# YaNFS: Yet Another Network File System

Zhao Ruizhe
School of EECS, Peking University

**Abstract**

As the cloud services become more and more famous recently, and personal user level data clouds become cheaper, developing practical techniques which could build a "virtual environment" for those applications is quite necessary. With the help of Secure File Transfer Protocol(SFTP), and the easy to understand and augment Linux Virtual File System(VFS) API, it's possible to conquer that task. In this report, a brief introduction to the file system, which has been named as YaNFS will be listed first, and the implementation details and the *communication protocal* between the *kernel mode* and the *user mode* will be illustrated next. Although this project is finished, there's still a chance to improve the performance. All the optimization techniques will be addressed in the end.

# Chapter 1

# Introduction

## 1.1 File System Architecture

### 1.1.1 Module and File System

Obviously, implementing a file system on Linux is just writing a *module*. Why is that? The idea is, you'll never know there's such a file system unless you read some files in it, write data to it, or listing the directories under the mounting point of that file system. All these functionalities are provided by the standard *virtual file system(VFS)* APIs, and the way we define a unique file system is by referring the fuction pointers to our self-defined functions. For example, if you want to count the number of executed read operations, you could define a read function:

```
// definition
static ssize_t read(struct file* flip, char __user *buf,
                        size_t len, loff_t *ppos)
{
        ...
        printk(KERN_INFO "Trying␣to␣read\n");
        ...
}
```

And refer to it:

```
// refering
struct file_operations file_operations = {
        .read = read,
        ...
};

file->f_op = file_operations;
```

How could Linux VFS know there's such a read function? *By defining it in the module.* So that's why I say: "implementing a file system is just writing a module."

### 1.1.2 Mounting

After we've defined the function pointers' reference, we need to put our file system under one directory and register the file system. *Mounting* is the procedure of making the files or directories under the *mounting point* directory use the VFS APIs provided by our file system, and registering is just a simple procedure to let the operating system know that there's such a file system.

Where to do the mounting and registering procedure? In the `init` call of our module:

```
1  static int __init init_fs(void )
2  {
3      int err;
4
5  // Here the register_filesystem() function will call
6  // our mounting function().
7      err = register_filesystem(&fs_type);
8
9      return err;
10 }
```

It's worthy to denote that there're 2 types of mounting: `mount_bdev` and `mount_nodev`. If you use the previous one, then all the `read` and `write` functions could operate on a block device, while the later one couldn't.

## 1.2 Kernel-User Communication

Linux VFS functions should be run under the *kernel mode*, which has the highest previledge, while at the same time, has less utilities. Several useful networking libraries are all under *user mode*. Due to this need, we have to build a communication mechanism between the *kernel mode* and the *user mode*.

Socket could do the job, but we want it to be simpler. Then we have *Netlink* module. *Netlink* is a special family of socket protocols, which is designed to provide communication between *kernel mode* and *user mode*.

## 1.3 SSH and SFTP

Secure Shell(SSH) is one of the most useful tool for programmers: It could create a remote shell for you to operate on a remote machine, at the same time, the communication is quite secure. Secure File Transfer Protocol(SFTP) is a secure version of FTP, which you could do file operations through the SSH session. SSH session will be created once you've logged on your remote machine.

LibSSH could help you do a lot of tasks based on SSH and SFTP, and our application is built on that library. It's mechanism is quite simple: After you've created a SSH session by doing the authetication and specifying the host and port, you could call any kind of remote operations afterwards.

# Chapter 2

# Implementation

## 2.1 Design Philosophy

The goal of our YaNFS has 3 points, which have been listed previously:

**Virtual Environment** The host machine couldn't know any other files on the remote machine, which means you need to build your own directory tree on the host("directory tree" is not a rigorous definition, it's just referring the inodes and dentries).

**Kernel-User Mode** The file system is built in kernel mode while the remote connection is built in user mode.

**SSH Connection** The remote connection is established through SSH and SFTP protocol.

And the whole project will be built on 3 different layers: *file system*, *communication* and *SSH connection*. It's rather good to understand this through an example.

