



Dept. of Computer Science and Engineering
University of Rajshahi
www.ru.ac.bd

Dr. Shamim Ahmad

Chapter 12: Indexing and Hashing

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.2

Chapter 12: Indexing and Hashing

- Basic Concepts
- Ordered Indices
- B+-Tree Index Files
 - B-Tree Index Files
- Hashing
 - Static Hashing
 - Dynamic Hashing
- Comparison of Ordered Indexing and Hashing
- Multiple-Key Access and Bitmap indices
- Index Definition in SQL
- Indexing in Oracle 10g

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.3

Basic Concepts

- Indexing mechanisms are used to speed up access to desired data.
- **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".

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.4

Index Evaluation Metrics

- Access time
- Insertion time
- Deletion time
- Space overhead
- 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.
 - This strongly influences the choice of index, and depends on usage.

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.

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.
 - **Only** 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	→

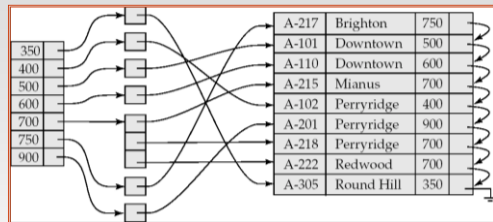
Index Update: Insertion

- 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

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, since the file is not sorted by the search key.

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 Node Structure

■ Typical node

P_1	K_1	P_2	\dots	P_{n-1}	K_{n-1}	P_n
-------	-------	-------	---------	-----------	-----------	-------

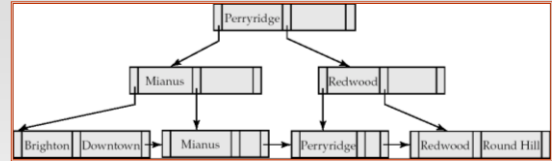
- K_i are the search-key values
- P_i are pointers to children (for non-leaf nodes) or pointers to records or buckets of records (for leaf nodes).

■ The search-keys in a node are ordered

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

■ Usually the size of a node is that of a block

Example of a B+-tree



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

B+-Tree Index File

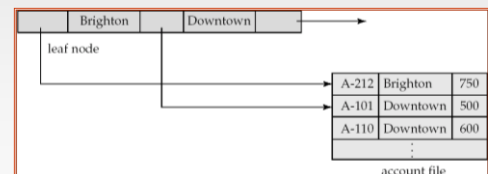
A B+-tree is a rooted tree satisfying the following properties:

- 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.

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

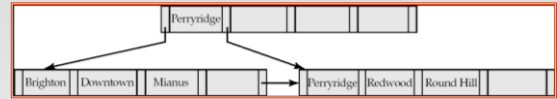


Non-Leaf Nodes in B+-Trees

- 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 $2 \leq i \leq m-1$, all the search-keys in the subtree to which P_i points have values greater than or equal to K_{i-1} and less than K_i
 - All the search-keys in the subtree to which P_m points have values greater than or equal to K_{m-1}



Example of B+-tree



B+-tree for accountfile ($n=5$)

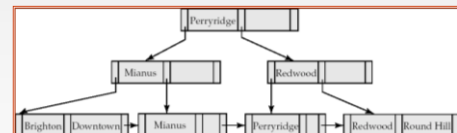
- Leaf nodes must have between 2 and 4 values ($\lceil (n-1)/2 \rceil$ and $n-1$, with $n=5$).
- Non-leaf nodes other than root must have between 3 and 5 children ($\lceil n/2 \rceil$ and n with $n=5$).
- Root must have at least 2 children.

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
 - Level below root has at least $2 * \lceil n/2 \rceil$ values
 - Next level has at least $2 * \lceil n/2 \rceil * \lceil n/2 \rceil$ values
 - .. etc.
- If there are K search-key values in the file, the tree height is no more than $\lceil \log_{n/2}(K) \rceil$
- thus searches can be conducted efficiently.
- Insertions and deletions to the main file can be handled efficiently, as the index can be restructured in logarithmic time (as we shall see some details, and more in the book).


Queries on B+-Trees

- Find all records with a search-key value of k .
 - $N = \text{root}$
 - Repeat
 - Examine N for the smallest search-key value $> k$.
 - If such a value exists, assume it is K_i . Then set $N = P_i$
 - Otherwise $k \geq K_{m-1}$. Set $N = P_m$
 Until N is a leaf node
 - If for some i , key $K_i = k$ follow pointer P_i to the desired record or bucket.
 - Else no record with search-key value k exists.



