

Major Project - OS

Team:

Swarna Swapna Ilamathy - si2150 - Anubis submission id : 1dbe14c2411724d6c759

Vikram Badhan - vb2174

TASKS:-

1. Get Familiar with Linux strace

i) Executing the “strace echo hello” command shows us the list of system calls. The same can be seen in the screenshot below:-

(ii) Executing the “strace ls” command shows us the following list of system calls:-

```

Problems      anubis@anubis-ide:~/playground-3d0d78ee-vb2174 ×
mprotect(0x7f1132a32000, 4096, PROT_READ) = 0
mprotect(0x7f1132a3c000, 4096, PROT_READ) = 0
mprotect(0x7f1132ad4000, 4096, PROT_READ) = 0
mprotect(0x7f1132cc3000, 4096, PROT_READ) = 0
mprotect(0x55fdc53d000, 4096, PROT_READ) = 0
mprotect(0x7f1132fcfc000, 4096, PROT_READ) = 0
munmap(0x7f1132cc9000, 34804)           = 0
set_tid_address(0x7f1132a156e0)         = 659
set_robust_list(0x7f1132a156e0, 24)     = 0
rt_sigaction(SIGRTMIN, {sa_handler=0x7f1132a1d690, sa_mask=[], sa_flags=SA_RESTORER|SA_SIGINFO, sa_restorer=0x7f1132a2a140}, NULL, 8) = 0
rt_sigaction(SIGRT_1, {sa_handler=0x7f1132a1d730, sa_mask=[], sa_flags=SA_RESTORER|SA_RESTART|SA_SIGINFO, sa_restorer=0x7f1132a2a140}, NULL, 8) = 0
rt_sigprocmask(SIG_BLOCK, [RTMIN RT_1], NULL, 8) = 0
prlimit64(0, RLIMIT_STACK, NULL, {rlim_cur=8192*1024, rlim_max=RLIM64_INFINITY}) = 0
statfs("/sys/fs/selinux", 0xfffff8347d410) = -1 ENOENT (No such file or directory)
statfs("/selinux", 0xfffff8347d410)       = -1 ENOENT (No such file or directory)
brk(NULL)                                = 0x55fdc686000
brk(0x55fdc6a7000)                      = 0x55fdc6a7000
openat(AT_FDCWD, "/proc/filesystems", O_RDONLY|O_CLOEXEC) = 3
fstat(3, {st_mode=S_IFREG|0444, st_size=0, ...}) = 0
read(3, "nodev\tsysfs\nnodev\vtroot\tnodev\tr\...", 1024) = 379
read(3, "", 1024)                         = 0
close(3)                                  = 0
access("/etc/selinux/config", F_OK)        = -1 ENOENT (No such file or directory)
openat(AT_FDCWD, "/usr/lib/locale/locale-archive", O_RDONLY|O_CLOEXEC) = 3
fstat(3, {st_mode=S_IFREG|0644, st_size=3041456, ...}) = 0
mmap(NULL, 3041456, PROT_READ, MAP_PRIVATE, 3, 0) = 0x7f113272d000
close(3)                                  = 0
ioctl(1, TCGETS, {B38400 opost isig icanon echo ...}) = 0
ioctl(1, TIOCGWINSZ, {ws_row=40, ws_col=148, ws_xpixel=0, ws_ypixel=0}) = 0
openat(AT_FDCWD, ".", O_RDONLY|O_NONBLOCK|O_CLOEXEC|O_DIRECTORY) = 3
fstat(3, {st_mode=S_IFDIR|0755, st_size=4096, ...}) = 0
getdents64(3, 0x55fdc68c9f0 /* 15 entries */, 32768) = 472
getdents64(3, 0x55fdc68c9f0 /* 0 entries */, 32768) = 0
close(3)                                  = 0
fstat(1, {st_mode=S_IFCHR|0620, st_rdev=makedev(0x88, 0), ...}) = 0
write(1, "compile_flags.txt Dockerfile k\...", 98compile_flags.txt Dockerfile kernel launch.json tmpl LICENSE Makefile mkfs README.md user ) = 98
close(1)                                  = 0
close(2)                                  = 0
exit_group(0)                            = ?

```

The above screenshots show us the list of system calls when we execute the “strace ls” command. A system call is a programmatic means for a program to ask the kernel for a service, and strace is a useful tool for tracing the thin layer between user programs and the Linux kernel.

