

实验二 添加Linux系统调用

实验目的

- 学习如何使用Linux系统调用：实现一个简单的shell
- 学习如何添加Linux系统调用：实现一个简单的top

实验环境

- OS: Ubuntu 20.04.4 LTS
- Linux内核版本: 4.9.263

实验时间安排

注：此处为实验发布时的安排计划，请以课程主页和课程群内最新公告为准

- 4.1晚实验课，讲解实验 (19:00开始) 及检查
- 4.8晚实验课，检查实验
- 4.15晚实验课，检查实验

补检查分数照常给分，但会有标记记录此次检查未按时完成，标记会在最后综合分数时作为一种参考。

友情提示

- 本次实验以实验一为基础。一些步骤（如如何启动qemu运行Linux内核）不会在实验说明中详述。如果有不熟悉的步骤，请复习实验一。
- 请注意单词拼写、大小写，如common、asmlinkage、user等。
- 在本次实验中，如果你写出的代码出现了数组越界，则会出现很多奇怪的错误（如段错误（Segmentation Fault）、输出其他数组的值或一段乱码（注意，烫烫烫是Visual C的特性，Linux下没有烫烫烫）、其他数组的值被改变等）。提问前请先排查是否出现了此类问题。
- 如果同学们遇到了问题，请先查询在线文档。在线文档地址：
<https://docs.qq.com/sheet/DR1dZTnFRTURHc051>

第一部分 续·Linux与Shell指令教学

1.1 Linux指令

1.1.1 echo

输出一个字符串。用法：`echo string`。

同样也可以输出一个文件的内容，如 `echo file` 将文件file的内容以纯文本形式读并输出。

经常在Shell脚本中使用该指令，以打印指定信息。

1.1.2 top

`ps` 命令可以一次性给出当前系统中进程状态，但使用此方式得到的信息缺乏时效性，并且，如果管理员需要实时监控进程运行情况，就必须不停地执行 `ps` 命令，这显然是缺乏效率的。

为此，Linux 提供了`top`命令。`top` 命令动态地持续监听进程地运行状态，与此同时，该命令还提供了一个交互界面，用户可以根据需要，人性化地定制自己的输出，进而更清楚地了进程的运行状态。直接输入 `top` 即可使用。

参数	含义
-d 秒数	指定 <code>top</code> 命令每隔几秒更新。默认是 3 秒。
-b	使用批处理模式输出。一般和 <code>-n</code> 选项合用，用于把 <code>top</code> 命令重定向到文件中。
-n 次数	指定 <code>top</code> 命令执行的次数。一般和 <code>-b</code> 选项合用。
-p 进程PID	仅查看指定 ID 的进程。
-u 用户名	只监听某个用户的进程。

在 `top` 命令的显示窗口中，还可以使用如下按键，进行以下交互操作：

按键 (注意大小写)	含义
? 或 h	显示交互模式的帮助。
P	按照 CPU 的使用率排序，默认就是此选项。
M	按照内存的使用率排序。
N	按照 PID 排序。
k	按照 PID 给予某个进程一个信号。一般用于中止某个进程，信号 9 是强制中止的信号。
q	退出 <code>top</code> 命令。

1.1.3 sleep

`sleep n` 可以使当前终端暂停n秒。常见于shell脚本中，用于实现两条指令间的等待。

1.1.4 grep

筛选并高亮指定字符串。尝试在终端下运行 `grep a`，它会等待用户输入。若输入一行不带a的字符串并按回车，它什么都不会输出。若输入多行字符串，其中某些行带a，它会筛选出带a的行，并将该行的a高亮后输出。

1.1.5 wc

英文全拼：wordcount

常见用法：`wc [-lw] [filename]`，用于统计字数。

参数	含义
-l	统计行数。
-w	统计字数。
filename	统计指定文件的行数/字数。若不指定，会从标准输入(C/C++的stdin/scanf/cin)处读取数据。

1.2 Shell指令

1.2.1 管道符 |

问题：我们想统计Linux下进程的数量，应该如何解决？

一个很麻烦的解决方法是：先 `ps aux` 输出所有进程的列表，然后复制它的输出，执行 `wc -l`，将之前 `ps` 输出的内容粘贴到标准输入，传递给 `wc` 以统计 `ps` 输出的行数，即进程数量。有没有方法可以免去复制粘贴的麻烦？

