

Chapter 10: Storage and File Structure

Database System Concepts, 6th Ed.

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



Magnetic Hard Disk Mechanism



NOTE: Diagram is schematic, and simplifies the structure of actual disk drives



Optimization of Disk-Block Access

Block – a contiguous sequence of sectors from a single track

- data is transferred between disk and main memory in blocks
- sizes range from 512 bytes to several kilobytes
 - Smaller blocks: more transfers from disk
 - Larger blocks: more space wasted due to partially filled blocks
 - Typical block sizes today range from 4 to 16 kilobytes

Disk-arm-scheduling algorithms order pending accesses to tracks so that disk arm movement is minimized

elevator algorithm:



Inner track

Optimization of Disk Block Access (Cont.)

- File organization optimize block access time by organizing the blocks to correspond to how data will be accessed
 - E.g. Store related information on the same or nearby cylinders.
 - Files may get **fragmented** over time
 - E.g. if data is inserted to/deleted from the file
 - Or free blocks on disk are scattered, and newly created file has its blocks scattered over the disk
 - Sequential access to a fragmented file results in increased disk arm movement
 - Some systems have utilities to defragment the file system, in order to speed up file access



File Organization

- The database is stored as a collection of *files*. Each file is a sequence of *records*. A record is a sequence of fields.
- One approach:

assume record size is fixed

- each file has records of one particular type only
- •different files are used for different relations

This case is easiest to implement; will consider variable length records later.

columns



Fixed-Length Records

- Simple approach:
 - Store record *i* starting from byte *n* ★ (*i* 1), where *n* is the size of each record.
 - Record access is simple but records may cross blocks
 - Modification: do not allow records to cross block boundaries
- Deletion of record i: alternatives:

 - move record n to i
 - do not move records, but link all free records on a free list

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000
record 11	98345	Kim	Elec. Eng.	80000



Deleting record 3 and compacting

record record

0	10101	Srinivasan	Comp. Sci.	65000
1	12121	Wu	Finance	90000
2	15151	Mozart	Music	40000
4	32343	El Said	History	60000
5	33456	Gold	Physics	87000
6	45565	Katz	Comp. Sci.	75000
7	58583	Califieri	History	62000
8	76543	Singh	Finance	80000
9	76766	Crick	Biology	72000
10	83821	Brandt	Comp. Sci.	92000
11	98345	Kim	Elec. Eng.	80000



Deleting record 3 and moving last record

record 0
record 1
record 2
record 11
record 4
record 5
record 6
record 7
record 8
record 9
record 10

10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
98345	Kim	Elec. Eng.	80000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
	10101 12121 15151 98345 32343 32343 33456 45565 58583 76543 76766 83821	10101 Srinivasan 12121 Wu 15151 Mozart 98345 Kim 32343 El Said 33456 Gold 45565 Katz 58583 Califieri 76543 Singh 76766 Crick 83821 Brandt	10101SrinivasanComp. Sci.12121WuFinance15151MozartMusic98345KimElec. Eng.32343El SaidHistory33456GoldPhysics45565KatzComp. Sci.58583CalifieriHistory76543SinghFinance76766CrickBiology83821BrandtComp. Sci.



Free Lists

- Store the address of the first deleted record in the file header.
- Use this first record to store the address of the second deleted record, and so on
- Can think of these stored addresses as pointers since they "point" to the location of a record.
- More space efficient representation: reuse space for normal attributes of free records to store pointers. (No pointers stored in in-use records.)

header				```	Ra 1
record 0	10101	Srinivasan	Comp. Sci.	65000	
record 1				<u>a</u>	$\rightarrow 4$
record 2	15151	Mozart	Music	40000	
record 3	22222	Einstein	Physics	95000	
 , r ecord 4					
record 5	33456	Gold	Physics	87000	\sum
record 6				×	
record 7	58583	Califieri	History	62000	
record 8	76543	Singh	Finance	80000	
record 9	76766	Crick	Biology	72000	
record 10	83821	Brandt	Comp. Sci.	92000	
record 11	98345	Kim	Elec. Eng.	80000	

Database System Concepts - 6th Edition

©Silberschatz, Korth and Sudarshan



Variable-Length Records

- Variable-length records arise in database systems in several ways:
 - Storage of multiple record types in a file.
 - Record types that allow variable lengths for one or more fields such as strings (varchar)
 - Record types that allow repeating fields (used in some older data models).
- Attributes are stored in order
- Variable length attributes represented by fixed size (offset, length), with actual data stored after all fixed length attributes
- Null values represented by null-value bitmap



©Silberschatz, Korth and Sudarshan







Variable-Length Records: Slotted Page Structure



End of Free Space

- Slotted page header contains:
 - number of record entries
 - end of free space in the block
 - Iocation and size of each record
- Records can be moved around within a page to keep them contiguous with no empty space between them; entry in the header must be updated.
- Pointers should not point directly to record instead they should point to the entry for the record in header.



