

1 Memory and Disk

Whenever a database processes data, that data must exist in memory. Accessing this data is relatively fast, but once the data becomes very large, it becomes impossible to fit all of it within memory. Disks are used to cheaply store all of a database's data persistently but incur a large cost whenever data is accessed or new data is written.

2 Files, Pages, Records

The basic unit of data for relational databases is a **record** (row). These records are organized into **relations** (tables) and can be modified, deleted, searched, or created in memory.

The basic unit of data for disk is a **page**, the smallest unit of transfer from disk to memory and vice versa. In order to represent relational databases in a format compatible with disk, each relation is stored in its own file and its records are organized into pages in the file. Based on the relation's schema and access pattern, the database will determine: (1) type of file used, (2) how pages are organized in the file, (3) how records are organized on each page, (4) and how each record is formatted.

3 Choosing File Types

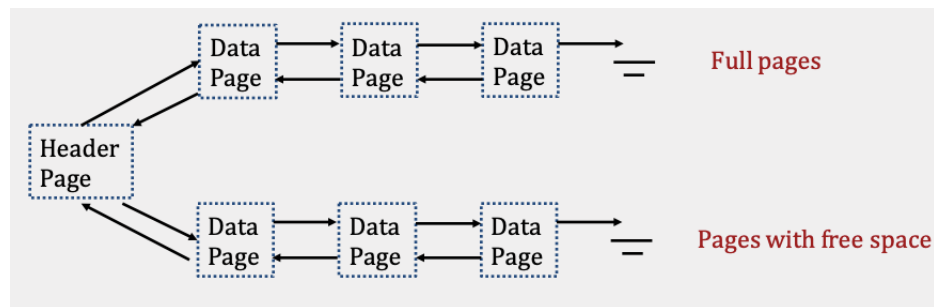
There are two main types of files: Heap Files and Sorted Files. The database chooses which file type to use for a relation based on the I/O cost associated with the relation's access pattern. **1 I/O** is equivalent to **1 page read from disk** or **1 page write to disk** and I/O calculations are made for each file type based on the insert, delete, and scan operations specified in the access pattern. The type that requires less I/Os is chosen.

4 Heap File

A **heap file** is a file type with no particular ordering of pages or records on pages and has two main implementations.

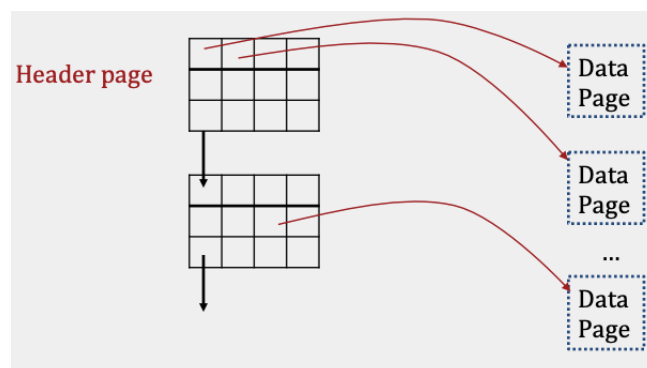
4.1 Linked List Implementation

In the linked list implementation, each data page contains **records**, a **free space tracker**, and **pointers** (byte offsets) to the next and previous page. There is one header page that serves as the start of the file and is used to separate the data pages into full pages and free pages. When space is needed, empty pages are allocated and appended to the free pages portion of the list. When free data pages become full, they are moved from the free space portion to the end of the full pages portion of the list.



4.2 Page Directory Implementation

The Page Directory implementation differs from the Linked List implementation by using a **linked list for only header pages**. Each header page contains a pointer (byte offset) to the next header page along with entries that contain both a **pointer to a data page** and **information about the free space within that page**. Since data pages only store records, they no longer need to track pointers to neighboring pages.

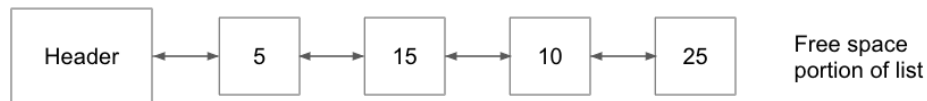


The main advantage of Page Directories over Linked Lists is that inserting records is often faster. Instead of performing I/Os to read each data page to determine whether there is enough space to insert, only header pages need to be read in order to make that decision.