Queries on B+-Trees (Cont.)

- If there are K search-key values in the file, the height of the tree is no more than $\lceil \log_{n/2}(K) \rceil$.
- A node is generally the same size as a disk block, typically 4Kbytes
 - 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.
 - ▶ I.e. at most 4 accesses to disk blocks are needed
- 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

 This image cannot currently be displayed.

Hashing

José Aferees

Versão modificada de Database System Concepts, 5th Ed.
©Silberschatz, Korth and Sudarshan

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

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

Example of Hash File Organization

Hash file organization of *account* file, using *branch_name* as key

bucket 0				bucket 5	A-100	Perryridge	400
					A-200	Perryridge	600
					A-210	Perryridge	700
bucket 1				bucket 6			
bucket 2				bucket 7	A-215	Miami	700
bucket 3	A-210	Brighton	750	bucket 8	A-210	Downtown	500
	A-305	Round Hill	350		A-110	Downtown	600
bucket 4	A-222	Redwood	700	bucket 9			

- 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. $h(\text{Perryridge}) = 5$ $h(\text{Round Hill}) = 3$ $h(\text{Brighton}) = 3$

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.25

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.

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.26

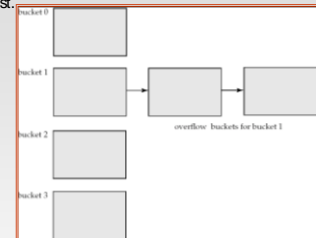
Handling of Bucket Overflows

- 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**.

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.27

Handling of Bucket Overflows (Cont.)

- **Overflow chaining** – the overflow buckets of a given bucket are chained together in a linked list.



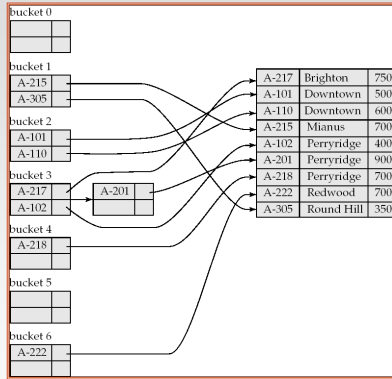
- Above scheme is called **closed hashing**.
 - An alternative, called **open hashing**, which does not use overflow buckets, is not suitable for database applications.

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.28

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.



José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.29

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.

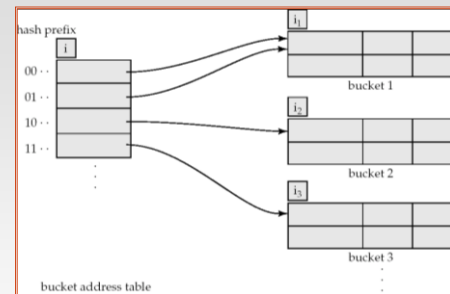
José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.30

Dynamic Hashing

- Good for database that grows and shrinks in size
- Allow 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 $b = 32$ (Note that 2^{32} is quite large!)
 - 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 same bucket. Thus, actual number of buckets is $< 2^i$
 - ▶ The number of buckets also changes dynamically due to coalescing and splitting of buckets.

José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.31

General Extendable Hash Structure



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

José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.32

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
 - ▶ Overflow buckets used instead in some cases

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.33

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.

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.34

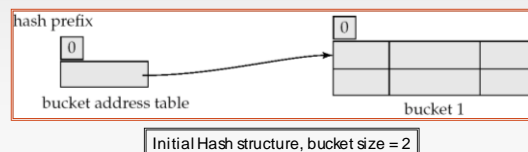
Deletion in Extendable Hash Structure

- 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

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.35

Use of Extendable Hash Structure: Example

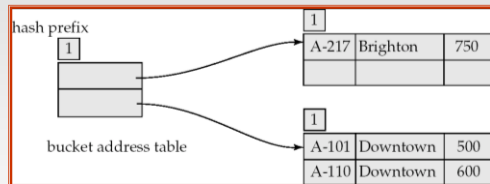
branch_name	$h(\text{branch_name})$
Brighton	0010 1101 1111 1011 0010 1100 0011 0000
Downtown	1010 0011 1010 0000 1100 0110 1001 1111
Mianus	1100 0111 1110 1101 1011 1111 0011 1010
Perryridge	1111 0001 0010 0100 1001 0011 0110 1101
Redwood	0011 0101 1010 0110 1100 1001 1110 1011
Round Hill	1101 1000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



José Aferes - Adaptado de Database System Concepts - 5th Edition 12.36

Example (Cont.)

