

Operating Systems 2015 Assignment 4: File Systems

Deadline: Tuesday, May 26 before 23:59 hours.

1 Introduction

A disk can be accessed as an array of disk blocks, often each block is 512 bytes in length. In order to store files and directories in such an array of blocks, we need to think about how to organize the data. In which of the blocks will we write a given file? How can we find out in which blocks the file's content is located? How do we store other data about a file, such as its permissions and time of last modification? Finally, how do we store directories?

Several "formats" to organize such data have been devised over the years, such as FAT, NTFS, ext2, HFS and XFS. We usually refer to these formats as *file systems*.

A file system can be split in roughly two parts: the actual data and the metadata about this data. The metadata contains a table of filenames and information about these filenames such as permissions, file size and more importantly a list of disk blocks where the contents of the file can be found.

In this final assignment we will have a thorough look at a simple implementation of a file system and extend this implementation. You may have noticed that a file system named WFS has been used throughout this lab. You will extend this file system with write support for files (as the file system currently is read only) and full support for subdirectories.

2 Requirements

You have to modify existing operating system source code written in C and therefore you will be writing C code. Your submissions must adhere to the following requirements:

- Submit the source code of the operating system with your functioning extended WFS implementation.
- Your implementation should contain full support for subdirectories:
 - The user must be able to create/remove directories with the provided *mkdir* and *rmdir* utilities.
 - It must be possible to create/remove directories inside the root directory as well as inside subdirectories.
 - The implementation must be able to read and modify subdirectories on a file system image provided by us.
 - Both the structures on disk as well as the VFS structures must be updated. Changes must be persistent after unmounting the file system and rebooting the system.
 - *fsck.wfs* must always pass when run on an unmounted file system, also after the file system has been modified.
- Your implementation must be able to write to existing files with non-zero length. This includes:
 - The file size must be updated in the WFS file system on disk as well as in the VFS data structures. Changes must be persistent after unmounting the file system and rebooting the system.
 - We have provided a utility called *overwr* which overwrites a given file. See also the section on Testing below.

- After using the *overwr* utility, the *cat* utility must be able to show the correct contents of the file and *chksun* must be able to compute the correct checksum before unmounting and after unmounting and restarting (of course without modification to these utilities).
 - The *overwr* utility must work on files in both the root directory as well as subdirectories.
 - *fsck.wfs* must always pass when run on an unmounted file system, also after the file system has been modified.
- *Important!* Make sure to use extensive error handling in your code so that your code detects various kinds of failures. For example, when writing to a file system, return appropriate error codes when the file system is full (ENOSPC), a given filename is invalid (too long or contains invalid characters, EINVAL), when you are asked to create a directory in a node which is not a directory, I/O errors (EIO), or accesses to non-existing entries (ENOENT), etc.

3 Submission and Grading

You may work in teams of at most 2 persons. A single tar file should be handed in of the modified operating system containing the improved WFS implementation. To hand in your work, remove any object fields and binaries by removing the build directory. *Make sure that the files you have modified contain your names and student IDs.* Create a gzipped tar file of the source code directory:

```
tar -czvf assignment4.tar.gz assignment4/
```

Mail your tar files to *krietvel (at) liacs (dot) nl* and make sure the subject of the e-mail **equals** “OS2015 Assignment 4”. Include your names and student IDs in the e-mail.

Deadline: We expect your submissions **before** Tuesday, May 26 before 23:59. No exceptions; deliveries after the deadline *will not be graded!*

The grade is determined based on whether the program correctly implements the functionalities listed in the specification above and whether the source code looks adequate: good structure, consistent indentation, error handling, correct memory handling and comments where these are required. Comments are usually required if the code is not immediately obvious, which often means you had to make a deliberate decision or trade-off. Document these decisions, trade-offs and why in the source code. Commenting on the obvious is superfluous and bad style. Note that we may always invite teams to elaborate on their submission in an interview in case parts of the source code need further explanation.

The maximum grade that can be obtained for this assignment is 10. The points are distributed as follows: Code Layout & Quality (1.0 / 10), Reading Subdirectories (1.5 / 10), Creating/Removing Subdirectories (3.5 / 10), Writing files (4.0 / 10).

4 Kernel

We will use the same kernel as with the second and third assignments, however, you will be provided with a new starting point and a new SD card image with additional utilities, see the lab assignment website. In this assignment the first partition (FAT16) is used as root partition instead of the second partition (WFS). Because of this the utilities must now be installed on the FAT16 partition instead. Once the system is booted you can mount and mount the second partition under the */mnt* directory using the *mountfs* and *umountfs* commands (no arguments required). The commands to access the partitions (*mttools* or *hdiutil* and *wfstool*) remain the same, see the getting started guide.

5 Limitations

As the file system code is still in its infancy, you should be aware of a couple of limitations. Concurrent access to the file system (by different processes) is not possible. This will be corrected in the future by introducing proper locking in the file system code. Also we do not support setting owning users, groups and permissions.