管道符 `|` 可以将前面一个命令的标准输出（对应C/C++的stdout/printf/cout）传递给下一个命令，作为它的标准输入（对应C/C++的stdin/scanf/cin）。

- 例1：在终端中运行 `ps aux | wc -l` 可以显示所有进程的数量。
- 例2：先打开Linux下的firefox，然后在终端运行 `ps aux | grep firefox` 可以显示所有进程名含firefox的进程。
- 例3：接例2，在终端运行 `ps aux | grep firefox | wc -l` 可以显示所有进程名含firefox的进程的数量。

思考：`wc -l` 的作用是统计输出的行数。所以 `ps aux | wc -l` 统计出的数字真的是进程数量吗？

1.2.2 重定向符 >/>>/<

重定向符 `>`，`>>`：它可以把前面一个命令的标准输出保存到文件内。如果文件不存在，则会创建文件。`>` 表示覆盖文件，`>>` 表示追加文件。

- 举例：`ps -aux > ps.txt` 可以把当前运行的所有进程信息输出到ps.txt。ps.txt的原有内容会被覆盖。
- 举例：`ps -aux >> ps.txt` 可以把当前运行的所有进程信息追加写到ps.txt。

重定向符 `<`：可以将 `<` 前的指令的标准输入重定向为 `<` 后文件的内容。

- 举例：`wc -l < ps.txt` 可以把ps.txt中记录的信息作为命令的输入，即统计ps.txt中的行数。

1.2.3 分隔符；

子命令：每行命令可能由若干个子命令组成，各子命令由 `;` 分隔，这些子命令会被按序依次执行。

如：`ps -a; pwd; ls -a` 表示：先打印当前用户运行的进程，然后打印当前shell所在的目录，最后显示当前目录下所有文件。

第二部分 实现一个Linux Shell

注意：

- 本部分的主要目的是锻炼大家使用 Linux系统调用编写程序的能力，不会涉及编写一个完整的系统调用。
- 本部分主要是代码填空，费力的部分已由助教完成，代码量不大，大家不必恐慌。需要大家完成的代码已经用注释 **TODO: 标记**，可以通过搜索轻松找到，使用支持TODO高亮编辑器（如vscode装TODO highlight插件）的同学也可以通过高亮找到要添加内容的地方。

2.1 Shell的基本原理

我们在各种操作系统中会遇到各式各样的用于输入命令、基于文字与系统交互的终端，Windows下的cmd和powershell、Unix下的sh、bash和zsh等等，这些都是不同的shell，我们可以看到它们的共同点——基于“命令”（command）与系统交互，完成相应的工作。

Shell处理命令的方式很好理解，每条输入的“命令”其实都对应着一个可执行文件，例如我们输入 `top` 时，shell程序会创建一个新进程并在新进程中运行可执行文件。具体而言，shell程序在特定的文件夹中搜索名为“top”的可执行文件，搜索到之后运行该可执行文件，并将进程的输出定向到指定的位置（如终端设备），这一过程也是我们本次实验要实现的内容。具体来讲：

- **创建新进程**：用到我们课上讲过的 `fork()`。
- **运行可执行文件**：使用课上讲的 `exec()` 系列函数。
 - `exec`开头的函数有很多，如 `execvp`、`execl` 等，可以使用man查看不同exec函数的区别及使用方法，最后选取合适的函数。
 - `exec`中需要指定可执行文件，因为用户的命令给出的只有文件名本身，没有给出文件的所在位置，所以在系统中通常会定义一个环境变量PATH来记录一组文件夹，所有命令对应的可执行文件通常存在于某个文件夹下，找不到则返回不存在。可以使用 `export -p PATH` 来查看你现在使用的shell都会去哪些文件夹下查找可执行程序。
- **结果输出**：命令默认会从STDIN_FILENO（默认值为0）中读输入，并输出到STDOUT_FILENO（默认值为1）中。在shell中也会有其他的输入来源和输出目标，如将一个命令的输出作为另一个命令的输入，在这种场合下就需要进行进程间的通信，比如采用我们课上讲过的管道pipe，执行两个命令的进程一个在写端输出内容，一个在读端读取输入。此外，命令还可以设置从文件输入和输出到文件，比如将执行命令程序中的STDOUT_FILENO使用文件描述符覆盖后，命令的执行结果就输出到文件中。