Organization of Records in Files

- Heap a record can be placed anywhere in the file where there is space
- Sequential store records in sequential order, based on the value of the search key of each record
- Hashing a hash function computed on some attribute of each record; the result specifies in which block of the file the record should be placed
- Records of each relation may be stored in a separate file. In a multitable clustering file organization records of several different relations can be stored in the same file
 - Motivation: store related records on the same block to minimize I/O



Sequential File Organization

- Suitable for applications that require sequential processing of the entire file
- The records in the file are ordered by a search-key

10101	Srinivasan	Comp. Sci.	65000	_	$ \rightarrow $
12121	Wu	Finance	90000	_	\checkmark
15151	Mozart	Music	40000	_	\checkmark
22222	Einstein	Physics	95000		\leq
32343	El Said	History	60000		\leq
33456	Gold	Physics	87000	l,	\leq
45565	Katz	Comp. Sci.	75000	I,	\leq
58583	Califieri	History	62000	I,	\leq
76543	Singh	Finance	80000	I.	\leq
76766	Crick	Biology	72000	I,	\leq
83821	Brandt	Comp. Sci.	92000		\leq
98345	Kim	Elec. Eng.	80000		

Sequential File Organization (Cont.)

- Deletion use pointer chains
- Insertion –locate the position where the record is to be inserted
 - if there is free space insert there
 - if no free space, insert the record in an overflow block
 - In either case, pointer chain must be updated
- Need to reorganize the file from time to time to restore sequential order

10101	Srinivasan	Comp. Sci.	65000	
12121	Wu	Finance	90000	
15151	Mozart	Music	40000	
22222	Einstein	Physics	95000	
32343	El Said	History	60000	
33456	Gold	Physics	87000	
45565	Katz	Comp. Sci.	75000	
58583	Califieri	History	62000	
76543	Singh	Finance	80000	
76766	Crick	Biology	72000	
83821	Brandt	Comp. Sci.	92000	
98345	Kim	Elec. Eng.	80000	
				/
22222	\$7.1	30.	40000	
32222	Verdi	Music	48000	



Multitable Clustering File Organization

Store several relations in one file using a **multitable clustering** file organization

	dept_	name	building		b	udget	
department	department Comp. Sci. Physics		Ta W	Taylor Watson		100000 70000	
	ID	name		dept_name		salary	
instructor	10101 33456 45565 83821	Srinivasaı Gold Katz Brandt	ı	Comp. Sci. Physics Comp. Sci. Comp. Sci.		65000 87000 75000 92000	
multitable clustering of <i>department</i> and <i>instructor</i>	Comp. Sci. 45564 10101 83821 Physics 22456			Faylor Katz Brinivasan Brandt Watson Gold		100000 75000 65000 92000 70000 87000	

Multitable Clustering File Organization (cont.)

- good for queries involving *department* \bowtie *instructor*, and for queries involving one single department and its instructors
- bad for queries involving only department
- results in variable size records
- Can add pointer chains to link records of a particular relation

Comp. Sci.	Taylor	100000	
45564	Katz	75000	
10101	Srinivasan	65000	
83821	Brandt	92000	
Physics	Watson	70000	
33456	Gold	87000	



Data Dictionary Storage

The **Data dictionary** (also called **system catalog**) stores **metadata**; that is, data about data, such as

- Information about relations
 - names of relations
 - names, types and lengths of attributes of each relation
 - names and definitions of views
 - integrity constraints
 - User and accounting information, including passwords
- Statistical and descriptive data
 - number of tuples in each relation
- Physical file organization information
 - How relation is stored (sequential/hash/...)
 - Physical location of relation
- Information about indices (Chapter 11)



Relational Representation of System Metadata

- Relational representation on disk
- Specialized data structures designed for efficient access, in memory

<u>relation_name</u>
number_of_attributes
storage_organization
location

Relation metadata

Index_metadata <u>index_name</u> <u>relation_name</u> index_type index_attributes

View_metadata <u>view_name</u> definition Attribute_metadata <u>relation_name</u> <u>attribute_name</u> domain_type position length

User_metadata

<u>user_name</u> encrypted_password group

Database System Concepts - 6th Edition



Storage Access

- A database file is partitioned into fixed-length storage units called blocks. Blocks are units of both storage allocation and data transfer.
- Database system seeks to minimize the number of block transfers between the disk and memory. We can reduce the number of disk accesses by keeping as many blocks as possible in main memory.
- Buffer portion of main memory available to store copies of disk blocks.
- Buffer manager subsystem responsible for allocating buffer space in main memory.