Here I'd like to read a file on the remote machine:

1. Call `read()` in kernel mode.

2. Translate from dentry hierachy links to file pathname string.

3. Build message string "read [filepath]", and stall the kernel mode `read()`.

4. SSH daemon come to ask if there's any command to execute, and the message string found.

5. According to the message string, call `sftp_read(filepath, buf, size)` to store the read data to `buf`.

6. Send back read data to the kernel.

7. Kernel mode `read()` found the message, and copy to the `buf` parameter provided by the function call by `copy_to_user()`.

## 2.2 File System Implementation

### 2.2.1 File System Type

It's required to provide a `file_system_type` parameter if you want to register your file system.

```
1  static struct file_system_type sshfs_fs_type = {
2      .name     = "kernelsshfs",
3      .mount    = sshfs_mount,
4      .kill_sb  = sshfs_kill_sb,
5      .fs_flags = FS_USERNS_MOUNT, // what for?
6  };
7
8  struct dentry *sshfs_mount(struct file_system_type *fs_type,
9                             int flags,
10                             const char *dev_name,
11                             void *data)
12  {
13      return mount_nodev(fs_type, flags, data, sshfs_fill_super);
14  }
15
16  static void sshfs_kill_sb(struct super_block *sb)
17  {
18      kill_litter_super(sb);
19  }
```

### 2.2.2 Super Block

Look at the `sshfs_mount` function, which will do `mount_nodev`. `mount_nodev` denotes our file system does not need a backend block device. And `mount_nodev` need a function `sshfs_fill_super` to fulfill the properties of the super block of our file system.

It's the definition of the *super block operations*. The functions are provided by default Linux VFS super block operations.

```
1  static const struct super_operations sshfs_ops = {
2      .statfs       = simple_statfs,
3      .drop_inode   = generic_delete_inode,
4      .show_options = generic_show_options,
5  };
```

Here the `sshfs_get_inode` will not be covered in this document, you could easily find it in the `sshfs_inode.c` file. So the basic operations `sshfs_fill_super` does are just setting:

1. `s_op`

2. `s_root`

3. `s_magic`

```
1  int sshfs_fill_super(struct super_block *sb, void *data,
2          int silent)
3  {
4      sb->s_blocksize      = PAGE_CACHE_SIZE;
5      sb->s_blocksize_bits = PAGE_CACHE_SHIFT;
6      sb->s_magic          = SSHFS_MAGIC;
7      sb->s_op             = &sshfs_ops;
8
9      inode = sshfs_get_inode(sb, NULL,
10                 S_IFDIR|SSHFS_DEFAULT_MODE, 0);
11
12     sb->s_root = d_make_root(inode);
13         ...
14 }
```

### 2.2.3   File and Inode Operations

In order to satisfy our "virtual environment" requirement, we need to:

1. build `inode` both on host and remote.

2. only do file operations on remote.

3. querying the host-only inodes.

And the code definition of inode operations, and only the first 4 operations are self-defined. All the other operations are default operations which means they will not visit the remote information.

**inode operations**
```
1  static const struct inode_operations
2        sshfs_dir_inode_operations = {
3      .mknod  = sshfs_mknod,
4      .mkdir  = sshfs_mkdir,
5      .rmdir  = sshfs_rmdir,
6      .create = sshfs_create,
7      .lookup = simple_lookup,
8      .link   = simple_link,
9      .unlink = simple_unlink,
10     .rename = simple_rename,
11 };
```

File operations are much more simplified, only `read`, `write` and `llseek` are defined.

**file operations**
```
1  const struct file_operations
2        sshfs_file_operations = {
3      .read   = sshfs_file_read,
4      .write  = sshfs_file_write,
5      .llseek = sshfs_llseek,
6      .fsync  = noop_fsync,
7  };
```

All the self-defined operations are based on one function, `netlink_send_command`, to send the command message like "read dir1/file1" to the netlink module. I'll cover this in the next section.

## 2.3   Netlink Communication

Only 2 APIs are important in this section.