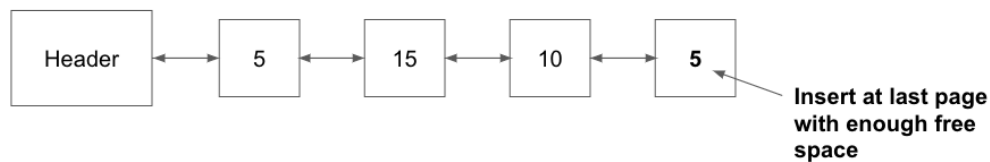
To highlight this point, consider the following example where a heap file is implemented as both a Linked List and a Page Directory. Each **page is 30 bytes** and a **20 byte record is being inserted** into the file:

Linked List

within box = # of free bytes



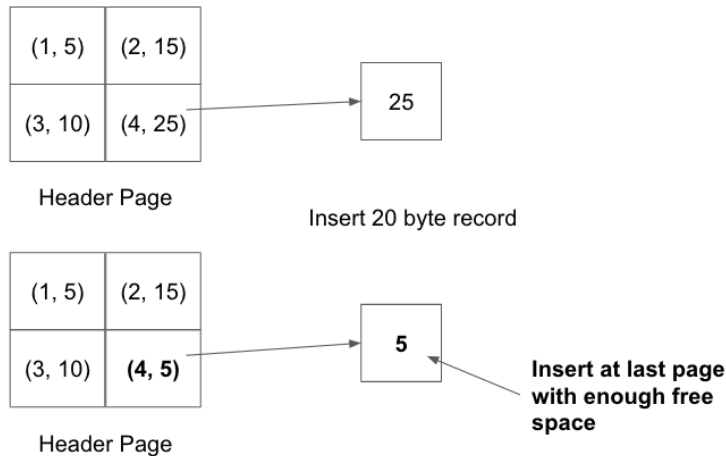
Insert 20 byte record



I/O Cost: 5 (read all pages) + 1 (write last data page) = **6 I/Os**

Page Directory

(#, #) within box = (page #, # of free bytes)



I/O Cost: 1 (read header) + 1 (read data) + 1 (write data) + 1 (write header) = **4 I/Os**

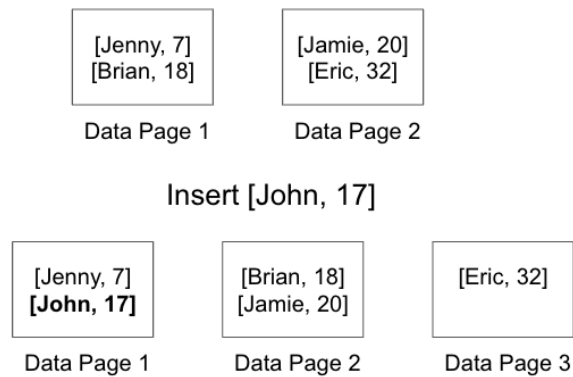
This is only a small example and as the number of pages increases, a scenario like this would cause insertion in a linked list to be many times more expensive than insertion in a page directory.

Regardless of the implementation used, heap files provide fast insertions as records can be added to any page with free space and finding a page with enough free space is often very cheap. However, searching for records within heap files requires a full scan every time. Every record on every page must be looked at because records are unordered, resulting in a linear cost of N I/Os for every search operation.

5 Sorted Files

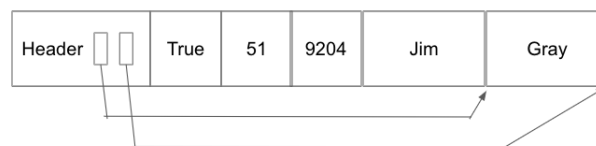
A sorted file is a file type where pages are ordered and records on pages are sorted by key(s).

These files are implemented using Page Directories and enforce an ordering upon data pages based on how records are sorted. Searching through sorted files takes $\log N$ I/Os where $N = \#$ of pages as binary search can be used to find the page containing the record. Meanwhile, insertion, in the worst case, takes $\log N + N$ I/Os as binary search is needed to find the page to write to and that inserted record could potentially cause all later records to be pushed back by one.



6 Record Types

Record types are completely determined by the relation's schema and come in 2 types: **Fixed Length Records (FLR)** and **Variable Length Records (VLR)**. FLRs only contain fixed length fields (integer, boolean, date, etc.) and FLRs with the same schema consist of the same number of bytes. Meanwhile, VLRs contain both fixed length and variable length fields (varchar), resulting in each VLR of the same schema having a potentially different number of bytes. VLRs store all fixed length fields before variable length fields and use a record header that contains pointers to the end of the variable length fields.



Regardless of the format, every record can be uniquely identified by its record id - **[page #, record # on page]**.