Given below are four amongst the many system calls that we got after running the strace command:-

a) “Nmap” system call:-

“nmap” stands for “Network Map” which is an open source tool which was basically designed for security auditing and network exploration. It was created to quickly scan big networks, but also works well with single hosts. Nmap analyzes raw IP packets in unique ways to figure out which host is currently using the network, which Operating system they are running, the number firewalls being used,etc.

b) “read’ system call:-

“read” system call is a type of call which is used to read from a file descriptor. It reads the input typed by the person on the keyboard and stores it in an array(buffer). It can only read 10 bytes at max, after that no character will be read.

c) “fstat” system call:-

“fstat” system call is a type of call which fetches us the status of a file which comprises all the useful information related to the file. The information can be about a file being a directory or a regular file, it may also tell us user Id of the file and it may also tell us other information like size, access, modification count of the file.

d) “openat” system call:-

“Openat” system call is a type of call which works the same way as “open” system call except for a few differences. In Open system call a pathname is specified and based on that the file is opened whereas an Openat system call allows you to specify a relative pathname that is interpreted relative to another directory represented by a file descriptor.

2. Building strace in xv6

i) Command : strace on

“strace on” command will set the trace flag to be on and while executing, since the trace flag is set, the syscall function will print the trace information of each process when each system call gets called by the process. The system call information is printed in the following format : pid, process name, system call name and return value of the system call.

```
$ strace on
$ echo hello
TRACE: pid = 3 | process name = sh | syscall = exec
hTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
eTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
oTRACE: pid = 3 | process name = echo | syscall = write | return val = 1

TRACE: pid = 3 | process name = echo | syscall = write | return val = 1
TRACE: pid = 3 | process name = echo | syscall = exit
```

ii) Command : strace off

“strace off” command will unset the trace flag and only if the flag is set, the syscall function in syscall.c will write the trace information. Unsetting the flag will switch off the trace.

```
$ strace on
$ echo hello
TRACE: pid = 3 | process name = sh | syscall = exec
hTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
eTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 3 | process name = echo | syscall = write | return val = 1
oTRACE: pid = 3 | process name = echo | syscall = write | return val = 1

TRACE: pid = 3 | process name = echo | syscall = write | return val = 1
TRACE: pid = 3 | process name = echo | syscall = exit
$ strace off
$ echo hello
hello
```

iii) Command strace run

"strace run" command will provide trace information only for the command given along with the strace run command. In this example, strace run echo hi, strace information will be shown only for echo hi

```
$ strace run echo hi
TRACE: pid = 2 | process name = sh | syscall = exec
hTRACE: pid = 2 | process name = echo | syscall = write | return val = 1
iTRACE: pid = 2 | process name = echo | syscall = write | return val = 1
|
TRACE: pid = 2 | process name = echo | syscall = write | return val = 1
TRACE: pid = 2 | process name = echo | syscall = exit
```

The next command won't have any strace information as shown below.

```
$ echo hi
hi
$
```

v) Trace child process

To show how the strace information will be when we fork the child process, I have forked two child process and shown how a child process is forked from a parent process and the system calls involved in that.

```
$ strace on
$ test_trace_chil
TRACE: pid = 9 | process name = sh | syscall = exec
TRACE: pid = 9 | process name = test_trace_chil | syscall = fork | return val = 10
Executing Child 1
TRACE: pid = 9 | process name = test_trace_chil | syscall = fork | return val = 11
PTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
Executing Child 2
aTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
rTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
nTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
tTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
lTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
wTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
aTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
iTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
tTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
iTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
nTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
gTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1

TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
TRACE: pid = 9 | process name = test_trace_chil | syscall = wait | return val = 10
TRACE: pid = 9 | process name = test_trace_chil | syscall = wait | return val = 11
CTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
hTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
iTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
```

```

dTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
rTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
nTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
cTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
oTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
mTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
pTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
lTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
tTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
dTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1

TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
PTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
aTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
rTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
nTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
tTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
ETRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
xTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
cTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
uTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1

tTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
iTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
nTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
gTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1

TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
PTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
aTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
rTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
nTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
tTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
eTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
xTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
iTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
tTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
iTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
nTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
gTRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
. TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1

TRACE: pid = 9 | process name = test_trace_chil | syscall = write | return val = 1
TRACE: pid = 9 | process name = test_trace_chil | syscall = exit

```

3. Building option for strace

To implement the options for trace, we have implemented a system call to set the flag in proc.c for proc->option for the specified option. Based on the flag information from proc->option

i) Option: -e <system call name>

For option e, we have set the flag proc->option as 1 and also set filter name in the proc->flag in the kernel space. So, we can access the flag name and option set and process based on that.

