

University of Cergy-Pontoise

REPORT

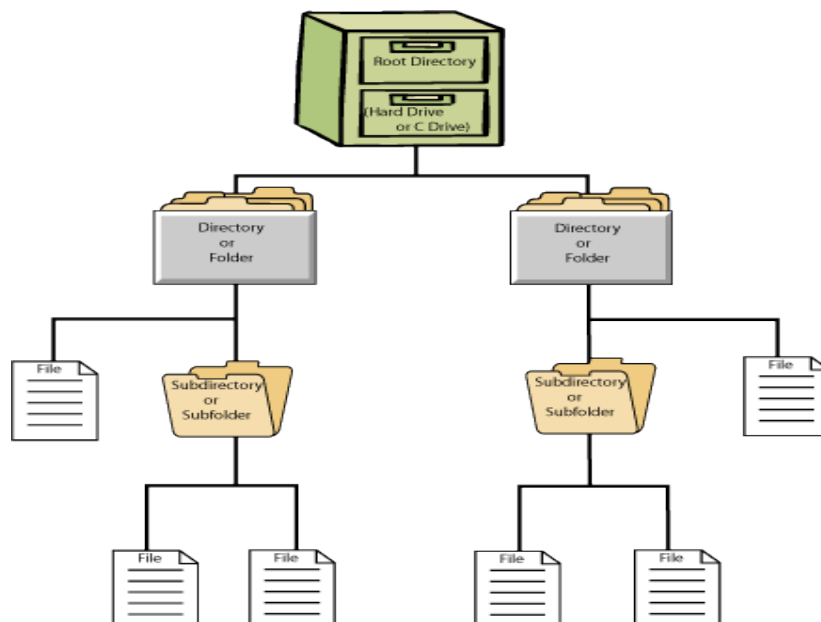
for the project of operating system
Bachelor of Computer Science third year

on the subject

file system manager

written by

AYAD Ishak and ABDELMOUMENE Djahid



March 2019

Table des matières

1	Introduction	2
2	Architecture	2
3	Implementation	5
3.1	block allocation	5
3.2	I/O lazy allocation	5
3.3	Directory entries	5

1 Introduction

In this project, we created a simplified version of the UNIX extended file system. To be mounted on an another UNIX file which was used as a virtual partition.

Our version of the file system followed a lot of the EXT standards with exception of the usage of block groups and the block group descriptors which were not used. The main file system API contains UNIX system call style functions such as open and write etc. However we did add an underscore at the end of all of these functions to avoid conflicts with the actual system.

2 Architecture

For the implementation of the file system we chose to divide different parts of the FS (file system) into five levels of abstraction :

The first one (in disk.c) separating the interaction with the UNIX file -used as a partition- from the rest of the FS. And for that a structure was created to facilitate the transportation of some useful information about the file to the rest of of the FS. This structure serves as the virtual partition and is called **fs_filesyst** and it holds the file descriptor to the partition file and the size. Which is used to simulate how a real partition has an immutable size. And lastly the number of blocks which can also be deduced from the size, these blocks have a fixed size of 4096 (FS_BLOCK_SIZE).

```
1 struct fs_filesyst {
2     uint32_t fd;          /**< file descriptor */
3     uint32_t tot_size; /**< total size of our file (partition) */
4     uint32_t nblocks; /**< number of blocks in disk image*/
5 };
```

Listing 1 – virtual filesystem structure

The second level of abstraction handles the formatting of the main parts of the FS. such as the super-block which is always in the first block and holds information about the FS in general as well as the offsets and sizes of other sections of the FS. the structure holding the superblock is as follows :

```
1 struct fs_super_block {
2     uint32_t magic;          /**< the filesystem magic number */
3
4     uint32_t data_bitmap_loc; /**< data bitmap location in block num */
5     uint32_t data_bitmap_size; /**< data bitmap size in blocks */
6     uint32_t inode_bitmap_loc; /**< inode bitmap location in block num */
7     uint32_t inode_bitmap_size; /**< inode bitmap size in blocks */
8
9     uint32_t inode_loc;      /**< location of inodes in block num */
10    uint32_t inode_count;     /**< no of inodes in blocks */
11    uint32_t data_loc;        /**< location of the data in blocks */
12    uint32_t data_count;      /**< no of data blocks */
13
14    uint32_t free_inode_count; /**< no of free inodes */
15    uint32_t free_data_count;  /**< no of free blocks */
16
17    uint32_t nreads;          /**< number of reads performed */
18    uint32_t nwrites;         /**< number of writes performed*/
19    uint32_t mounts;          /**< number of mounts*/
20
21    uint32_t mtime;           /**< time of mount of the filesystem */
22    uint32_t wtime;           /**< last write time */
23 };
```

Listing 2 – super block structure

To get the offset and sizes of the other sections, we set the ratio of the inode table as constant (1% of the total blocks), from this we can also get the bitmap size and we can leave the rest to the data blocks and its bitmap. From all of this we can get the offsets of all the sections. So at the end our file system will look like this :

