

CS380L: Advanced Operating Systems Lab #3

Zeyuan Hu ¹, iamzeyuanhu@utexas.edu

EID:zh4378 Spring 2019

1 Environment

Unless otherwise noted, we use a Linux server for all the experiments. The server has 4 Intel(R) Xeon(R) CPU E3-1220 v5 @ 3.00GHz processors and 16GB of memory, and runs Ubuntu 16.04.2 LTS (kernel version 4.11.0).

2 ELF format

3 execve implementation

To write our own loader program, we first need to understand how `execve` system call is implemented in Linux kernel ². Essentially, `execve` performs the heavylifting work of loading a program into memory and start the execution of program as a process in Linux. The main logic of `execve` is implemented in `do_execveat_common` in `fs/exec.c`. Inside `do_execveat_common`, we see the critical function call stack looks like below:

```
1 do_execveat_common
2   |- exec_binprm
3     |- search_binary_handler
4       |- load_binary
```

For binary with ELF format, `load_elf_binary` is registered with `load_binary` in `fs/binfmt_elf.c`. Thus, we locate the core function `load_elf_binary` to load binary with ELF format in the Linux kernel. Further study of the function reveals that the following critical steps are performed in order to load the program into memory and start execution:

- Read the ELF header and perform some simple consistency checks. For example, the function checks whether the binary file is executable file (`ET_EXEC`) or shared object file (`ET_DYN`).
- Read the program header table by calling `load_elf_phdrs`.

¹30 hours spent on this lab.

²we study the source code of Linux 4.8.12

- Walk through the entries of the program header table and perform specific actions based on program header table entry type (e.g., `p_type`). Note that `e_phnum` holds the number of entries in the program header table. Each entry of the table corresponds to a memory segment in the binary file. Kernel is only interested in three types of program header entries during the walk through of the program header table:
 - `PT_INTERP` entry, which identifies the run-time linker needed to assemble the complete program [1].
 - `PT_GNU_STACK` entry, which determines whether the program's stack should be executable [2].
 - `PT_LOPROC` ... `PT_HIPROC`, which are values reserved for processor-specific semantics [1].
- Once all the program header table entries have been processed, the function performs some checks based on `p_type` value it obtained from the previous step (e.g., check for interpreter, check for specific processor architecture).
- The function now is ready to set up the new program. The very first step is to call `flush_old_exec`, which clears up state in the kernel that refers to the previous program. Some minor setups (e.g., set up program's personality [3]) are performed immediately after.
- `setup_new_exec` is called to set up the kernel's internal state for the new program and new credentials for this executable is installed via calling `install_exec_creds`.
- `setup_arg_pages` is invoked to set up kernel's memory tracking structures (e.g., stack `vm_area_struct`). This step is part of the goal to set up virtual memory for the new program.
- The function now traverse the program header table again and look for entries with type `PT_LOAD`, which indicates loadable segments (i.e., areas of the new program's running memory). The entries contain code and data sections that come from the executable file and the size of a BSS section. For each `PT_LOAD` entry, the function maps it into the process' address space via `elf_map` call and sets up the new program's memory layout accordingly.
- `set_brk` is invoked to set up zero-filled pages that correspond to the program's BSS [4] segment.
- `create_elf_tables` is called to set up the rest of the new program's stack. Basically, the function puts `argc`, `argv`, `envp`, and auxiliary vectors. LWN article [2] provides an illustration of the content of the stack set up by the kernel. In addition, Figure 3.11 of System V ABI for `x86_64` manual [5] provides what initial process stack looks like.
- `start_thread` is invoked to start the execution of the new program.

4 User-space Loader

We leverage `libfuse` and `libssh` to implement a network file system. Filesystem in Userspace (FUSE) is a user-space file system framework. FUSE consists of a Linux kernel module and a user-level daemon. When a user application performs operations on a mounted FUSE file system, the operation will be routed to FUSE's kernel driver by VFS. The operations as requests will be maintained by a queue and user-level daemon will pick a request from the kernel queue and process the request. Daemon will write response back to kernel once it is done with processing. More information about FUSE can be seen in [6].

