# Implementing File Systems

In this chapter we discuss various methods for storing information on secondary storage. The basic issues are device directory, free space management, and space allocation on a disk.

**11.1** In cases where the user (or system) knows exactly what data is going to be needed. Caches are algorithm-based, while a RAM disk is user-directed.

**11.2** If all extents are of the same size, and the size is predetermined, then it simplifies the block allocation scheme. A simple bit map or free list for extents would suffice. If the extents can be of any size and are allocated dynamically, then more complex allocation schemes are required. It might be difficult to find an extent of the appropriate size and there might be external fragmentation. One could use the Buddy system allocator discussed in the previous chapters to design an appropriate allocator. When the extents can be of a few fixed sizes, and these sizes are predetermined, one would have to maintain a separate bitmap or free list for each possible size. This scheme is of intermediate complexity and of intermediate flexibility in comparison to the earlier schemes.

**11.3** Such a scheme would decrease internal fragmentation. If a file is 5 KB, then it could be allocated a 4 KB block and two contiguous 512-byte blocks. In addition to maintaining a bitmap of free blocks, one would also have to maintain extra state regarding which of the subblocks are currently being used inside a block. The allocator would then have to examine this extra state to allocate subblocks and coalesce the subblocks to obtain the larger block when all of the subblocks become free.

**11.4** The advantage is that while accessing a block that is stored at the middle of a file, its location can be determined by chasing the pointers stored in the FAT as opposed to accessing all of the individual blocks of the file in a sequential manner to find the pointer to the target block. Typically, most of the FAT can be cached in memory and therefore the pointers can be determined with just memory accesses instead of having to access the disk blocks.

**11.5**   The results are:

|      | Contiguous | Linked | Indexed |
|------|------------|--------|---------|
| a.   | 201        | 1      | 1       |
| b.   | 101        | 52     | 1       |
| c.   | 1          | 3      | 1       |
| d.   | 198        | 1      | 0       |
| e.   | 98         | 52     | 0       |
| f.   | 0          | 100    | 0       |

**11.6**   $(12 * 8 /\text{KB}/) + (2048 * 8 /\text{KB}) + (2048 * 2048 * 8 /\text{KB}/) + (2048 * 2048 * 2048 * 8 /\text{KB}) = 64$ terabytes

**11.7**   One issue is maintaining consistency of the name cache. If the cache entry becomes inconsistent, then either it should be updated or its inconsistency should be detected when it is used next. If the inconsistency is detected later, then there should be a fallback mechanism for determining the new translation for the name. Also, another related issue is whether a name lookup is performed one element at a time for each subdirectory in the pathname or whether it is performed in a single shot at the server. If it is perfomed one element at a time, then the client might obtain more information regarding the translations for all of the intermediate directories. On the other hand, it increases the network traffic as a single name lookup causes a sequence of partial name lookups.

**11.8**   Restores are easier because you can go to the last backup tape, rather than the full tape. No intermediate tapes need be read. More tape is used as more files change.

**11.9**   In case of system crash (memory failure) the free-space list would not be lost as it would be if the bit map had been stored in main memory.

**11.10**   Let $Z$ be the starting file address (block number).

a.   Contiguous. Divide the logical address by 512 with $X$ and $Y$ the resulting quotient and remainder respectively.
   i. Add $X$ to $Z$ to obtain the physical block number. $Y$ is the displacement into that block.
   ii. 1

b.   Linked. Divide the logical physical address by 511 with $X$ and $Y$ the resulting quotient and remainder respectively.
   i. Chase down the linked list (getting $X + 1$ blocks). $Y + 1$ is the displacement into the last physical block.
   ii. 4

c.   Indexed. Divide the logical address by 512 with $X$ and $Y$ the resulting quotient and remainder respectively.
   i. Get the index block into memory. Physical block address is contained in the index block at location $X$. $Y$ is the displacement into the desired physical block.
   ii. 2

**11.11**    Dynamic tables allow more flexibility in system use growth — tables are never exceeded, avoiding artificial use limits. Unfortunately, kernel structures and code are more complicated, so there is more potential for bugs. The use of one resource can take away more system resources (by growing to accommodate the requests) than with static tables.

**11.12**    For a file system to be recoverable after a crash, it must be consistent or must be able to be made consistent. Therefore, we have to prove that logging metadata updates keeps the file system in a consistent or able-to-be-consistent state. For a file system to become inconsistent, the metadata must be written incompletely or in the wrong order to the file system data structures. With metadata logging, the writes are made to a sequential log. The complete transaction is written there before it is moved to the file system structures. If the system crashes during file system data updates, the updates can be completed based on the information in the log. Thus, logging ensures that file system changes are made completely (either before or after a crash). The order of the changes is guaranteed to be correct because of the sequential writes to the log. If a change was made incompletely to the log, it is discarded, with no changes made to the file system structures. Therefore, the structures are either consistent or can be trivially made consistent via metadata logging replay.

**11.13**    Relocation of files on secondary storage involves considerable overhead — data blocks have to be read into main memory and written back out to their new locations. Furthermore, relocation registers apply only to *sequential* files, and many disk files are not sequential. For this same reason, many new files will not require contiguous disk space; even sequential files can be allocated noncontiguous blocks if links between logically sequential blocks are maintained by the disk system.

**11.14**    a.  In order to reconstruct the free list, it would be necessary to perform "garbage collection." This would entail searching the entire directory structure to determine which pages are already allocated to jobs. Those remaining unallocated pages could be relinked as the free-space list.

   b.  The free-space list pointer could be stored on the disk, perhaps in several places.

   c.  TBA

**11.15**    This method requires more overhead then the standard contiguous allocation. It requires less overhead than the standard linked allocation.

**11.16**    The primary difficulty that might arise is due to delayed updates of data and metadata. Updates could be delayed in the hope that the same data might be updated in the future or that the updated data might be temporary and might be deleted in the near future. However, if the system were to crash without having committed the delayed updates, then the consistency of the file system is destroyed.