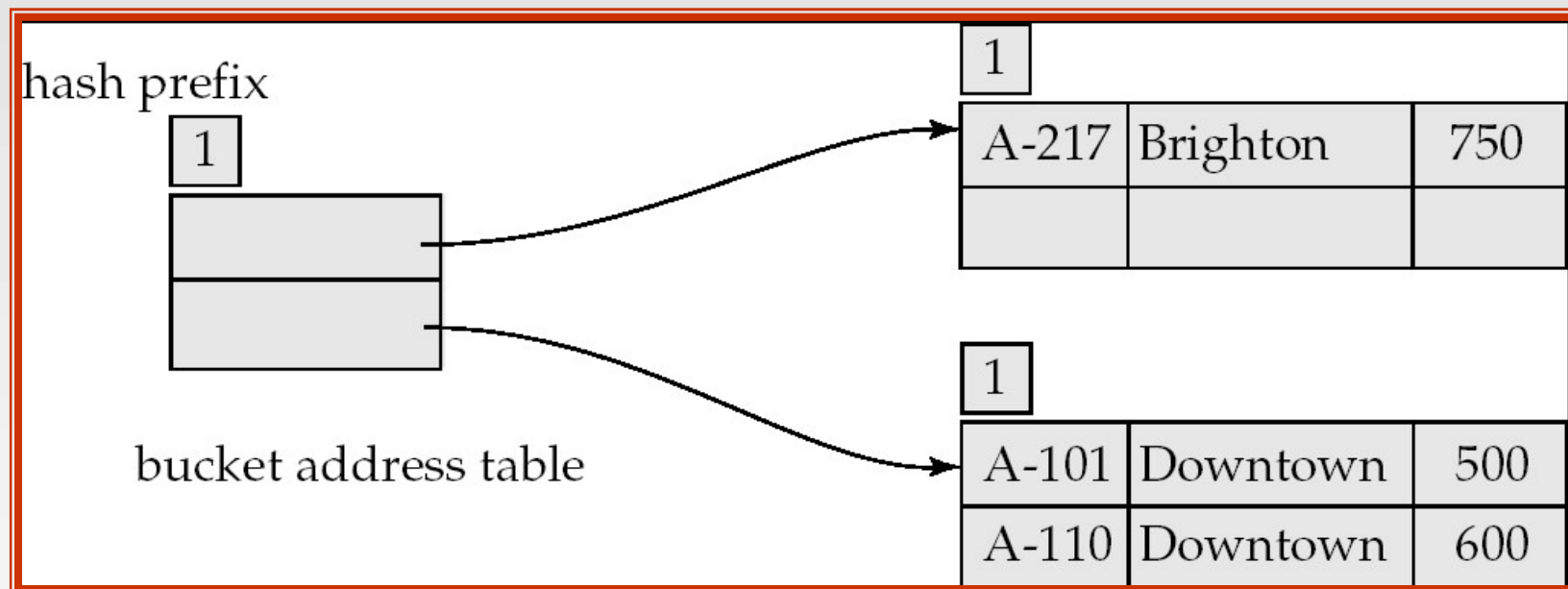
Initial Hash structure, bucket size = 2





Example (Cont.)

