Friday, April 09, 2004 11:33 AM

Hash Indexing

Disadvantages of Sequential File Organization Must use an index and/or binary search to locate data

File organization based on hashing allow us to avoid accessing an index structure.

Also provides a way to construct indices.

- A bucket is a drame 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.
- Include the second s



bucket 0 bucket 5

Lectures Desktop (C) Page 1

Bank Account Example

- 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
 - \succ E.g. h(Perryridge) = 5 h(Round Hill) = 3 h(Brighton) = 3

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			bucket 8		
A-217	Brighton	750	A-101	Downtown	500
A-305	Round Hill	350	A-110	Downtown	600
bucket 4			bucket 9		
A-222	Redwood	700			

Hash Functions

- Worst has 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.

- Key = 'x₁ x₂ ... x_n' *n* byte character string
- Have *b* buckets
- h: add x₁ + x₂ + x_n
 - compute sum modulo *b*

Insert Delete h(e) = 1 f c

Example with 2 records per bucket



Bucket Overflow

- Bucket overflow can occur because of
 - Insufficient buckets
 - Skew in distribution of records. This can occur due to two reasons:
 - ★ multiple records have same search-key value
 - ★ chosen hash function produces non-uniform distribution of key values
- Although the probability of bucket overflow can be reduced, it cannot be eliminated; it is handled by using overflow buckets.
- Overflow chaining the overflow buckets of a given bucket are chained together in a linked list.
- Above scheme is called closed hashing.
 - An alternative, called open hashing, which does not use overflow buckets, is not suitable for database applications.



organization may already be sequential?

Use 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.

Hash Index Example



- 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.

Thus, the topic of dynamic hashing ...

- Good for database that grows and shrinks in size
- Allows the hash function to be modified dynamically
- Extendable (aka extensible) 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, $0 \le i \le 32$.
 - > 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.

So two ideas combine in extendable hashing:

(a) Use *i* of *b* bits output by hash function

$$h(K) \rightarrow 00110101$$

use $i \rightarrow$ grows over time....

(b) Use directory



Example: h(k) is 4 bits; 2 keys per bucket



Buckets



Lectures Desktop (C) Page 8







Extensible hashing: deletion

- No merging of blocks
- Merge blocks
 and cut directory if possible
 (Reverse insert procedure)