Buffer Manager

- Programs call on the buffer manager when they need a block from disk.
 - 1. If the block is already in the buffer, buffer manager returns the address of the block in main memory
 - 2. If the block is not in the buffer, the buffer manager
 - 1. Allocates space in the buffer for the block
 - 1. Replacing (throwing out) some other block, if required, to make space for the new block.
 - 2. Replaced block written back to disk only if it was modified since the most recent time that it was written to/fetched from the disk.
 - 2. Reads the block from the disk to the buffer, and returns the address of the block in main memory to requester.



Buffer-Replacement Policies

- Most operating systems replace the block least recently used (LRU strategy)
- Idea behind LRU use past pattern of block references as a predictor of future references
- Queries have well-defined access patterns (such as sequential scans), and a database system can use the information in a user's query to predict future references
 - LRU can be a bad strategy for certain access patterns involving repeated scans of data
 - For example: when computing the join of 2 relations r and s by a nested loops for each tuple *tr* of *r* do for each tuple *ts* of *s* do if the tuples *tr* and *ts* match ...
 - Mixed strategy with hints on replacement strategy provided by the query optimizer is preferable

Buffer-Replacement Policies (Cont.)

- Pinned block memory block that is not allowed to be written back to disk.
- Toss-immediate strategy frees the space occupied by a block as soon as the final tuple of that block has been processed
- Most recently used (MRU) strategy system must pin the block currently being processed. After the final tuple of that block has been processed, the block is unpinned, and it becomes the most recently used block.
- Buffer manager can use statistical information regarding the probability that a request will reference a particular relation
 - E.g., the data dictionary is frequently accessed. Heuristic: keep data-dictionary blocks in main memory buffer
- Buffer managers also support **forced output** of blocks for the purpose of recovery (more in Chapter 16)



Static Hashing

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



Example of Hash File Organization

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



Example of Hash File Organization

bucket 0

bucket 1

15151	Mozart	Music	40000

bucket 2

32343	El Said	History	80000
58583	Califieri	History	60000

bucket 3

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

bucket 4

12121	Wu	Finance	90000
76543	Singh	Finance	80000

bucket 5

76766	Crick	Biology	72000

bucket 6

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

bucket 7



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

Database System Concepts - 6th Edition

©Silberschatz, Korth and Sudarshan



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.



Database System Concepts - 6th Edition

©Silberschatz, Korth and Sudarshan



Hash Indices

- Hashing can be used not only for file organization, but also for indexstructure 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



Database System Concepts - 6th Edition



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, $0 \le i \le 32$.
 - 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 (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





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_i* bits.
- To locate the bucket containing search-key K_i :
 - 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_i
 - 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_i :

- If $i > i_i$ (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 $i = i_i$ (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_i Now $i > i_i$ so use the first case above.

Database System Concepts - 6th Edition

Deletion in Extendable Hash Structure

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



Use of Extendable Hash Structure: Example

dept_name

Biology Comp. Sci. Elec. Eng. Finance History Music Physics

h(*dept_name*)

0010 1101 1111 1011 0010 1100 0011 0000 1111 0001 0010 0100 1001 0011 0110 1101 0100 0011 1010 1100 1100 0110 1101 1111 1010 0011 1010 0000 1100 0110 1001 1111 1100 0111 1110 1101 1011 1111 0011 1010 0011 0101 1010 0110 1100 1001 1110 1011 1001 1000 0011 1111 1001 1100 0000 0001





Initial Hash structure; bucket size = 2



Database System Concepts - 6th Edition



Hash structure after insertion of "Mozart", "Srinivasan", and "Wu" records





Hash structure after insertion of Einstein record





Hash structure after insertion of Gold and El Said records





Example (Cont.)

Hash structure after insertion of Katz record







Database System Concepts - 6th Edition

©Silberschatz, Korth and Sudarshan

92000





Database System Concepts - 6th Edition

©Silberschatz, Korth and Sudarshan