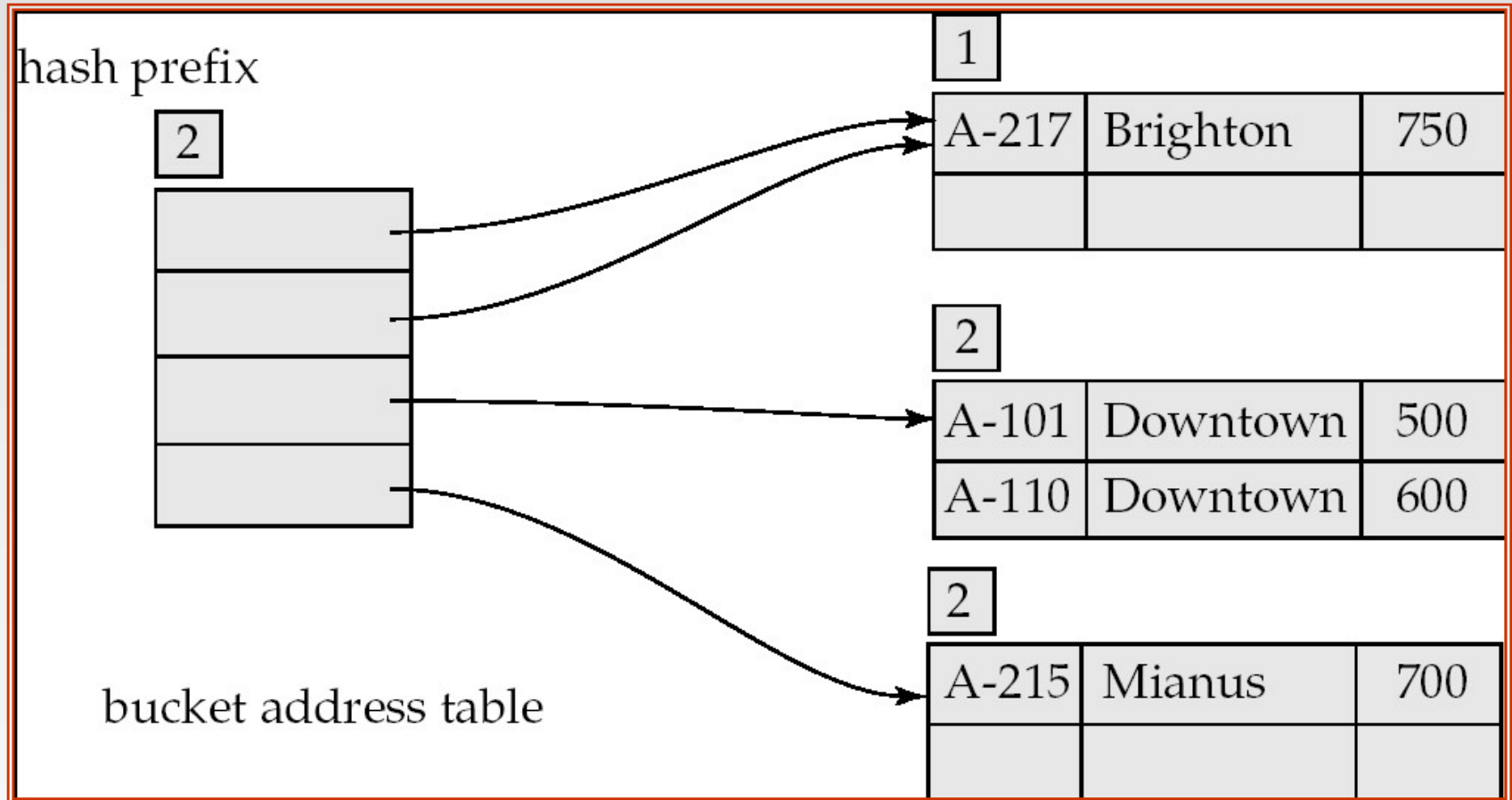
- Hash structure after insertion of one Brighton and two Downtown records





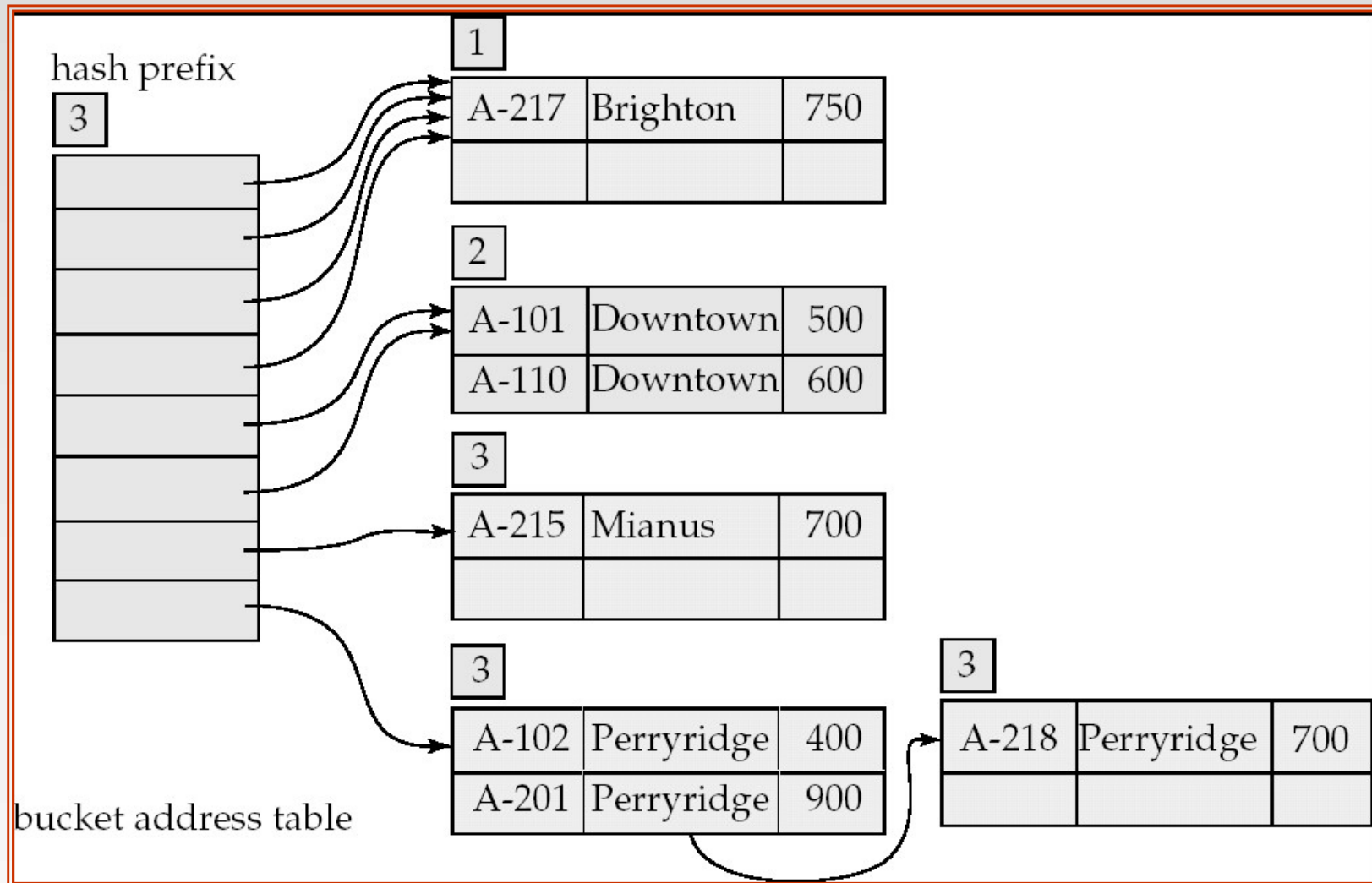
Example (Cont.)