### 2.3.1   Send Command and Stall

```
static int netlink_send_command(cmdtype cmd,
                                char *params,
                                struct iovec *iov)
{
        ...
    cmd_msg_buffer.cmd    =
                kzalloc(sizeof(char) * cmdstr_len, GFP_KERNEL);
    cmd_msg_buffer.cmdlen = cmdstr_len;
    cmd_msg_buffer.exec   = false;
    cmd_msg_buffer.msglen = 0;
    cmd_msg_buffer.msg    = NULL;

    strcpy(cmd_msg_buffer.cmd, cmdstr);

    cmd_msg_buffer.ready  = true;
    while(!cmd_msg_buffer.exec)
        msleep(1000);
        ...
    iov->iov_len = cmd_msg_buffer.msglen+1;
    iov->iov_base =
        (char *) kzalloc((iov->iov_len), GFP_KERNEL);
    strcpy(iov->iov_base, cmd_msg_buffer.msg);


        ...
    kzfree(cmd_msg_buffer.msg);
    cmd_msg_buffer.msg = NULL;

    return 0;
}
```

In this function, we'll first fulfill the `cmd` buffer in the global parameter: `cmd_msg_buffer`, and set the `ready` bool to `true`. Then this procedure will stall until the `exec` is `true`.

After the stall is finished, this function will copy back the message from the `msg` buffer. How to use that message is specified by different methods. Take the `sshfs_file_read` as an instance:

**file read**

```
1   static ssize_t sshfs_file_read(struct file* flip,
2           char __user *buf, size_t len, loff_t *ppos)
3   {
4           ...
5       netlink_send_command(Cmd_read, params, iov);
6
7       if(copy_to_user(buf, (void*)iov->iov_base, iov->iov_len))
8           return -EFAULT;
9
10      *ppos += iov->iov_len;
11      ret = (ssize_t) iov->iov_len - 1;
12
13      kzfree(params);
14      kzfree(iov);
15      return ret;
16  }
```

This one just copy the received message to the buffer required by the `read()` VFS API.

## 2.3.2   Message Receiving

It's necessary to issue that there're 3 steps in this module:

1. First of all, this method will be called only if there's a message sent from SSH daemon. The message is contained in **struct** `sk_buff *skb`.

2. This method will parse the message from SSH daemon next, there're 3 kind of status:

   | Status | Meaning |
   |--------|---------|
   | READY | SSH daemon is working and has nothing to do. |
   | message | SSH daemon is transferring message and |
   | | the content of received message have not yet finished. |
   | DONE | SSH daemon has finished transferring. |

3. At last, send the response to SSH Daemon by `nlmsg_unicase`.

   | Status | Meaning |
   |--------|---------|
   | NOTHING | If SSH daemon is ready and the file system has no workload |
   | COPIED | If SSH daemon is transferring message and the data is copied. |
   | OK | SSH daemon has finished transferring and response to it. |

The codes are listed below:

**Message Receiving**

```
1  static void netlink_recv_msg(struct sk_buff *skb)
2  {
3      struct nlmsghdr *nlh;
4      int pid;
5          ...
6
7      nlh = (struct nlmsghdr*) skb->data;
8      pid = nlh->nlmsg_pid;
9      recvmsg = (char *) nlmsg_data(nlh);
10
11     if (!strcmp(recvmsg, recv_ready_str))
12         ... // When SSH Daemon is Ready
13     else if (!strcmp(recvmsg, recv_done_str))
14         ... // When SSH Daemon is Done
15     else
16         ... // When SSH Daemon is transferring data.
17     skb_out = nlmsg_new(msg_size, 0);
18     nlh = nlmsg_put(skb_out, 0, 0, NLMSG_DONE, msg_size, 0);
19     NETLINK_CB(skb_out).dst_group = 0;
20     strncpy(nlmsg_data(nlh), msg, msg_size);
21
22     res = nlmsg_unicast(nl_sk, skb_out, pid);
23 }
```

## 2.4   SSH Daemon

The SSH daemon is just like a bridge, which lies between the kernel mode `netlink` module, and the user mode `libssh` library.