4.0.1 System Architecture

We implement a network file system. Like NFS, it supports a server with multiple clients. On the client side, the FUSE client program will mount a user-space file system with remote user, remote host, remote path, local cache path, and local mount point provided by the user. Remote user and remote host specify the server that our network file system want to contact. Remote path specifies the location on the server we want to store and fetch files for the clients. Our system requires user to specify a path for local cache as we serve clients' requests from local cache as much as we can. In other words, some file operations are directly served by the local cache without further contact with server. Communications between clients and the server are done via SSH protocol. Thus, there is no server-side implementation in our system.

4.0.2 FUSE Calls Implemented

We implement the following operations: `getattr`, `readdir`, `create`, `open`, `read`, `write`, `fsync`, `release`, `mkdir`, `unlink`, and `rmdir` in our system. Our operation implementation is based on `libfuse` version is 2.9.4 ³ API. In `libfuse`, there are two sets of APIs: a "high-level", synchronous API, and a "low-level" asynchronous API. The key difference between "high-level" and "low-level" API is that "high-level" allows us to work with file names and path instead of inodes in synchronous fashion. Thus, for the simplicity, our implementation adopts the "high-level" API.

`getattr` operation is used to get attributes of a file or a directory. Function signature of `libfuse's getattr` operation is `int getattr(const char * path, struct stat * stbuf);`. Our implementation should fill `stbuf` to contain attributes of file or directory indicated by `path`. Since `struct stat` is the same structure used in `stat` call [7], we can leverage `stat` command [8]. Specifically, in our implementation, we execute `stat` command remotely via SSH to collect attributes of the desired target (specified by `path`). The output of the `stat` command is parsed on the client side of the file system. We fill `st_size`, `st_blocks`, `st_mode`, and `st_nlink` fields of `struct stat` from

³check via `fusermount -V`

the parse result. `st_uid` and `st_gid` are filled with UID and GID of the file system client process. `st_mode` is filled with 644 (i.e., `rw-r--r--`) if the target is file and 755 (i.e., `rw-r-xr-x`) if the target is directory.

`readdir` operation is used to get entries in a directory (i.e., read directory). There are two modes of operation for the implementation: whether we keep track of offset provided as one of function signature argument. If we ignore the offset, the content of whole directory will be read in a single operation. For the simplicity, we ignore the offset in our implementation. This can be troublesome if the content of directory is greater than the supplied buffer size. The key to our implementation is the implementation of filler (with type `fuse_fill_dir_t`), which is used to add directory entry to the supplied buffer. In our case, `getattr` is invoked whenever filler is called.

`create` and `open` operations correspond to `creat` and `open` system calls. `libfuse`'s `open` operation has signature `int open(const char* path, struct fuse_file_info* fi);`. `path` specifies the file to be opened and `fi->flags` indicates the open flags (same as `flags` in `open` system call [9]). `fi->fh` represents file handle, which may be filled for future usage. In our implementation, when `open` is invoked, the target file will be downloaded from server via SCP and saved in the local cache path. Operations like `read`, `write`, and `fsync` will be served from local cache copy. The local cache copy is opened with the same flags and corresponding `fd` is saved in `fi->fh` for future use. `create` operation is similar to `open` except it will first create the file remotely if the file does not exist.

`read`, `write`, and `fsync` operations are served by local cache file copy via `fi->fh`. Same as NFSv2, we implement flush-on-close semantic for update visibility. Specifically, we do not upload the file to the server on `write` and `fsync`, instead we update server's copy when the file is closed. Doing so removes network communication overhead between each file update, but we may face file inconsistency if there are multiple clients updating the same file competitively.

`release` operation is invoked by `libfuse` when closing a file descriptor. This is the place where we implement flush-on-close semantic. We first close `fi->fh`, which is the file descriptor the local cache copy. Then, we upload the local cache copy to the server via SCP.

`mkdir`, `rmdir`, and `unlink` operations are performed by executing `mkdir` and `rm` command remotely via SSH.

4.0.3 Limitation

Our implementation has coarse-grained access control in the sense that all the files and directories have the same UID, GID, and permission mask. To enable a more fine-grained access control, we could use a file similar to `/etc/exports` in NFS to indicate what ip address can have what access (read, write, or both) to what directories and files. Doing so requires us to implement a server-side code as a guard to perform identity check. However, access control is orthogonal to our experiment goal and we left this feature as future work.

