

VE482, LAB 8 REPORT, GROUP 9

1. What is a filesystem?

A filesystem is the methods and data structures that an operating system uses to keep track of files on a disk or partition; that is, the way the files are organized on the disk. The word is also used to refer to a partition or disk that is used to store the files or the type of the filesystem¹.

2. How is the Linux VFS working?

In Linux, all files are accessed through VFS. This is a layer of code which implements generic filesystem actions and vectors requests to the correct specific code to handle the request². Two main types of code modules take advantage of the VFS services, which are device drivers and filesystems.

The role of the VFS is:

- Keep track of available filesystem types
- Associate (and dissociate) devices with instances of the appropriate filesystem.
- Do any reasonable generic processing for operations involving files.
- When filesystem-specific operations become necessary, vector them to the filesystem in charge of the file, directory, or inode in question.

The interaction between the VFS and specific filesystem types occurs through two main data structures, the `super_block` structure and the `inode` structure, and their associated data structures, including `super_operations`, `inode_operations`, `file_operations`, and others.

The source code for the VFS is in the `fs/` subdirectory of the Linux kernel source, along with a few other related pieces, such as the buffer cache and code to deal with each executable file format. Each specific filesystem is kept in a lower subdirectory.

It uses `init_name_fs()` to register filesystem with the registry kept in `fs/super.c`. It's also possible for a filesystem to support more than one type of filesystem. For instance, in `fs/sysv/inode.c`, three possible filesystem types are supported by one filesystem with the following codes:

```
1 static struct file_system_type sysv_fs_type[3] = {
2     {sysv_read_super, "xenix", 1, NULL},
3     {sysv_read_super, "sysv", 1, NULL},
4     {sysv_read_super, "coherent", 1, NULL}
5 };
6
7 int init_sysv_fs(void)
8 {
9     int i;
10    int ouch;
11
12    for (i = 0; i < 3; i++) {
13        if ((ouch = register_filesystem(&sysv_fs_type[i])) != 0)
14            return ouch;
15    }
16    return ouch;
17 }
```

To connect the filesystem to a disk, it will call a function(for example, `ext2_read_super()` for `ext2` filesystem) to see if a device can be mounted to the system, then if it succeeds in reading the superblock and is able to mount the filesystem, it fills in the `super_block` structure with information that includes a pointer to a structure called `super_operations`, which contains pointers to functions which do common operations related to superblocks. When a filesystem is mounted, `do_umount()` calls `read_super`, which

¹<https://www.tldp.org/LDP/sag/html/filesystems.html>

²<https://www.tldp.org/LDP/khg/HyperNews/get/fs/vfstour.html>

ends up calling the reading function, which returns the superblock, which points to that structure of pointers to functions in `_sops`.

The `inode_operations` structure is:

```
1 struct inode_operations {
2     struct file_operations * default_file_ops;
3     int (*create) (struct inode *, const char *, int, int, struct inode **);
4     int (*lookup) (struct inode *, const char *, int, struct inode **);
5     int (*link) (struct inode *, struct inode *, const char *, int);
6     int (*unlink) (struct inode *, const char *, int);
7     int (*symlink) (struct inode *, const char *, int, const char *);
8     int (*mkdir) (struct inode *, const char *, int, int);
9     int (*rmdir) (struct inode *, const char *, int);
10    int (*mknod) (struct inode *, const char *, int, int, int);
11    int (*rename) (struct inode *, const char *, int, struct inode *, const char *, int);
12    int (*readlink) (struct inode *, char *, int);
13    int (*follow_link) (struct inode *, struct inode *, int, int, struct inode **);
14    int (*readpage) (struct inode *, struct page *);
15    int (*writepage) (struct inode *, struct page *);
16    int (*bmap) (struct inode *, int);
17    void (*truncate) (struct inode *);
18    int (*permission) (struct inode *, int);
19    int (*smmap) (struct inode *, int);
20};
```