6 Virtual File System

Like many other kernels, this kernel makes use of a Virtual File System (VFS). The Virtual File System is an in-memory representation of the system's file system. This representation has a data structure in the form of a tree.

A VFS is used for two main reasons. Firstly, performance, by storing file nodes in a VFS tree we do not have to read file entries from disk each time we want to list a directory or get the size of a file. Secondly, because of support for mount points. You can mount a file system at a given directory in another file system. For example, in our kernel a "devfs" file system is mounted under `/dev`. Thirdly for flexibility. Using VFS, we have an overview of the entire file system of the system (which comprises multiple file systems) as well as a generic interface for accessing files and directories. The VFS will take care to call the functions of the correct file system implementation in order to carry out the desired operation on the file or directory. The interface makes it easy to implement file systems which is essentially done by implementing the functions in the VFS file ops and file sys structures.

The VFS consists of a number of fundamental types: `vfs_vnode_t` that represents any file in the VFS, including normal files, mount points, directories and devices. `vfs_file_t` that represents an open file, that is, it has a VFS node and an offset where the current read / write pointer is located in the file. Further, each file system implements a number of file operations. These are stored in the type `vfs_fileops_t`. This structure contains what would be called virtual functions in some programming languages. Essentially, when calling the `vfs_open` function, the call will be dispatched to the open function, whose address is stored in the file ops structure for the given vnode. The next important type is `vfs_filesys_t` that implements the basic interface for mounting and unmounting the file system; it also contains the function that returns the root node of a file system.

7 The WFS File System

The WFS file system has been inspired by FAT, but has been greatly simplified. The file system has a fixed size, fixed amount of entries per directory, no support for user/group settings and file permissions and does not make use of a superblock.

The file system spans a fixed length of 8425488 bytes. At the start of this area are 16-bytes of magic numbers that are used for file system identification. This is followed by an area which contains the file entries of the root directory. There are 64 entries in total. Each file entry has the following format:

```
typedef struct
{
    char filename[58];
    uint16_t start_block;

    /* Given that the maximum size of a file is about 8MiB, we use the top
     * 4 bits of the size field for flags.
     */
    uint32_t size;
} __attribute__((__packed__)) wfs_file_entry_t;
```

The size of each entry is 64 bytes. Storing 64 of such entries requires 4096 bytes, or 8 512-byte disk blocks.

The 4 high bits of the `size` field are reserved for special information. Bit 32 is set when the entry describes a directory (instead of a file). The `start_block` field points out the first disk block (in the data section) that contains the file's content. Subsequent blocks can be found in the block table.

A filename is restricted to 58 bytes and may only contain: A-Z, a-z, 0-9 and dots (“.”). Make sure to verify this when entering new files in the file system.

We have limited the file system to support 16384 blocks. This means that at most 8 megabytes of information can be stored in a WFS file system. In order to know in which blocks the contents of a file are stored, we make use of the `start_block` field and the block table. The block table is indexed by `block - 1`, so `block_table[block - 1]` gives you the block that follows after `block`. There are two special block codes: `0x0` means that the block is free and can be allocated, `0xfffe` means that no block will follow the current block. So, to read a full file, you walk through the block table starting at the start block until the end of file code `0xffff` is detected. A similar scheme is used in the FAT file system.

Given that we support at most 16384 blocks, the size of the block table is 2 bytes (`sizeof(uint16_t)`) times 16384, which is 32768 bytes. See also the `WFS_BLOCK_TABLE_SIZE` define in `wfs.h`. The table starts after the table with file entries for the root directory, see also the define `WFS_BLOCK_TABLE_START`.

After the block table follows the data area. The data area is large enough to hold 16384 blocks of 512 bytes. Remember that block 0 has a special meaning and is not stored, so to compute the offset of a block, you use `WFS_DATA_START + (block - 1) × WFS_BLOCK_SIZE`.

Schematically, the file system has the following layout:

Magic numbers 16 bytes
Root directory entries 4096 bytes
Block table 32768 bytes
Data area 8388608 bytes

8 Subdirectory Support

We have seen that the entries for the root directory are stored at a special place in the file system. By setting a certain bit in the size field, an entry can be made to indicate a directory instead of a file. An entry that indicates a directory is essentially a subdirectory and what is missing here is a place to store the file entries of that subdirectory. For subdirectories, the file entries should be stored in a specially allocated disk block. We will keep it easy and limit the amount of file entries in a subdirectory to 8. Given that 8 times 64 makes 512, this fits exactly in a single disk block. The block number is stored in the `start_block` field in the subdirectory's file entry. Make sure to terminate the chain for future compatibility.

8.1 Reading Subdirectories

In order to implement support for reading existing subdirectories, you need to modify the function `wfs_directory_readdir`. You need to determine whether the given node is the root directory of the file system, or a subdirectory. This is important, because you need to know how many directory entries the directory has and where these entries are located. You can find out the directory type by looking at the node's parent, if this node is of type `VFS_MOUNT` you know this node is the root directory, otherwise it is a subdirectory. Finally, you modify the function to read the correct number of file entries from the correct location.