2.2 Shell内建指令

`cd` 命令可以用来切换shell的当前目录。但需要指出的是，不同于其他命令（比如 `ls`，我们可以在/bin下面找到一个名为 `ls` 的可执行文件），`cd` 命令其实是一个shell内置指令。由于子进程无法修改父进程的参数，所以若不使用内建命令而是fork出一个子进程并且在子进程中exec一个 `cd` 程序，因为子进程执行结束后会回到了父shell环境，而父shell的路径根本没有被改变，最终无法得到期望的结果。同理，不仅是 `cd`，**改变当前shell的参数（如 `source` 命令、`exit` 命令）基本都是由shell内建命令实现的。**

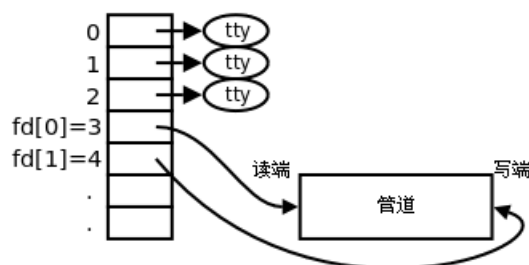
2.3 有关管道的背景知识

“一切皆文件”是Unix/Linux的基本哲学之一。普通文件、目录、I/O设备等在Unix/Linux都被当做文件来对待。虽然他们的类型不同，但是linux系统为它们提供了一套统一的操作接口，即文件的open/read/write/close等。当我们调用Linux的系统调用打开一个文件时，系统会返回一个文件描述符，每个文件描述符与一个打开的文件唯一对应。之后我们可以通过文件描述符来对文件进行操作。管道也是一样，我们可以通过类似文件的read/write操作来对管道进行读写。为便于理解，本次实验使用匿名管道。匿名管道具有以下特点：

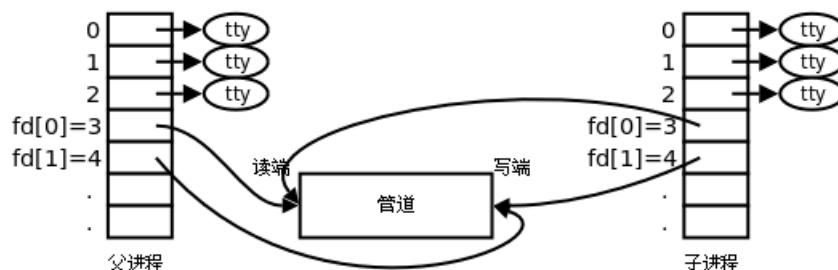
1. 只能用于父子进程等有血缘关系的进程；
2. 写端不关闭，并且不写，读端读完，继续等待，此时阻塞，直到有数据写入才继续（就好比你的C程序在scanf，但你一直什么都不输入，程序会停住）；
 - 尤其地，假如一条管道有多个写端，那么只有在所有写端都关闭之后（管道的引用数降为0），读端才会解除阻塞状态。
3. 读端不关闭，并且不读，写端写满管道buffer，此时阻塞，直到管道有空位才继续；
4. 读端关闭，写端在写，那么写进程收到信号SIGPIPE，通常导致进程异常中止；
5. 写端关闭，读端在读，那么读端在读完之后再读会返回0；
6. 匿名管道的通信通常是一次性的，如果需要反复通信，可以使用命名管道。

一般来说，匿名管道的使用方法是：

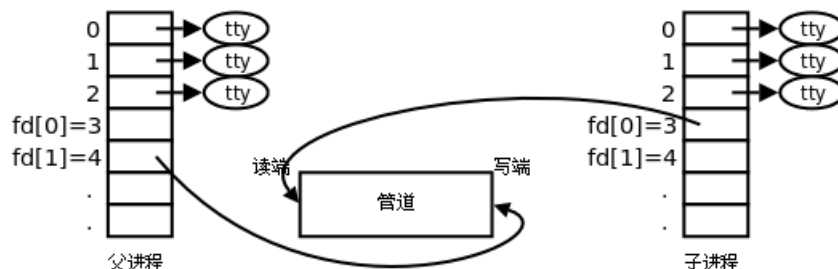
1. 父进程创建管道



2. 父进程fork出子进程



3. 父进程关闭fd[0]，子进程关闭fd[1]