Incomplete file metadata also impacts how we implement the file consistency mechanism. In our implementation, we do not keep creation time, access time, and modification time of directories and files. As a result, we cannot selectively perform SCP on **open**: if we maintain file stats, we can perform stats comparison between remote copy and cached copy to see whether we need to perform expensive data transmission over the network. In addition, since we only SCP file during **open**, file can change between **open** and **read**. NFS periodically issues **getattr** to server to ensure cache consistency. Since we do not support file timestamp, we cannot issue **getattr** periodically and let **read** directly read from server when the local cache is invalidated.

Some other issue might exist regarding file consistency. For example, we implement the last-writer-wins policy: if a file is updated by the multiple users, the last one who close the file will keep its change to the file on server. However, this might be troublesome. A more sophisticated method is to automatically merge change to the file whenever possible and maintain multiple versions of files on the server with each version associated with its owner. Other versions may not be visible the user and only a specific system command issued will make those versions visible. We also allows a file can be opened multiple times. This is troublesome as later open operations can overwrite changes made by previous file descriptors. One possible fix is to only download copy file from server when there is no local cache copy or implement a copy-on-write mechanism: each open will lead to a unique version of the file and it is up to user to resolve potential conflicts.

4.1 Evaluation

In this section, we compare our network file system with NFS under three workloads.

4.1.1 Experiment Setup

For our network file system, we use **thoothukudi-lom** as client and **erode-lom** as server. On the client side, our system is mounted under **/tmp/barfs** and the local cache is under **/tmp/barfs_cache**. For NFS, to make a fair comparison, we also use memory as storage location for NFS server and we require the NFS to reply requests only after changes have been committed to stable storage (**sync**) [10]. Our **/etc/exports** looks like below:

```
1 /tmp/nfs *(rw, sync, no_subtree_check, no_root_squash)
```

On the client side, we mount our NFS client under **/tmp/nfs** as following:

```
1 sudo mount -t nfs -o sync 192.168.1.120:/tmp/nfs /tmp/nfs
```

As pointed out earlier, 192.168.1.120 refers to **thoothukudi-lom**.

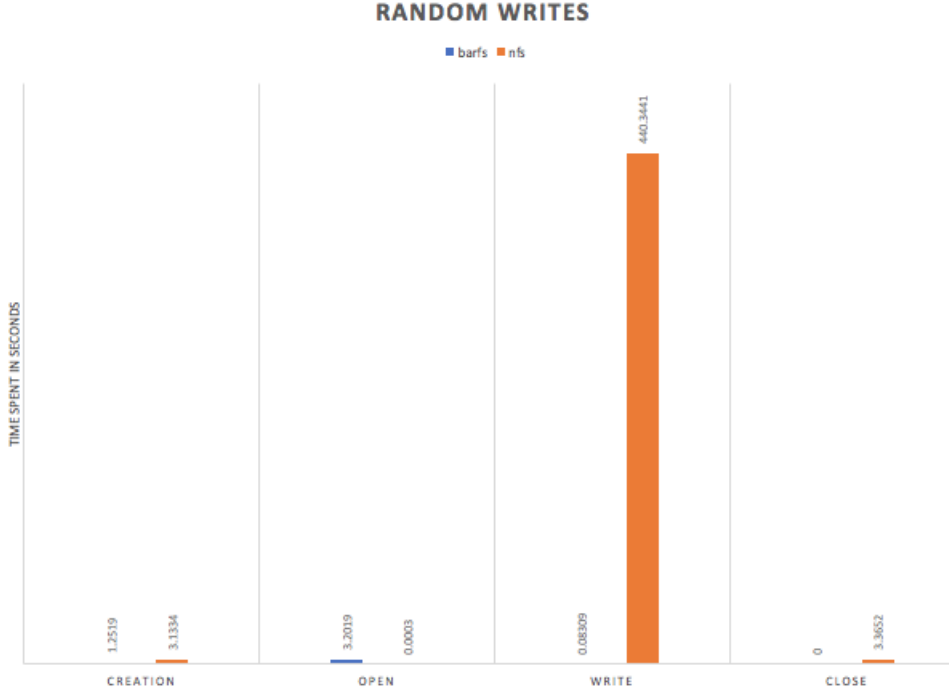


Figure 1: Performance of NFS and our network file system under random writes workload

4.1.2 Workloads and Results

We use three workloads to compare our network file system with NFS. We measure the total amount of time spent on file creation (**creation**), file open (**open**), file write (**write**), and file close (**close**). The detailed description of workloads and corresponding results follow:

1. Random writes. We create a 100MB file with random bytes in it. Then, we perform 10^5 4KB random writes to the created file. Note that by random writes, we mean the offset of the file is chosen randomly when write. The performance of our network file system and NFS is shown in Figure 1.