Most of these functions map directly to system calls. There are also many system calls related to files and directories which are in the `file_operations` structure. This is the same one used when writing device drivers and contains operations that work specifically on files, as the following codes show:

```
1 struct file_operations {
2     int (*lseek) (struct inode *, struct file *, off_t, int);
3     int (*read) (struct inode *, struct file *, char *, int);
4     int (*write) (struct inode *, struct file *, const char *, int);
5     int (*readdir) (struct inode *, struct file *, void *, filldir_t);
6     int (*select) (struct inode *, struct file *, int, select_table *);
7     int (*ioctl) (struct inode *, struct file *, unsigned int, unsigned long);
8     int (*mmap) (struct inode *, struct file *, struct vm_area_struct *);
9     int (*open) (struct inode *, struct file *);
10    void (*release) (struct inode *, struct file *);
11    int (*fsync) (struct inode *, struct file *);
12    int (*fasync) (struct inode *, struct file *, int);
13    int (*check_media_change) (kdev_t dev);
14    int (*revalidate) (kdev_t dev);
15};
```

3. What is FUSE, and how does it interact with the VFS? Can you sketch it quickly to make it clearer?

Fuse is the abbreviation for Filesystem in Userspace, which is a software interface for Unix-like computer operating systems that lets non-privileged users create their own files systems without editing kernel code³. It provides a "bridge" between the file system running in user space with the kernel interfaces.

The FUSE project consists of two components: the fuse kernel module(maintained in the regular kernel repositories) and the libfuse userspace library. libfuse provides the reference implementation for communicating with the FUSE kernel module⁴.

A FUSE file system is typically implemented as a standalone application that links with libfuse. libfuse provides functions to mount the file system, unmount it, read requests from the kernel, and send responses back. libfuse offers two APIs: a "high-level", synchronous API, and a "low-level" asynchronous API. In both cases, incoming requests from the kernel are passed to the main program using callbacks. When using the high-level API, the callbacks may work with file names and paths instead of inodes, and processing of a request finishes when the callback function returns. When using

³https://en.wikipedia.org/wiki/Filesystem_in_Userspace

⁴<https://github.com/libfuse/libfuse>

the low-level API, the callbacks must work with inodes and responses must be sent explicitly using a separate set of API functions. Figure 1 shows how it works⁵.

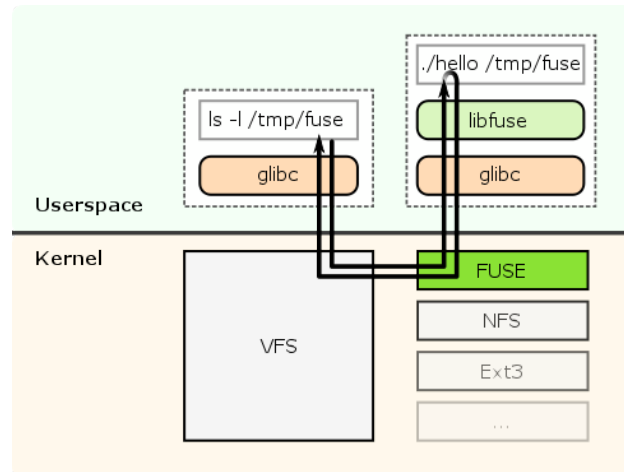


Figure 1: A flow-chart diagram illustrating how FUSE works

This is the illustration corresponds to the “hello world” file system on the FUSE website. At a high level the “hello world” file system is compiled to create a binary called “hello”. This binary is executed in the upper right hand corner of the illustration with a file system mount point of /tmp/fuse. Then the user executes an ls -l command against the mount point (ls -l /tmp/fuse). This commands goes through glibc to the VFS. The VFS then goes to the FUSE module since the mount point corresponds to a FUSE based file system. The FUSE kernel module then goes through glibc and libfuse (libfuse is the FUSE library in user space) and contacts the actual file system binary (“hello”). The file system binary returns the results back down the stack to the FUSE kernel model, back through the VFS, and finally back to the ls -l command. There are some details I have glossed over of course, but the figure illustrates the general flow of operations and data.

⁵<http://www.linux-mag.com/id/7814/>