File System Structure

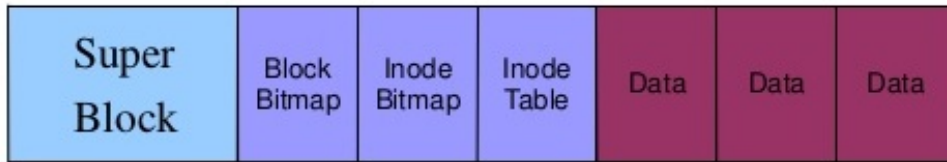


FIGURE 1 – Simplified filesystem structure

in this level we also have functions to allocate, free, read and write data and inode blocks. For the inodes this structure was used :

```
1 struct fs_inode {  
2     uint16_t mode;           /**< file type and permissions */  
3     uint16_t uid;           /**< id of owner */  
4     uint16_t gid;           /**< group id of owners */  
5     uint32_t atime;          /**< last access time in seconds since the epoch */  
6     uint32_t mtime;          /**< last modification time in seconds since the epoch*/  
7     uint32_t size;           /**< size of the file in bytes */  
8     uint32_t hcount;         /**< the number of hard links pointing to this inode*/  
9     uint32_t direct[FS_DIRECT_POINTERS_PER_INODE]; /**< direct data blocks */  
10    uint32_t indirect;        /**< indirect data blocks */  
11 };
```

Listing 3 – inode structure

And there is also the structure that holds the different representations of a data blocks :

```
1 union fs_block {  
2     struct fs_super_block super;           /**< super block */  
3     struct fs_inode inodes[FS_INODES_PER_BLOCK]; /**< array of inodes */  
4     uint32_t pointers[FS_POINTERS_PER_BLOCK]; /**< array of pointers */  
5     uint8_t data[FS_BLOCK_SIZE];          /**< array of data bytes */  
6 };
```

Listing 4 – union of a block structure

Note that a single block can hold up to 64 inodes and 1024 pointers to other blocks (used in indirect blocks).

For the third level of abstraction we no longer need to worry about the allocation of blocks and formatting, so it handles the actual input/output operations on the inodes, as well as the allocation of file descriptors. For the latter we used the following structures :

```

1 struct io_filedesc {
2     int is_allocated;
3     uint32_t offset;
4     uint32_t mode;
5     uint32_t inodenum;
6 };
7
8 struct io_filedesc_table {
9     struct io_filedesc fds[IO_MAX_FILEDESC];
10 };

```

Listing 5 – files descriptors structure

The file descriptors store the information needed to do the I/O operations such as the inode numbers as well as their current offsets -the offset from the start to write/read to/from-

For the fourth level of abstraction, we handle the directories and operations on them. And for that we used the same approach as UNIX, that is with directory entries, meaning that directories themselves are files containing 'entries' of inode numbers and filename pairs. For this we used these structures :

```

1 struct dirent {
2     uint32_t d_ino;           /* the inode number */
3     int d_type;               /* File type */
4     char d_name[256];        /* File name */
5 };

```

Listing 6 – dirent structure

```

1 typedef struct {
2     int fd;                  /* the directory file descriptor */
3     int size;                /* the number of the directory entries stored */
4     int idx;                 /* the current index in the reading (used in readdir) */
5     struct dirent* files;    /* the directory entries */
6 } DIR_;

```

Listing 7 – dir structure

the first one is a single directory entry containing a fixed size filename of 256 and the type of the file, and the inode number. As for the second structure, it holds information about the directory as a file (its file descriptor etc).

For the fifth and last level of abstractions, we used the previous two level to construct the main API's functions which are the following :

```

1 int initfs(char* filename, size_t size, int format); // to initialize the file system
2 int open_(const char* filename, int creat, uint16_t perms); // open a file
3 int close_(int fd); // close a file
4 int write_(int fd, void* data, int size); // write data to an fd
5 int read_(int fd, void* data, int size); // read data from an fd
6 int lseek_(int fd, uint32_t newoff); // change the current offset
7 int rm_(const char* filename); // remove a file
8 DIR_* opendir_(const char* dirname, int closedir(DIR* dir); // close a directory
9 int creat, uint16_t perms); // open a directory
10 struct dirent readdir_(DIR* dir); // read an entry from the directory
11 int rmdir_(const char* filename, int recursive); // remove a directory
12 int ln_(const char* src, const char* dest); // create a hard link
13 int cp_(const char* src, const char* dest); // copy a file
14 int mv_(const char* src, const char* dest); // move or rename
15 void closefs(); // close the file system

```

Listing 8 – main API's functions

3 Implementation

As for the implementation the main algorithms are those of the data blocks and inode allocations, and the reading and writing to files as well as the handling of the directory entries.

3.1 block allocation

for the allocation of data blocks and inodes, the algorithm simply parses the respective bitmaps of the wanted block type, and finds the first null bit, which meant the use of bit wise operations because C doesn't allow for the direct access of bits. after this is done the bit is set to 1 and then the necessary calculations are done to get the offsets.

3.2 I/O lazy allocation

When writing to an inode (file) we have direct and indirect blocks to write to, which are stored in the in the inode itself. When implementing this we optimised the writing so that it takes less space, that is by implementing a lazy allocation system. What it does is that it delays the allocation of a block until the very last moment it's needed (when the file is written to), meaning that if we don't write to a block it won't get allocated, even if the block is between two other allocated blocks.

An example of when this optimisation saves space is when a user writes a byte of data at the start of the file and then changes the offset to 100000 and then writes another byte of data. The idea is that only the first and last blocks are allocated. whereas the others are always not allocated.

To achieve this we have a function that figures out the blocks that need to be allocated when writing - called **io_lazy_allocation**- which takes the the start offset and the size to allocated and calculates the blocks that need to be allocated and then allocates them.

3.3 Directory entries

To store the directory entries, we used a sorted list to store the entries. the sorting is done in order of the file name, to insert a file we simply put it at the end of the directory and then we 'bubble' it into its sorted place. for the search we used a standard binary search. and for the deletion we just find the entry and then delete it, and then adjust the entries that follow.