



OS Lab 1

System Calls

王世烜

PB20151796

October 26, 2022

Part 1: 实验要求

本次实验要求我们完成 System Calls 的某些功能的添加，了解用户态与内核态间的联系。

简要介绍一下系统调用：

系统调用的功能就是用户态的程序需要执行等级更高的指令时，需要先切换到内核态，在内核态执行完后，再切换到用户态。过程如下图所示：

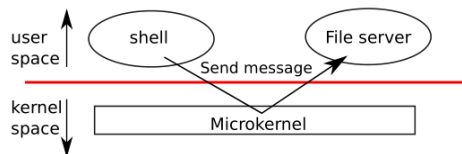


图 1: A microkernel with a file-system server

切换时，使用的是 cpu 特定的中断符号。所有的系统调用都是通过该中断进入到内核态中，并将所需要的执行的指令码传入到中断内核态中。内核态的程序取出指令码，执行相应的指令，然后返回到用户态。

主要涉及到的文件有：trace.c user.h usys.pl syscall.c sysproc.c proc.h proc.c kalloc.c defs.h

Part 2: System call tracing

2.1 问题简述

该部分需要我们实现一个功能 trace，其作用是根据参数 mask 追踪相应的系统调用。用户态中的代码 user/trace.c 已经给出，需要我们实现的是内核态中的 trace 功能。

2.2 实验步骤

(1) Makefile

首先根据实验提示，首先在 Makefile 文件中的 UPROGS 部分中添加 \$U/_trace:

```
1 UPROGS=\
2     $U/_cat\
3     $U/_echo\
4     $U/_forktest\
5     $U/_grep\
6     $U/_init\
7     $U/_kill\
8     $U/_ln\
9     $U/_ls\
10    $U/_mkdir\
11    $U/_rm\
12    $U/_sh\
13    $U/_stressfs\
14    $U/_usertests\
15    $U/_grind\
16    $U/_wc\
17    $U/_zombie\
18    $U/_trace \
```

(2) 添加声明

可以观察到 user/trace.c 中调用了 trace() 函数，但是它并未被定义，所以在 user/user.h 中加入声明：int trace(int);

同时根据给出的提示，为了生成进入中断的汇编文件，需要在 user/usys.pl 添加进入内核态的入口函数的声明，以便使用中断指令进入内核态；entry("trace");

```
1 sub entry {
2     my $name = shift;
3     print ".global _$name\n";
4     print "${name}:\n";
5     print "_li_a7, _SYS_${name}\n";
6     print "_ecall\n";
7     print "_ret\n";
8 }
```

根据该文件的内容可以判断，汇编语言将指令码放到了 a7 寄存器中。在内核态 kernel/syscall.c 的 syscall 函数中，取出寄存器中的指令码，然后调用对应的函数。

所以，在 kernel/syscall.h 中添加相应指令码。

```
1 // System call numbers
```

```

2  #define SYS_fork      1
3  #define SYS_exit      2
4  #define SYS_wait      3
5  #define SYS_pipe      4
6  #define SYS_read      5
7  #define SYS_kill      6
8  #define SYS_exec      7
9  #define SYS_fstat      8
10 #define SYS_chdir      9
11 #define SYS_dup       10
12 #define SYS_getpid    11
13 #define SYS_sbrk      12
14 #define SYS_sleep     13
15 #define SYS_uptime    14
16 #define SYS_open      15
17 #define SYS_write     16
18 #define SYS_mknod     17
19 #define SYS_unlink    18
20 #define SYS_link      19
21 #define SYS_mkdir     20
22 #define SYS_close     21
23 #define SYS_trace     22

```

(3) 实现 sys_trace() 函数

根据提示，在 kernel/proc.h struct proc 中 (即 PCB) 加入一个变量 mask, 用于储存用户态传进来的 mask。

```

1  struct proc {
2      struct spinlock lock;
3
4      // p->lock must be held when using these:
5      enum procstate state;          // Process state
6      struct proc *parent;           // Parent process
7      void *chan;                    // If non-zero, sleeping on chan
8      int killed;                    // If non-zero, have been killed
9      int xstate;                     // Exit status to be returned to parent's wait
10     int pid;                         // Process ID
11
12     // these are private to the process, so p->lock need not be held.
13     uint64 kstack;                   // Virtual address of kernel stack
14     uint64 sz;                       // Size of process memory (bytes)
15     pagetable_t pagetable;           // User page table
16     struct trapframe *trapframe;     // data page for trampoline.S