Hash structure after insertion of Mianus record





Example (Cont.)



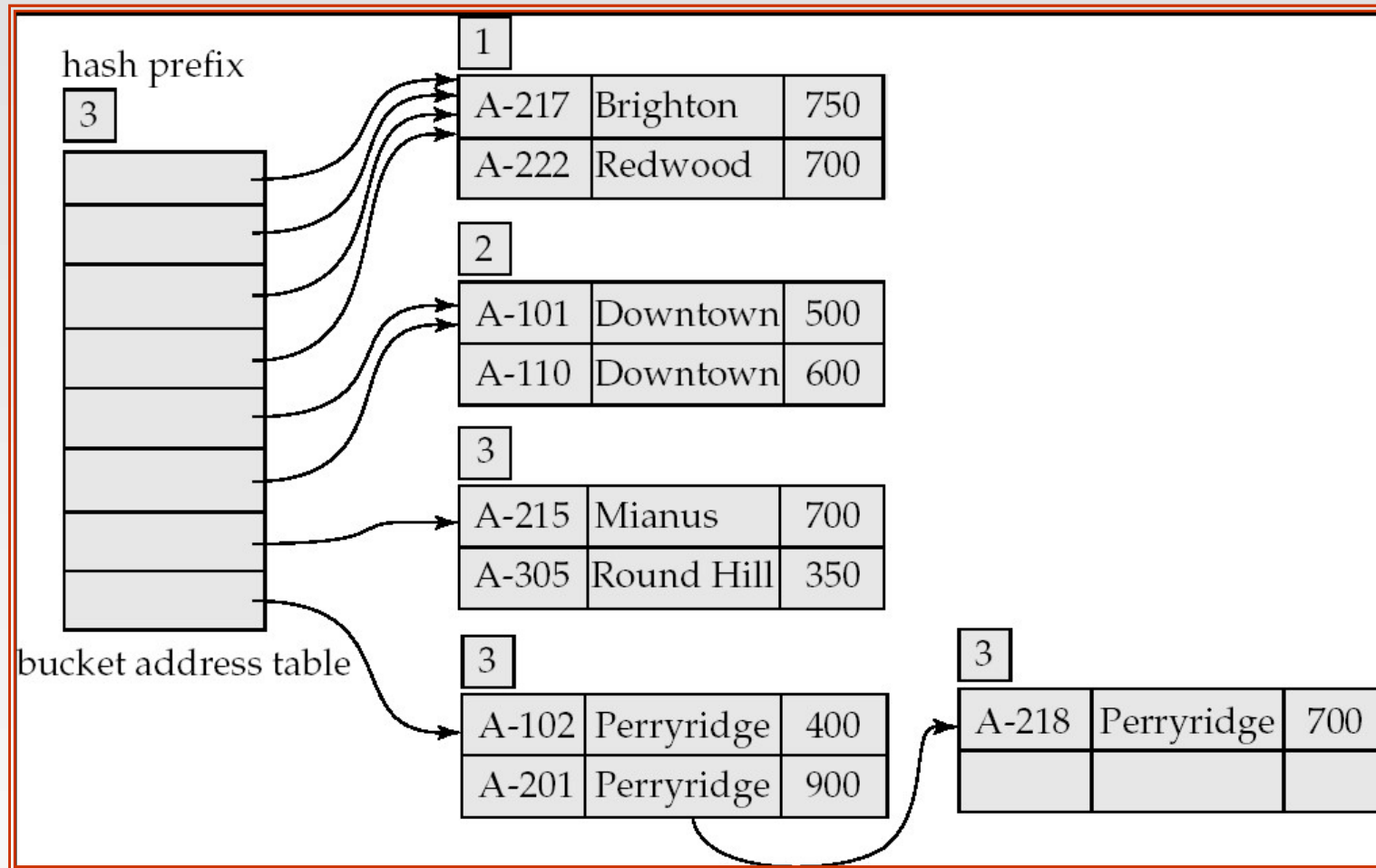
Hash structure after insertion of three Perryridge records





Example (Cont.)

- Hash structure after insertion of Redwood and Round Hill records





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** *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>

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