As shown in the figure, we observe our networked file system spent more time when open the file (3 seconds vs. almost instant in NFS). This is expected as we download file from remote each time we open a file. In contrast, NFS seems to only create local copy on open when there is no modification to the file yet. There is a dramatic performance difference when come to write. Compared to our network file system, NFS spent around 440 seconds to finish writes. This difference is reasonable as we require NFS to synchronously update the local change with the remote server whereas ours only performs write into local copy. Our network file system also performs well on close due to the synchronization work need to be done by NFS.

2. Sequential writes. In this workload, we first create an empty file. Then we repeatedly append 4KB data to file (i.e., sequential writes) until the file size reaches 500MB. The performance of two systems is shown in Figure 2.

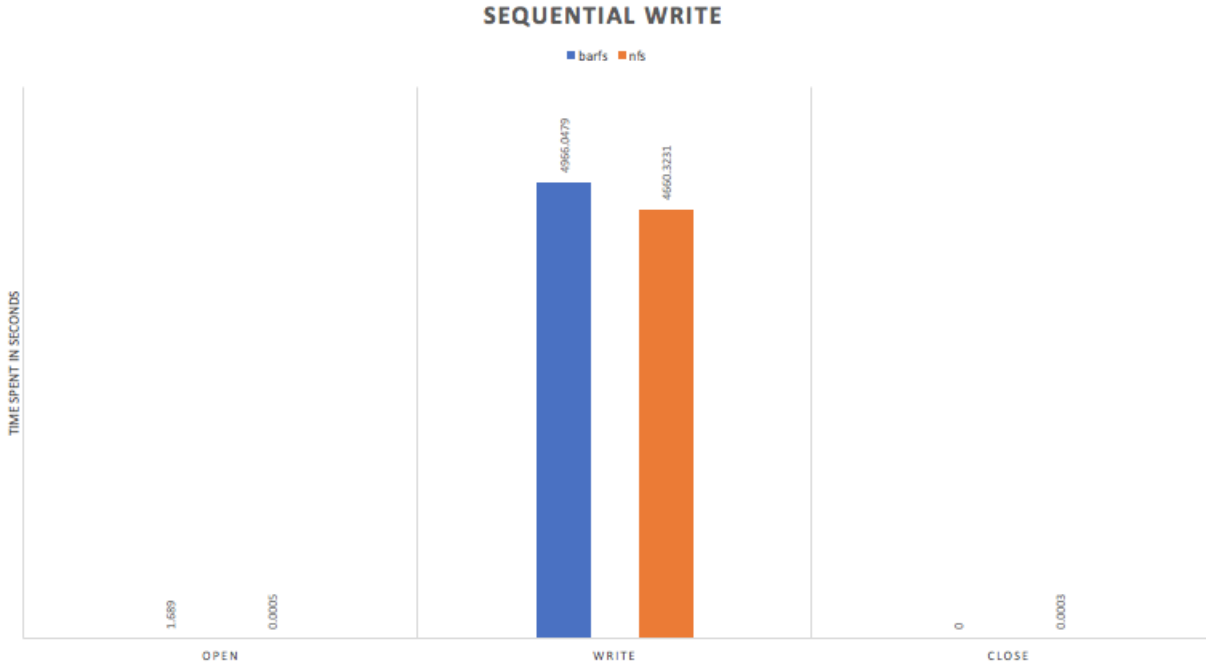


Figure 2: Performance of NFS and our network file system under sequential writes workload