8.2 Creating and Removing Subdirectories

To add the ability to create and remove subdirectories, the functions `wfs_directory_mkdir` and `wfs_directory_rmdir` should be implemented and registered in the `wfs_directory_fileops` structure. Look for the exact prototypes of these two functions in `vfs.h`.

The `mkdir` function has a parent node (which must be a directory, do not forget to verify this) and a name string as argument. A subdirectory with the name string as name should be created in the file system. Remember to update the metadata on the disk itself as well as the VFS data structures. For the latter, the function `vfs_dirp_addnode` will come in handy.

Implementing the `rmdir` function is easier. A single node is given as argument, which is the directory to remove. Look for this node in the table of directory entries in the node's parent and remove the node. Again update the metadata on disk as well as the VFS data structures and for `rmdir` the function `vfs_dirp_removenode` will be of use when updating the VFS data structures.

Note that the implementations of the `mkdir` and `rmdir` system calls (in `syscalls.c`) already perform various checks (for instance whether the directory to be removed is empty). You do not have to replicate checks that are already done in the system call or VFS function, so investigate what is and what is not already checked when the control flow reaches your code.

9 File Write Support

To implement write support for files, you will have to modify the function `wfs_file_write` in `wfs.c`. Before you start writing the code, we recommend that you study the code of `wfs_file_read` first. As argument you get a pointer to a `vfs_file_t` structure, which has a `vnode` field pointing out the node's type and `fs_data` and an `off` field that indicates the current read/write offset into the file.

The purpose of the function is to write `len` bytes from the buffer pointed to by `buff` to the given `file` starting at the offset indicated by `file->off`. When implementing this, it is important to be aware of the following:

1. `len` can be any length and does not have to be aligned on disk block boundaries.
2. You have to allocate new disk blocks using the block table when this is necessary.
3. Make sure to set the block table entry of the last block to `WFS_BLOCK_EOF` to indicate the end of a chain.
4. Your implementation of `wfs_file_write` must be able to cope with non-block-aligned writes.
5. Your implementation of `wfs_file_write` must be able to cope with seeks that are interleaved with write actions. What happens if a seek is performed past the end of the file followed by a write?
6. You only have to support writes to existing files with non-zero length. So, this guarantees that at least one block is already allocated to the file. Make sure to check this and bail out with an error if this is not the case! As a consequence of this, you do not have to write special code to deal with writes to new files that previously did not exist.
7. You do *not* have to support file truncation. If you write 20 bytes to the beginning of a 200-byte file, you can leave the remaining 180 bytes unmodified and do not have to throw these away. In other words, you only add support for growing of files and not for shrinking of files.

If you have modified the start block or the file size, update the fields in the VFS node (`file->vnode`) and FS data (`file->vnode->fs_data`) and write a function `wfs_update_node(vfs_vnode_t *vnode)` which will update the given `vnode` in the file system metadata on disk. This update function should be called when the file is closed (`wfs_file_close`).

10 Testing

We have provided several utilities to test your code. We strongly recommended to start the assignment by implementing subdirectory support. To test support for reading subdirectories, use the initialized SD card image from the website. You should be able to inspect a hierarchy of subdirectories using *ls* and *cd*. A textual representation of the hierarchy as stored on the WFS partition is available from the website.

To test the support for creating and removing subdirectories, you can make use of the *mkdir* and *rmdir* utilities. These take the directory to create/remove as an argument. Note that the utilities do not accept absolute paths, only relative paths. You should test creating subdirectories in the root directory, but also creating subdirectories in subdirectories (nesting).

For testing file write support, you can use the *overwr* utility. You need to specify an *existing* file name as an argument (for now it does *not* accept absolute paths). The utility will overwrite the specified file with a particular pattern. By default this pattern is 3550 bytes in length (so files that were originally larger than 3550 bytes are only partially overwritten!).

Note that the initialized SD card image contains five files named from **file1.txt** to **file5.txt**. These files also contain a specific pattern. When you specify one of these files to *overwr* (*overwr* compares the filename), a special action will be performed that is defined for that specific file. Ensure to test your code by using *overwr* on each of the five test files. After overwriting the file, inspect the new file size with *ls* and its contents with *cat* and *chksun*. The following table specifies the correct checksums (Adler-32) and size for the overwritten files:

file1.txt	a36485b5	1278 bytes
file2.txt	fca9dcaf	43310 bytes
file3.txt	a5e94a20	16330 bytes
file4.txt	7d197fc2	16330 bytes
file5.txt	907842ae	15265 bytes

Important: Also unmount the file system (and optionally reboot) and mount it again to verify the made changes are persistent.

To verify the file system is not corrupted after making changes to it, you can use the *fsck.WFS* utility. This is a file system checker that is provided with the initialized SD card image for this assignment. This utility ensures that the data on disk is correct and not corrupted. It should report no errors and that the file system is “clean”. Though, remember that the file system check might not catch every possible error!