- 首先，父进程调用pipe函数创建一个匿名管道。pipe的函数原型是 `int pipe(int pipefd[2])`。我们传入一个长为2的一维数组 `pipefd`，Linux会将 `pipefd[0]` 设为读端的文件描述符，并将 `pipefd[1]` 设为写端的文件描述符。（注：此时管道已被打开，相当于已调用了 `open` 函数打开文件）但是需要注意：此时管道的读端和写端接在了同一个进程上。如果你此时往 `pipefd[1]` 里写入数据，这些数据可以从 `pipefd[0]` 里读出来。不过这种“我传我自己”(原地tp) 通常没什么意义，我们接下来要把管道应用于进程通信。
- 其次，使用 `fork` 函数创建一个子进程。`fork` 完成之后，数组 `pipefd` 也会被复制。此时，子进程也拥有了对管道的控制权。若目的是父进程向子进程发送数据，那么父进程就是写端，子进程就是读端。我们应该把父进程的读端关闭，把子进程的写端关闭，进而便于数据从父进程流向子进程。

- 如果不关闭子进程的写端，子进程会一直等待（参考2.3.2）。
- 因为匿名管道是单向的，所以如果想实现从子进程向父进程发送数据，就得另开一个管道。
- 父子进程调用 `write` 函数和 `read` 函数写入、读出数据。
 - `write` 函数的原型是：`ssize_t write(int fd, const void * buf, size_t count);`
 - `read` 函数的原型是：`ssize_t read(int fd, void * buf, size_t count);`
 - 如果你要向管道里读写数据，那么这里的 `fd` 就是上面的 `pipefd[0]` 或 `pipefd[1]`。
 - 这两个函数的使用方法在此不多赘述。如有疑问，可以百度。

注意：如果你的数组越了界，或在`read/write`的时候`count`的值比`buf`的大小更大，则会出现很多奇怪的错误（如段错误（Segmentation Fault）、其他数组的值被改变、输出其他数组的值或一段乱码（注意，烫烫烫是Visual C的特性，Linux下没有烫烫烫）等）。提问前请先排查是否出现了此类问题。

2.4 输入/输出的重定向

注：本节描述的是如何将一个程序产生的标准输出转移到其他非标准输出的地方，不特指'>'和'>>'符号。

在使用 '|', '>', '>>', '<' 时，我们需要将程序输出的内容重定向为文件输出或其他程序的输入。在本次实验中，为方便起见，我们将shell作为重定向的中转站。

- 当出现前三种符号时，我们需要把前一指令的标准输出重定向为管道，让父进程（即shell）截获它的标准输出，然后由shell决定将前一指令的输出转发到下一进程、文件或标准输出（即屏幕）。
- 当出现'|' 和 '<' 时，我们需要把后一指令的标准输入重定向为管道，让父进程（即shell）把前一进程被截获的标准输出/指定文件读出的内容通过管道发给后一进程。

重定向使用 `dup2` 系统调用完成。其原型为：`int dup2(int oldfd, int newfd);`。该函数相当于将 `newfd` 标识符变成 `oldfd` 的一个拷贝，与 `newfd` 相关的输入/输出都会重定向到 `oldfd` 中。如果 `newfd` 之前已被打开，则先将其关闭。举例：下述程序在屏幕上没有输出，而在文件输出"hello!goodbye!"。屏幕上不会出现"hello!"和"goodbye!"。

```
#include<stdio.h>
#include<stdlib.h>
#include<unistd.h>
#include<sys/types.h>
#include<sys/stat.h>
#include<fcntl.h>

int main(void)
{
    int fd;

    fd = open("./test.txt", O_RDWR | O_CREAT | O_TRUNC, 0666);

    // 此代码运行在dup2之前，本应显示到屏幕，但实际上是暂时放到屏幕输出缓冲区
    printf("hello!");

    // 程序默认是输出到STDOUT_FILENO即1中，现在我们让1重定向到fd。
    // 也就是将标准输出的内容重定向到文件，因此屏幕输出缓冲区以及后面的printf
    // 语句本应输出到屏幕（即标准输出，fd为STDOUT_FILENO）的内容重定向到文件中。
    dup2(fd, 1);

    printf("goodbye!\n");
    return 0;
}
```