### 2.4.1   Communicate with Netlink

Kernel mode netlink module could only response once SSH daemon send message to it. And SSH send message function follows these steps:

1. Initialize **struct** `nlmsghdr` by assigning `nlmsg_len`, `nlmsg_pid` and `nlmsg_flags`.

2. Initialize message itself. Send the message and stall.

3. Once received the message from the kernel, do the parsing work to get what to do next.

**Communicate with Netlink**

```
1  char * send_recv_netlink_portal(char *msg_str, int msg_len)
2  {
3      nlh = (struct nlmsghdr *) malloc(NLMSG_SPACE(MAX_PAYLOAD));
4      memset(nlh, 0, NLMSG_SPACE(MAX_PAYLOAD));
5      nlh->nlmsg_len   = NLMSG_SPACE(MAX_PAYLOAD);
6      nlh->nlmsg_pid   = getpid();
7      nlh->nlmsg_flags = 0;
8
9      strncpy(NLMSG_DATA(nlh), msg_str, msg_len);
10
11     iov.iov_base    = (void *)nlh;
12     iov.iov_len     = nlh->nlmsg_len;
13
14     msg.msg_name    = (void *)&dest_addr;
15     msg.msg_namelen = sizeof(dest_addr);
16     msg.msg_iov     = &iov;
17     msg.msg_iovlen  = 1;
18
19     sendmsg(sock_fd, &msg, 0);
20         \\ will stall here
21     recvmsg(sock_fd, &msg, 0);
22
23     return (char *)NLMSG_DATA(nlh);
24 }
```

There's a infinite loop in the `main()` function of SSH daemon:

**Main loop**

```
1  while (1)
2  {
3      char *recv =
4          send_recv_netlink_portal("READY", strlen("READY"));
5      if (!strcmp(recv, "ERROR"))
6          ...
7      else if (!strcmp(recv, "NOTHING"))
8          sleep(1);
9      else
10     {
11         struct iovec *iov =
12             sftp_execute(ssh_s, sftp_s, recv);
13         ...
14         recv =
15             send_recv_netlink_portal("DONE", strlen("DONE"));
16         if (strcmp(recv, "OK"))
17             exit(1);
18         sleep(1);
19     }
20 }
```

It's clear that this daemon will first send READY, and do different tasks by specifying the different responses. If SSH daemon is sure that it'll do some work on the remote machine, it'll call `sftp_execute`.

## 2.4.2 Communicate with SSH Library

There's no need to show some details about the libSSH library, so here I'll just take one function for instance: `sftp_execute_mkdir`. This function will be called only if the command sent by netlink module is something like: `mkdir dir1`.

**Remote Mkdir**

```
1  int sftp_execute_mkdir(ssh_session ssh,
2                         sftp_session sftp,
3                         char *dirname)
4  {
5      int rc;
6
7      rc = sftp_mkdir(sftp, dirname, S_IRWXU);
8      if (rc != SSH_OK)
9      {
10         if (sftp_get_error(sftp) != SSH_FX_FILE_ALREADY_EXISTS)
11         {
12             fprintf(stderr, "Can't␣create␣directory:␣%s\n",
13                     ssh_get_error(ssh));
14             return rc;
15         }
16     }
17
18     return rc;
19 }
```

As you can see, the key function is `sftp_mkdir`.
Another example is `sftp_execute_read`:

**Remote Read**

```
1  int sftp_execute_read(ssh_session ssh,
2                        sftp_session sftp,
3                        char *params, char **msg)
4  {   ....
5      file = sftp_open(sftp, filename, access_type, 0);
6      sftp_seek(file, offset);
7      for (;;)
8      {
9          nbytes = sftp_read(file, buffer, sizeof(buffer));
10         ....
11         if (nbytes == 0) break;
12         else if (nbytes < 0) return SSH_ERROR;
13
14         int buflen = strlen(buffer);
15         int msglen = (*msg != NULL) ? strlen(*msg) : 0;
16         msglen += buflen;
17         msglen ++;
18         ....
19     }
20     rc = sftp_close(file);
21     return rc;
22 }
```

The data received from remote operation `sftp_read` will be inserted into local buffer, and send back to netlink once the transfer is finished.

After the netlink received the message, it will copy the content to the `read()` API's buffer.