Since we create an empty file, there is no noticeable difference between NFS and our network file system on **open**. In addition, since the major performance bottleneck **release** is called after **close**, our network file system **close** is done almost immediately. For NFS close case, since the write is already synchronized to the server, there is no much work need to be done. The major difference between two systems under this workload is the **write** part. NFS takes less time finish writing than our networked file system (4660 seconds vs. 4966 seconds). The major overhead from our network file system is the FUSE framework which cross kernel boundary multiple times whereas NFS sits in the kernel space. However, the gap is offset by our design: **fsync** after each write forces NFS to synchronize with the remote server. However, in our system, we only synchronize to the remote on **close**. This is a tradeoff we play when we design our system: we trade file availability with the performance. Unlike NFS, when **fsync** is called, our system only synchronizes writes to local disk storage. In other words, without closing file, there is only one copy sitting on the client side. If client is crashed and un-recoverable, the file is gone. However, for NFS, the file is still safe as it already synchronizes with the server. Another motivation for our design of **fsync** is to reduce conflict in the concurrent writes. we think all the writes before **close** are "unofficial" and should not become visible to other users by uploading it to the server.

3. File creation. In this workload, we create 100 4KB files. The performance of two systems is shown in Figure 3. As shown by the figure, our system takes much more timer creating file than NFS does (49 seconds vs. almost 0 seconds). This meets our expectation because of overhead imposed by the remote file downloading. However, NFS will only create files in local cache and synchronize



Figure 3: Performance of NFS and our network file system under file creation workload

changes at the background. When file is closed, NFS will initiate the synchronization to the remote.

As we can see from workloads, our network file system has performance advantage over NFS when write to the file at the sacrifice of potential file loss. However, our network file system performs poorly when `open` operation is involved as remotely downloading file from server is inevitable.

4.1.3 vs. NFS

From the design and actual experiment, our network file system shares great similarities with NFS: both of them are stateless and are designed under the assumption of not many concurrent accesses in a one-server-multiple-clients setting. Both of systems implement last-writer-wins policy.

NFS has advantage in terms of file consistency and expensive remote file downloading is avoided with the help of attribute cache and file cache. Attributes cache periodically checks with remote and ensure the file cache is valid at the best effort (client can still with stale copy if the file is changed during the interval between attribute check). In addition, NFS consistently synchronize local writes to the remote to avoid serious consequence of local client power loss. On the other hand, since we follow the principle of frugality (use the least powerful solution to a given problem) [11], our network file system can reap some performance gain during read and write. Specifically, we only synchronize local copy with remote on close and work with local cache most of time, we avoid overhead during the basic I/O operations.


```
Script started on Fri 01 Mar 2019 04:37:34 PM CST
ESC[1;34m zeyuanhu @ ESC[0;36mHotDog(thoothukudi-lom)ESC[0mESC[1;34m ~ESC[0m
ESC[0;36m Fri Mar 01 16:37:34 $ ESC[0;39mstrace cat - > new_file
execve("/bin/cat", ["cat", "-"], [/* 37 vars */]) = 0
brk(NULL)                                = 0x2280000
access("/etc/ld.so.nohwcap", F_OK)       = -1 ENOENT (No such file or directory)
access("/etc/ld.so.preload", R_OK)       = -1 ENOENT (No such file or directory)
```

Raw output of script, which contains Unicode character.

```
Script started on Fri 01 Mar 2019 04:37:34 PM CST
zeyuanhu @ HotDog(thoothukudi-lom) ~
Fri Mar 01 16:37:34 $ strace cat - > new_file
execve("/bin/cat", ["cat", "-"], [/* 37 vars */]) = 0
brk(NULL)                                = 0x2280000
access("/etc/ld.so.nohwcap", F_OK)       = -1 ENOENT (No such file or directory)
access("/etc/ld.so.preload", R_OK)       = -1 ENOENT (No such file or directory)
```

Script output after cleanup

Figure 4: script output before and after cleanup

5 System Tools Exercise

5.1 strace

`script` command allows user to record terminal printout into a file [?]. Per the lab instruction, we use `strace` to trace the syscalls and signals of a target process [?]. In our case, we trace the process involving `cat`. One thing I notice is that `script` contains some unicode as shown in Figure 4. Thus, we use the following code to clean up the output:

```
1 cat $FILE | perl -pe 's/\e([^\[\]]|\[.?[a-zA-Z]|\[.??\a)//g' | col -b > $FILE-processed
```

The result is shown in Figure 4 ⁴.