易混淆的地方：向文件/管道内写入数据实际是程序的一个输出过程。

2.5 一些shell命令的例子和结果分析

本部分旨在帮助想要进一步理解真实shell行为的同学，希望下述示例及结果分析可以为同学们带来启发。

这一节中，我们给出一些shell命令的例子，并分析命令的输出结果，让大家能够更直观的理解shell的运行。**我们并不要求实验中实现的shell行为和真实的shell完全一致，但你自己需要能够清楚地解释自己实现的处理逻辑。**但在实现shell时，可以参考真实shell的这些行为。这些实验都是在bash下测试的，大家也可以自己在Ubuntu等的终端中重复这些实验。在这一节的代码中，每行开头为 `$` 的，是输入的shell命令，剩下的行是shell的输出结果。如：

```
$ cd /bin
$ pwd
/bin
```

表示在shell中先运行了 `cd /bin`，然后运行了 `pwd`，第一条命令没有任何输出，第二条命令输出为 `/bin`。`pwd` 命令会显示shell的当前目录，我们先用 `cd` 命令进入了 `/bin` 目录，所以 `pwd` 输出 `/bin`。

2.5.1 多个子命令

由 `;` 分隔的多个子命令，和在多行中依次运行这些子命令效果相同。

```
$ cd /bin ; pwd
/bin
$ echo hello ; echo my ; echo shell
hello
my
shell
```

`echo` 会将它的参数打印到标准输出，如 `echo hello world` 会输出 `hello world`。这个实验用 `;` 分隔了三个shell命令，结果和在三行中分别运行这些命令一致。

2.5.2 管道符

这个实验展示了shell处理管道时的行为，相同的命令在管道中行为可能和单独运行时不同。

```
$ echo hello | echo my | echo shell
shell
```

如上，虽然管道中的上一条命令的输出被重定向至下一条命令的输入，但是因为 `echo` 命令本身不接受输入，所以前两个 `echo` 的结果不会显示。如果实验检查中，你提供了这样的测试样例，是不能证明正确实现了管道的。

```
$ cd /bin ; pwd
/bin
$ cd /etc | pwd
/bin
```

如上，`cd /etc` 并没有改变shell的当前目录，这是因为**管道中的内置命令也是在新的子进程中运行的**，所以不会改变当前进程（shell）的状态。而在**不包含管道的命令中，内置命令在shell父进程中运行，外部命令在子进程中运行**，所以，不包含管道的内置命令能够改变当前进程（shell）的状态。可以用 `type` 命令检查命令是否是内部命令。

```
$ type cd
cd is a shell builtin
$ type cat
cat is hashed (/bin/cat)
```

说明 `cd` 是内置命令，而 `cat` 则是调用 `/bin/cat` 的外置命令。接下来我们测试管道中命令的运行顺序。

```
$ sleep 0.02 | sleep 0.02 | sleep 0.02 | ps | grep sleep
15840 tty1      00:00:00 sleep
$ sleep 0.02 | sleep 0.02 | sleep 0.02 | ps | grep sleep
15845 tty1      00:00:00 sleep
15846 tty1      00:00:00 sleep
15847 tty1      00:00:00 sleep
```

可以看到，运行了两次相同的命令，却得到了不同的结果，多次运行该命令，每次输出的sleep行数不同，从0行到3行都有可能（根据电脑速度和核数，可能需要将sleep后面的数字增大或缩小来重复该实验）。这说明**管道中各个命令是并行执行的**，ps 命令运行时，前面的 sleep 命令可能执行结束，也有可能仍在执行。

2.5.3 重定向

```
$ cd /tmp ; mkdir test ; cd test
$ echo hello > a >> b > c
$ ls
a b c
$ cat a
$ cat b
$ cat c
hello
```

当一个命令中出现多个输出重定向时，虽然所有文件都会被建立，但是只有最后一个文件会真正被写入命令的输出。

```
$ cd /tmp ; mkdir test ; cd test
$ echo hello > out | grep hello
$ cat out
hello
```

当输出重定向和管道符同时使用时，命令会将结果输出到文件中，而管道中的下一个命令将接收不到任何字符。

```
$ cd /tmp ; mkdir test ; cd test
$ > out2 echo hello ; cat out2
hello
```