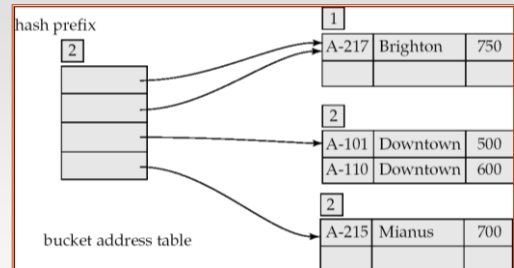
- Hash structure after insertion of one Brighton and two Downtown records



José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.37

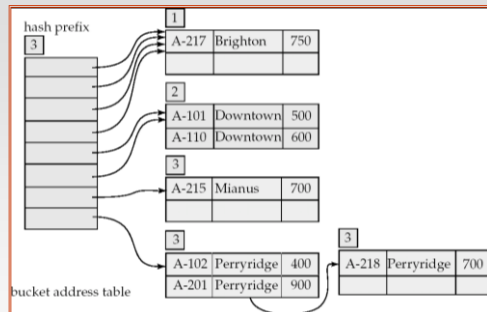
Example (Cont.)

Hash structure after insertion of Mianus record



José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.38

Example (Cont.)

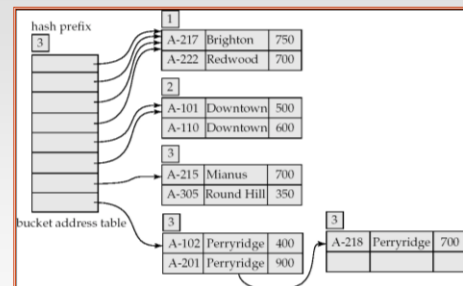


Hash structure after insertion of three Perryridge records

José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.39

Example (Cont.)

- Hash structure after insertion of Redwood and Round Hill records



José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.40

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 area 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

José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.41

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
 - ↳ Consider e.g. query with *where* $A \geq v_1$ *and* $A \leq v_2$
- In practice:
 - PostgreSQL supports hash indices, but discourages use due to poor performance
 - Oracle supports static hash organization, but not hash indices
 - SQL Server supports only B*-trees

José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.42

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
 - 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

José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.43

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	name	gender	address	income_level	Bitmaps for gender		Bitmaps for income_level	
					m	f	L1	L2
0	John	m	Perryridge	L1	1	0	1	0
1	Diana	f	Brooklyn	L2	0	1	0	1
2	Mary	f	Jonestown	L1	0	1	0	0
3	Peter	m	Brooklyn	L4	1	0	0	0
4	Kathy	f	Perryridge	L3	0	1	0	0

José Aferes - Adaptado de Database System Concepts - 5ª Edition 12.44

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)
 - Example query with *where gender='m' and income_level='L1'*
- 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
 - It doesn't even require accessing the file!

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.45

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
 - ▶ $\text{not}(A=v): (\text{NOT bitmap-}A-v) \text{ AND ExistenceBitmap}$
- Should keep bitmaps for all values, even null value
 - To correctly handle SQL null semantics for $\text{NOT}(A=v)$:
 - ▶ intersect above result with $(\text{NOT bitmap-}A\text{-Null})$

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.46

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 and-ed with just 31,250 instructions
- 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

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.47

Index Definition in SQL standard

- 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.

José Aferes - Adaptado de Database System Concepts - 5th Edition 12.48

Indexing in Oracle

- Oracle supports B+-Tree indices as a default for the create index SQL command
- A new non-null attribute *row-id* is added to all indices, so as to guarantee that all search keys are unique.
 - indices are supported on
 - attributes, and attribute lists,
 - on results of function over attributes
 - or using structures external to Oracle (*Domain indices*)
- Bitmap indices are also supported, but for that an explicit declaration is needed:
create bitmap index <index-name>
on <relation-name> (<attribute-list>)

José Añeres - Adaptado de Database System Concepts - 5ª Edition 12.49

Hashing in Oracle

- Hash indices are not supported
- However (limited) static hash file organization is supported for partitions
create table ... partition by hash(<attribute-list>)
partitions <N>
stored in (<tables>)
- Index files can also be partitioned using hash function
create index ... global partition by hash(<attribute-list>)
partitions <N>
 - This creates a global index partitioned by the hash function
- (Global) indexing over hash partitioned table is also possible
- Hashing may also be used to organize clusters in multi-table clusters

José Añeres - Adaptado de Database System Concepts - 5ª Edition 12.50