5.2 lsof

`lsof` lists all open files [?]. `lsof | grep /dev` shows all the open devices used by user-space programs. On our machine, we have the following opened devices:

- `/dev/null`: null device
- `/dev/pts/*` and `/dev/tty`: terminal devices
- `/dev/urandom`: kernel random number source device
- `/dev/ptmx`: a character file to create a pseudoterminal master

⁴raw output and cleanup output comes with the report as `session.record` and `session.record-processed` respectively

6 Network Tools

`ifconfig` command lists all the network interfaces the machine is using to communicate externally. On our machine, interface for Ethernet is `eno1`. We can find IP address, gateway address, and subnet mask from the output.

`tcpdump` command can dump traffic on a network interface. We use the `tcpdump` output provided by the lab instruction to answer the questions below.

a. Are DHCP messages sent over UDP or TCP?

We use `tcpdump -nn -r tcpdump.out.1 | grep -i dhcp` to filter out the DHCP messages from the dump. `-nn` ensures that we can see the actual port number instead of the port name. One line of the output is:

```
10:19:24.525962 IP 0.0.0.0.68 > 255.255.255.255.67: BOOTP/DHCP, Request from
a8:20:66:3b:66:51, length 300
```

The first field shows the time that the packet was traveling. The second field shows the source host address and port, followed by the destination host address and port. The third field shows the protocol the packet was using. From DHCP [?], we know DHCP messages sent over UDP. As shown by the printout, the messages are sent between port 68 (client) and port 67 (server).

b. What is the link-layer (e.g., Ethernet) address of your host? (Feel free to obscure the last couple bytes for privacy's sake)

We use the same `tcpdump` command as above with extra `-e` option to show link-layer header. The following printout contains DHCP messages for acquiring IP address:

```
10:19:24.525962 a8:20:66:3b:66:51 > ff:ff:ff:ff:ff:ff, ethertype IPv4 (0x0800
), length 342: 0.0.0.0.68 > 255.255.255.255.67: BOOTP/DHCP, Request from
a8:20:66:3b:66:51, length 300
10:19:24.566258 00:21:9b:fb:61:0c > a8:20:66:3b:66:51, ethertype IPv4 (0x0800
), length 342: 128.83.158.2.67 > 128.83.158.160.68: BOOTP/DHCP, Reply,
length 30
```

From the printout we can see that the link-layer address (MAC address) of the host is `a8:20:66:3b:66:51`.

c. What is the IP address of your DHCP server?

From the printout above, we can see the IP address of DHCP server is `128.83.158.2` and the new IP address acquired from DHCP server is `128.83.158.160`.

d. What is the purpose of the DHCP release message?

DHCP release message is used to release IP address.

e. Does the DHCP server issue an acknowledgment of receipt of the client's DHCP request?

DHCP server does not issue an acknowledgment of receipt of the client's release message.

f. What would happen if the client's DHCP release message is lost?

If DHCP release message is lost, the DHCP server has to wait for the lease to timeout before assigning it to other clients.

References

- [1] "Linker and libraries guide." <https://docs.oracle.com/cd/E19957-01/806-0641/6j9vuqujs/index.html#chapter6-71736>.
- [2] "How programs get run: Elf binaries." <https://lwn.net/Articles/631631/>.
- [3] "personality(2) - linux man page." <http://man7.org/linux/man-pages/man2/personality.2.html>.
- [4] ".bss." <https://en.wikipedia.org/wiki/.bss>.
- [5] "System v application binary interface." <https://software.intel.com/sites/default/files/article/402129/mpx-linux64-abi.pdf>.
- [6] B. K. R. Vangoor, V. Tarasov, and E. Zadok, "To {FUSE} or not to {FUSE}: Performance of user-space file systems," in *15th {USENIX} Conference on File and Storage Technologies ({FAST} 17)*, pp. 59–72, 2017.
- [7] "stat(2) - linux man page." <https://linux.die.net/man/2/stat>.
- [8] "stat(1) - linux man page." <https://linux.die.net/man/1/stat>.
- [9] "open(2) - linux man page." <http://man7.org/linux/man-pages/man2/open.2.html>, 2018.
- [10] "exports(5) - linux man page." <http://man7.org/linux/man-pages/man5/exports.5.html>.
- [11] H. Massalin and C. Pu, "Threads and input/output in the synthesis kernel," in *ACM SIGOPS Operating Systems Review*, vol. 23, pp. 191–201, ACM, 1989.