```

```

17 struct context context;      // swtch() here to run process
18 struct file *ofile[NOFILE]; // Open files
19 struct inode *cwd;           // Current directory
20 char name[16];               // Process name (debugging)
21 int mask;
22 };

```

然后在 kernel/sysproc.c 中添加 sys_trace() 函数为 mask 赋值为该字段进行赋值，赋值的 mask 为系统调用传过来的参数，放在了 a0 寄存器中：

```

1 uint64
2 sys_trace(void)
3 {
4     int mask;
5     if(argint(0, &mask) < 0)
6         return -1;
7     myproc()->mask = mask;
8     return 0;
9 }

```

其中 argint(int n, int *ip) 函数位于 kernel/syscall.c 中，作用是将 trapframe 中 a_n 寄存器中的值取出来。myproc() 函数的作用是获取当前进程的 PCB。

(4) 追踪子进程

需要跟踪所有 trace 进程下的子进程，在 kernel/proc.c 的 fork() 代码中，添加子进程的 mask。

```

1 int
2 fork(void)
3 {
4     int i, pid;
5     struct proc *np;
6     struct proc *p = myproc();
7
8     // Allocate process.
9     if((np = allocproc()) == 0){
10         return -1;
11     }
12
13     // Copy user memory from parent to child.
14     if(uvmcopy(p->pagetable, np->pagetable, p->sz) < 0){
15         freeproc(np);
16         release(&np->lock);
17         return -1;
18     }

```

```

19     np->sz = p->sz;
20
21     np->parent = p;
22
23     //copy the mask
24     np->mask = p->mask;
25
26     .
27     .
28     .
29
30     release(&np->lock);
31
32     return pid;
33 }

```

(5) 打印

先叙述一下代码运行的流程：user/trace.c 中调用 trace() 函数，引发中断，user/usys.pl 生成汇编文件 usys.S，是用户态系统调用接口，将 SYS_trace 的值放入 a7 寄存器，然后使用 ecall 命令进入内核，跳转到 kernel/syscall.c 中的 syscall() 函数处运行，于是可以这样修改：

```

1 void syscall(void)
2 {
3     int num;
4     char *syscalls_name[] = {
5         [SYS_fork]  "fork",
6         [SYS_exit]  "exit",
7         [SYS_wait]  "wait",
8         [SYS_pipe]  "pipe",
9         [SYS_read]  "read",
10        [SYS_kill]   "kill",
11        [SYS_exec]   "exec",
12        [SYS_fstat]  "fstat",
13        [SYS_chdir]  "chdir",
14        [SYS_dup]    "dup",
15        [SYS_getpid] "getpid",
16        [SYS_sbrk]   "sbrk",
17        [SYS_sleep]  "sleep",
18        [SYS_uptime] "uptime",
19        [SYS_open]   "open",
20        [SYS_write]  "write",
21        [SYS_mknod]  "mknod",

```

```

22     [SYS_unlink]  "unlink",
23     [SYS_link]   "link",
24     [SYS_mkdir]  "mkdir",
25     [SYS_close]  "close",
26     [SYS_trace]  "trace",
27 };
28 struct proc *p = myproc();
29
30 num = p->trapframe->a7;
31 if (num > 0 && num < NELEM(syscalls) && syscalls[num])
32 {
33     p->trapframe->a0 = syscalls[num]();
34     if ((1 << num) & p->mask)
35     {
36         printf("%d: %s %s->%d\n", p->pid, syscalls_name[num],
37             p->trapframe->a0);
38     }
39 }
40 else
41 {
42     printf("%d %s: %s %s %d\n",
43         p->pid, p->name, num);
44     p->trapframe->a0 = -1;
45 }
46 }

```

并且在该函数上面的 `static uint64 (*syscalls[])(void)` 中添加 `[SYS_trace] sys_trace`，在上面加入 `extern uint64 sys_trace(void)`；用于将 `sys_trace(void)` 函数扩展进来。这样，即完成了第一部分，实现了 `trace` 的功能。