重定向符可以写在命令前。

2.6 总结：本次实验中shell执行命令的流程

- 第一步：打印命令提示符（类似shell ->）。
- 第二步：把分隔符 ; 连接的各条命令分割开。（多命令选做内容）
- 第三步：对于单条命令，把管道符 | 连接的各部分分割开。
- 第四步：如果命令为单一命令没有管道，先根据命令设置标准输入和标准输出的重定向（重定向选做内容）；再检查是否是shell内置指令：是则处理内置指令后进入下一个循环；如果不是，则fork出一个子进程，然后在fork出的子进程中exec运行命令，等待运行结束。（如果不fork直接exec，会怎么样？）
- 第五步：如果只有一个管道，创建一个管道，并将子进程1的**标准输出重定向到管道写端**，然后fork出一个子进程，根据命令重新设置标准输入和标准输出的重定向（重定向选做内容），在子进程中先检查是否为内置指令，是则处理内置指令，否则exec运行命令；子进程2的**标准输入重定向到管道读端**，同子进程1的运行思路。（注：2.4的样例分析中我们得出，管道的多个命令之间，虽然某个命令可能会因为等待前一个进程的输出而阻塞，但整体是没有顺序执行的，即并发执行。所以我们为了让多个内置指令可以并发，需要在fork出子进程后才执行内置指令）

- 第六步：如果有多个管道，参考第三步，n个进程创建n-1个管道，每次将子进程的标准输出重定向到管道写端，父进程保存对应管道的读端（上一个子进程向管道写入的内容），并使得下一个进程的标准输入重定向到保存的读端，直到最后一个进程使用标准输出将结果打印到终端。（多管道选做内容）
- 第七步：根据第二步结果确定是否有剩余命令未执行，如果有，返回第三步执行（多命令选做内容）；否则进入下一步。（分隔符多命令和管道连接的多命令实现方式上有什么区别？为什么？）
- 第八步：打印新的命令提示符，进入下一轮循环。

2.7 任务目标

代码填空实现一个shell。这个shell不要求在qemu下运行。功能包括：

必做部分

- 实现运行单条命令（非shell内置命令）。
- 支持一条命令中有单个管道符 `|`。
- 实现exit, cd两个shell内置指令。
- 在shell上能显示当前所在的目录。如：`{ustc}shell: [/home/ustc/exp2/] >`（后面是用户的输入）

二选一部分

- 实现子命令符 `;` 和重定向符 `>`, `>>`, `<`。
- 支持一条命令中有多个管道符 `|`。如果你确信你的多管道功能可以正常实现单管道功能，代码填空里的单管道可以不做。

其他说明

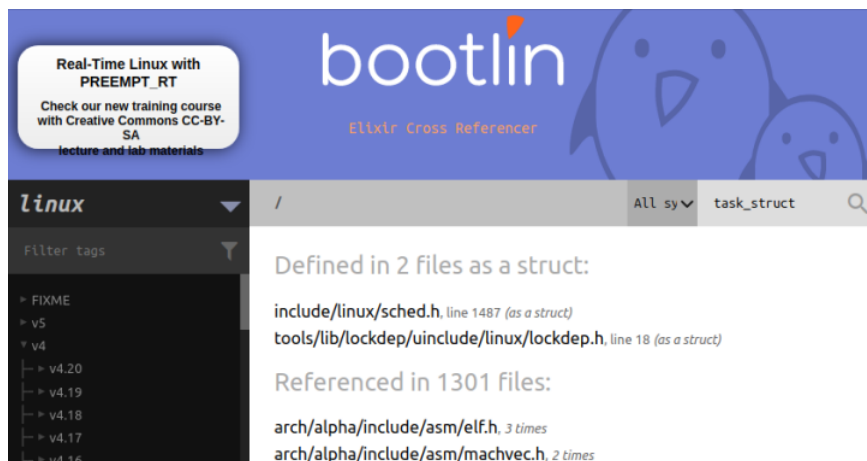
- 我们使用的是Linux系统调用，请不要尝试在Windows下运行自己写的shell程序。那是不可能的。
- 需要自行设计测试样例以验证你的实现是正确的。如测试单管道时使用 `ps aux | wc -l`，与自带的shell输出结果进行比较。
- 不限制分隔符、管道符、重定向符的符号优先级。你可以参考我们代码框架中实现、提示的优先级。
- 我们提供了本次实验使用的系统调用API的范围，请同学们自行查询它们的使用方法。你在实验中可能会用到它们中的一部分。你也可以使用不属于本表的系统调用。
 - read/write
 - open/close
 - pipe
 - dup/dup2
 - getpid/getcwd
 - fork/vfork/clone/wait/exec
 - mkdir/rmdir/chdir
 - exit/kill
 - shutdown/reboot
 - chmod/chown
- 如果你想问如何编译、运行自己写的代码，请参考lab1的2.3.17和2.2.3。
- 得分细则见文档4.2。