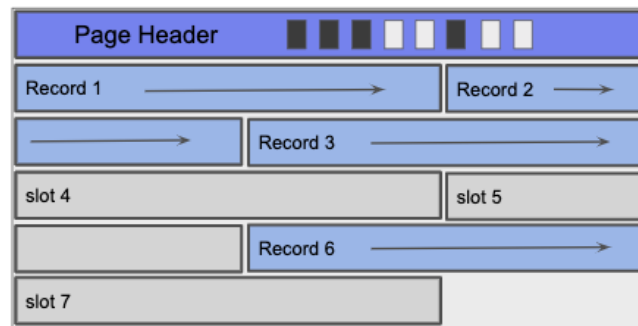
7 Page Formats

7.1 Pages with Fixed Length Records

Pages containing FLRs always use page headers to store the number of records currently on the page.

If the page is packed, there are no gaps between records. This makes insertion easy as we can calculate the next available position within the page using the # of existing records and the length of each record. Once this value is calculated, we insert the record at the computed offset. Deletion is slightly trickier as it requires moving all records after the deleted record forward by one position to keep the page packed.

If the page is unpacked, the page header typically stores an additional **bitmap** that breaks the page into slots and tracks which slots are open or taken.



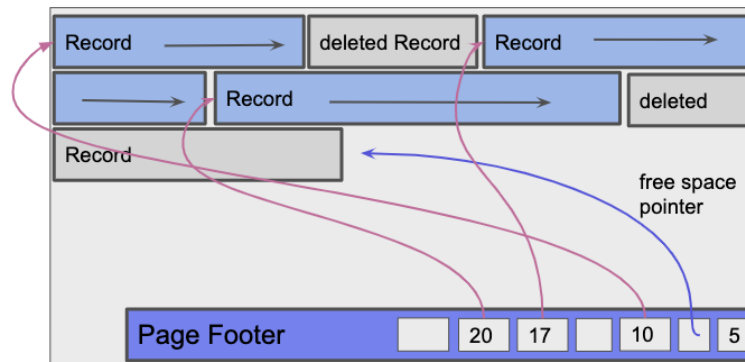
Using the bitmap, insertion involves finding the first open bit, setting the new record in the corresponding slot, and then setting the bit for that slot. With deletion, we clear the deleted record's corresponding bit so that future inserts can overwrite that slot.

7.2 Pages with Variable Length Records

The main difference between variable length records and fixed length records is that we no longer have a guarantee on the size of each record. To work around this, each page uses a **page footer** that maintains a **slot directory** tracking **slot count**, a **free space pointer**, and **entries**. The footer starts from the bottom of the page rather than the top so that the slot directory has room to grow when records are inserted.

The slot count tracks the number of used or empty slots and the free space pointer points to the next free position within the page. Each entry in the slot directory consists of a [**record pointer**, **record length**] pair.

If the page is unpacked, deletion involves finding the record's entry within the slot directory and setting both the record pointer and record length to null.



For future insertions, the record is inserted into the page at the free space pointer and a new [pointer, length] pair is set in any available null entry. In the case where there are no null entries, a new entry is added to the slot directory for that record. Periodically, records will be reorganized into a packed state where deleted records are removed to make space for future insertions.

If the page is packed, deletion involves removing the record's entry within the slot directory. Additionally, records after the deleted record must be moved forward within the page and their corresponding slot directory entries moved one position back. For insertion, the record is inserted at the free space pointer and a new entry is added every time.

8 Practice Questions

1. Given a heap file implemented as a Page Directory, what is the I/O cost to insert a record in the worst case? The directory contains 4 header pages and 3 data pages for each header page. Assume that at least one data page has enough space to fit the record.
2. What is the smallest size, in bytes, a record of the following schema can be given that the record header is 5 bytes? (boolean = 1 byte, date = 8 bytes)

name VARCHAR
student BOOLEAN
birthday DATE
state VARCHAR

3. 4 VLRs are inserted into an empty page. What is the size of the slot directory? (int = 4 bytes)

9 Solutions

1. In the worst case, the only data page with enough free space is on the very last header page. Therefore, the cost is 6 I/Os.

$$3 \text{ (read header pages)} + 1 \text{ (read data)} + 1 \text{ (write data)} + 1 \text{ (write last header)} = 6$$

2. The smallest size of the VLR is 14 bytes and occurs when both name and state are null.

$$5 \text{ (record header)} + 1 \text{ (boolean)} + 8 \text{ (date)} = 14$$

3. The slot directory contains a slot count, free space pointer, and entries, which are record pointer, record size pairs. Since pointers are just byte offsets within the page, the size of the directory is 40 bytes.

$$4 \text{ (slot count)} + 4 \text{ (free space)} + (4 \text{ (record pointer)} + 4 \text{ (record size)}) * 4 \text{ (\# records)} = 40$$