Part 2: Sysinfo

2.1 问题简述

该部分让我们实现一个系统调用 `sysinfo`，它收集有关正在运行的系统信息，记录空闲内存的字节数与状态不是 `UNUSED` 的进程数。

2.2 实验步骤

(1) Makefile

同上，添加 `$U/_sysinfotest` 到 Makefile 文件中

(2) 添加声明

同样需要添加一些声明才能进行编译，启动 qemu。需要以下几步：

在 user/user.h 文件中加入函数声明，同时添加结构体声明；

在 user/usys.pl 添加进入内核态的入口函数的声明；

同时在 kernel/syscall.h 中添加系统调用的指令码。

(3) 获取内存信息

可以在 kernel/sysinfo.h 中查看结构体 struct sysinfo，其中只有两个字段，一个是保存空闲内存信息，一个是保存正在运行的进程数目。

两个字段的的信息都需要自己写函数调用来获取，先来获取内存信息。内存信息的处理都写在 kernel/kalloc.c 文件中了，内存信息以链表的形式存储，每个节点存储一个物理内存页。

从 kfree 函数中可以发现，每次创建一个页时，将其内容初始化为 1，然后将它的下一个节点指向当前节点的 freelist，更新 freelist 为这个新创建的页。也就是说，freelist 指向最后一个可以使用的内存页，它的 next 指向上一个可用的内存页。

因此，我们可以通过遍历所有的 freelist 来获取可用内存页数，然后乘上页大小即可。

```
1 uint64 free_mem(void)
2 {
3     struct run *r = kmem.freelist;
4     uint64 n = 0;
5     while (r) {
6         n++;
7         r = r->next;
8     }
9     return n * PGSIZE;
10 }
```

(4) 获取进程数目

所有的进程有关的操作都保存在 /kernel/proc.c 文件中，其中的 proc 数组保存了所有进程。进程有五种状态，我们只需要遍历 proc 数组，计算不为 UNUSED 状态的进程数目即可。

```
1 int proc_size(void)
2 {
3     int i;
4     int n = 0;
5     for (i = 0; i < NPROC; i++)
6     {
7         if (proc[i].state != UNUSED)
8         {
9             n++;
10        }
```

```
10     }
11 }
12 return n;
13 }
```

(4) 声明和调用

在 kernel/defs.h 中添加上面这两个函数的声明。

然后在 kernel/sysproc.c 中的 sys_sysinfo 函数进行调用。

需要注意的是，这里使用 copyout 方法将内核空间中，相应的地址内容复制到用户空间中。这里就是将 info 的内容复制到进程的虚拟地址内，具体是哪个虚拟地址，由函数传入的参数决定（addr 读取第一个参数并转成地址的形式）。

```
1  uint64
2  sys_sysinfo(void)
3  {
4      struct sysinfo info;
5      uint64 addr;
6      struct proc* p = myproc();
7      if(argaddr(0, &addr) < 0) {
8          return -1;
9      }
10     info.freemem = free_mem();
11     info.nproc = proc_size();
12     if (copyout(p->pagetable, addr, (char*)&info, sizeof(info)) < 0) {
13         return -1;
14     }
15     return 0;
16 }
```

Part 3: 实验结果

运行 make grade，检验结果：


```

riscv64-linux-gnu-objdump -t kernel/kernel | sed '1,/SYMBOL TABLE/d; s/ .*'
make[1]: 离开目录"/home/ubuntu/xv6-labs-2020"
== Test trace 32 grep ==
$ make qemu-gdb
trace 32 grep: OK (4.6s)
== Test trace all grep ==
$ make qemu-gdb
trace all grep: OK (2.1s)
== Test trace nothing ==
$ make qemu-gdb
trace nothing: OK (2.0s)
== Test trace children ==
$ make qemu-gdb
trace children: OK (23.2s)
== Test sysinfotest ==
$ make qemu-gdb
sysinfotest: OK (7.0s)
== Test time ==
time: FAIL
    Cannot read time.txt
Score: 34/35
make: *** [Makefile:228: grade] 错误 1
ubuntu@VM5153-05VM:~/xv6-labs-2020$

```

图 2: Result

本次实验完成了系统调用某些功能的实现。在实验过程中遇到了以下问题：

- 对于每个文件都不是很理解，无法弄懂其中的原理
- 各种函数调用有些复杂，调试十分不方便
- 对于用户态和内核态的理解十分不深刻，亟待加强