第三部分 编写系统调用实现一个Linux top

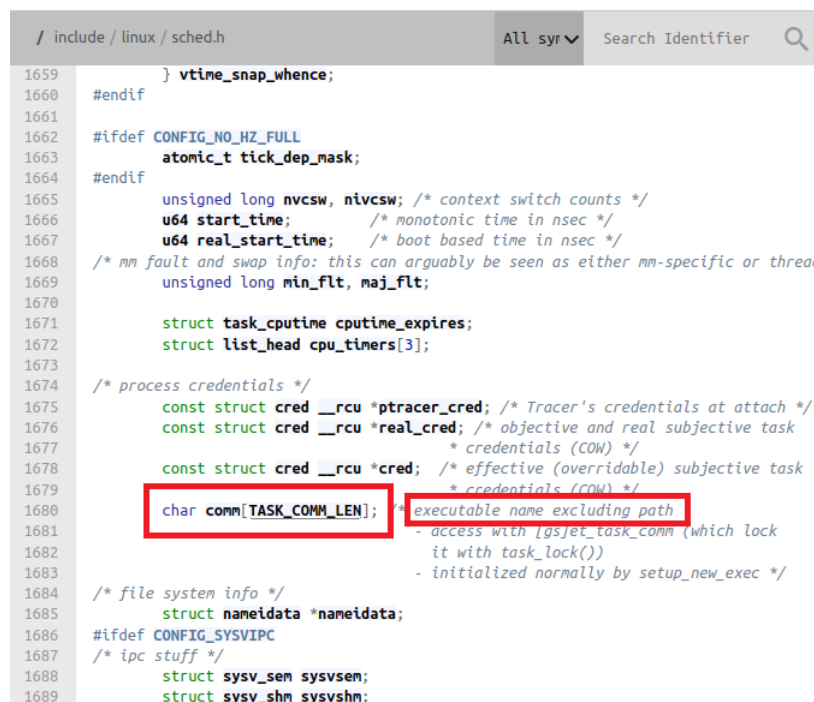
3.0 如何阅读Linux源码

因为直接在VSCode、Visual Studio等本地软件中阅读Linux源码会让电脑运行速度变得很慢，所以我们建议使用在线阅读。一个在线阅读站是 <https://elixir.bootlin.com/linux/v4.9.263/source>。左栏可以选择内核源码版本，右边是目录树或代码，右上角有搜索框。下面举例描述一次查阅Linux源码的过程。

假定我们要查询Linux中进程名最长有多长。首先，我们知道进程在内存中的数据结构是结构体 `task_struct`，所以我们首先查找 `task_struct` 在哪。搜索发现在 `include/linux/sched.h`。（下面那个从名字来看像个附属模块，所以大概可以确定结构体的声明应该在 `sched.h`。）



进入 `sched.h` 查看 `task_struct` 的声明。我们不难猜到，进程名应该是个一维char数组。所以我们开始找 `char` 数组。看了几十个成员之后，我们找到了 `char comm[TASK_COMM_LEN]` 这个东西。右边注释写着 `executable name excluding path`，应该就是这个了。所以它的长度就应该是 `TASK_COMM_LEN`。



于是我们点击一下这个 `TASK_COMM_LEN`，找到这个宏在哪定义。发现定义也在 `sched.h`。点开就能看到宏定义了。至于它的长度到底是多少，留作习题，给大家自己实践查询。

Defined in 3 files as a macro:

include/linux/sched.h, line 316 (as a macro)
samples/bpf/offwaketime_user.c, line 36 (as a macro)
samples/bpf/trace_event_user.c, line 49 (as a macro)

Referenced in 50 files:

arch/arc/kernel/unaligned.c, line 206
arch/arm/kernel/traps.c, line 283
arch/arm64/kernel/traps.c, line 248
arch/unicore32/kernel/traps.c, line 198
... ..