For this, only the system call mentioned in the command will be shown in the trace information.

```
$ strace -e write
$ echo hello
TRACE: pid = 12 | process name = sh | syscall = exec
hTRACE: pid = 12 | process name = echo | syscall = write | return val = 1
eTRACE: pid = 12 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 12 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 12 | process name = echo | syscall = write | return val = 1
oTRACE: pid = 12 | process name = echo | syscall = write | return val = 1

TRACE: pid = 12 | process name = echo | syscall = write | return val = 1
$ echo hello
TRACE: pid = 13 | process name = sh | syscall = exec
hTRACE: pid = 13 | process name = echo | syscall = write | return val = 1
eTRACE: pid = 13 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 13 | process name = echo | syscall = write | return val = 1
lTRACE: pid = 13 | process name = echo | syscall = write | return val = 1
oTRACE: pid = 13 | process name = echo | syscall = write | return val = 1

TRACE: pid = 13 | process name = echo | syscall = write | return val = 1
TRACE: pid = 13 | process name = echo | syscall = exit
```

ii) Option: -s

Set the proc->option for 2 and based on that only the successful system call information will be printed.

```
$ strace -s
$ echo hi
TRACE: pid = 15 | process name = sh | syscall = exec
hTRACE: pid = 15 | process name = echo | syscall = write | return val = 1
iTRACE: pid = 15 | process name = echo | syscall = write | return val = 1

TRACE: pid = 15 | process name = echo | syscall = write | return val = 1
```

iii) Option: -f

Set the proc->option for 2 and based on that only the failed system call information will be printed.

```
$ strace -f
$ cat temp
TRACE: pid = 22 | process name = cat | syscall = open | return val = -1
```

iv) Option runs only once

From the below screenshot you can see that the strace option command runs only once. The next time you execute a command, all the system calls are shown, not just the failed system calls.

```

$ strace -f
$ cat temp
TRACE: pid = 3 | process name = cat | syscall = open | return val = -1
cat: cannot open temp
$ cat temp
TRACE: pid = 4 | process name = sh | syscall = exec
TRACE: pid = 4 | process name = cat | syscall = open | return val = -1
cTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
aTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
tTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
:TRACE: pid = 4 | process name = cat | syscall = write | return val = 1
:TRACE: pid = 4 | process name = cat | syscall = write | return val = 1
cTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
aTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
nTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
nTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
bTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
tTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
:TRACE: pid = 4 | process name = cat | syscall = write | return val = 1
bTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
bTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
eTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
nTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
:TRACE: pid = 4 | process name = cat | syscall = write | return val = 1
tTRACE: pid = 4 | process name = cat | syscall = write | return val = 1
eTRACE: pid = 4 | process name = cat | syscall = write | return val = 1

```

5. Applications of Strace

i) executed a memory leak command to see trace information on xv6

```

$ trace_unexpected
TRACE: pid = 3 | process name = sh | syscall = exec
mTRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
eTRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
mTRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
:TRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
tTRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
eTRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
sTRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
tTRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1

TRACE: pid = 3 | process name = trace_unexpecte | syscall = write | return val = 1
TRACE: pid = 3 | process name = trace_unexpecte | syscall = getpid | return val = 3
TRACE: pid = 3 | process name = trace_unexpecte | syscall = fork | return val = 4

allocuvvm out of memory
allocuvvm out of memory
couldn't allocate mem?!!
TRACE: pid = 3 | process name = trace_unexpecte | syscall = wait | return val = 4

```

ii) memory leak on linux

iii) As we already know that strace is a diagnostic tool which is used to debug a program when we do not have access to the source code of that program. In operating systems like Linux, strace function is already inbuilt for the user to take advantage of as opposed to XV6 where there is no strace function present. When we run strace on Linux Operating system, we are able to view information about the various system calls running on the OS in much greater detail.

For example:- Using the “strace -c pwd” command we can get a summary of the system calls. It shows us the various system calls on the OS ordered by the time spent on each one of them, the number of times each call was made and the percentage of time each used in running the command.

And so, we notice that the Linux strace provides us with additional information about the various system calls in our Operating System and so it is quite better at debugging the problem.