3.1 若干名词解释

3.1.1 用户空间、内核空间

参考PPT。

3.1.2 系统调用

系统调用 (System call, Syscall)是操作系统提供给**用户程序访问内核空间的合法接口**。

系统调用**运行于内核空间，可以被用户空间调用**，是内核空间和用户空间划分的关键所在，它保证了两个空间必要的联系。从用户空间来看，系统调用是一组统一的抽象接口，用户程序**无需考虑接口下面是什么**；从内核空间来看，用户程序不能直接进行敏感操作，所有对内核的操作都必须通过功能受限的接口间接完成，保证了**内核空间的稳定和安全**。

我们平时的编程中**为什么没有意识到系统调用的存在**？这是因为应用程序现在一般**通过应用编程接口 (API)来间接使用系统调用**，比如我们如果想要C程序打印内容到终端，只需要使用C的标准库函数API `printf()` 就可以了，而不是使用系统调用 `write()` 将内容写到终端的输出文件 (`STDOUT_FILENO`) 中。

3.1.3 glibc

glibc是GNU发布的c运行库。glibc是linux系统中最底层的api，几乎其它任何运行库都会依赖于glibc。glibc除了封装linux操作系统所提供的系统服务外，它本身也提供了许多其它一些必要功能服务的实现，其内容包罗万象。glibc内含的档案群分散于系统的目录树中，像支架一般撑起整个操作系统。

`stdio.h` , `malloc.h` , `unistd.h` 等都是glibc实现的封装。

3.2 系统调用是如何执行的

1. 调用应用程序调用**库函数 (API)**

Linux使用的开源标准C运行库glibc (GNU libc) 有一个头文件 `unistd.h` , 其中声明了很多**封装好的**系统调用函数，如 `read/write` 。这些API会调用实际的系统调用。

2. API将**系统调用号存入EAX**，然后**触发软中断**使系统进入内核空间。

32位x86机器使用汇编代码 `int $0x80` 触发中断。具体如何触发软中断可以去参考Linux的实现源码 (如 `arch/x86/entry/entry_32.S`) 。因为这部分涉及触发中断的汇编指令，这部分代码是用汇编语言写的。

3. 内核的中断处理函数根据系统调用号，调用对应的**内核函数** (也就是课上讲的的系统调用)。前面的 `read/write` , 实际上就是调用了内核函数 `sys_read/sys_write` 。

总结：添加系统调用需要**注册中断处理函数和对应的调用号**，内核才能找到这个系统调用，并执行对应的内核函数。

4. 内核函数完成相应功能，将返回值存入EAX，返回到中断处理函数；

注1：系统执行相应功能，调用的是C代码编写的函数，而不是汇编代码，减轻了实现系统调用的负担。系统调用的函数定义在 `include/linux/syscalls.h` 中。因为汇编代码到C代码的参数传递是通过栈实现的，所以，所有系统调用的函数前面都使用了 `asm linkage` 宏，它意味着编译时限制只使用栈来传递参数。如果不这么做，用汇编代码调用C函数时可能会产生异常。

注2：用户空间和内核空间各自的地址是不能直接互相访问的，需要借助函数拷贝来实现数据传递，后面会说明。

5. 中断处理函数返回到API中；
6. API将EAX返回给应用程序。

注：当glibc库没有封装某个系统调用时，我们就没办法通过使用封装好的API来调用该系统调用，而使用 `int 0x80` 触发中断又需要进行汇编代码的编写，这两种方法都不适合我们在添加过新系统调用后再编写测试代码去调用。因此我们后面采用的是第三种方法，使用glibc提供的 `syscall` 库函数，这个库函数只需要传入调用号 and 对应内核函数要用到的参数，就可以直接调用我们新写的系统调用。具体详见 3.4.1.

3.3 添加系统调用的流程

Linux 4.9的文档 [linux-4.9.263/Documentation/adding-syscalls.txt](#) 中有说明如何添加一个系统调用。本实验文档假设所有同学使用的平台都是x86，采用的是x86平台的系统调用添加方法。其他平台的同学可以参考上述文档给出的其他平台系统调用添加方法。

在实现系统调用之前，要先考虑好添加系统调用的函数原型，确定函数名称、要传入的参数个数和类型、返回值的意义。在实际设计中，还要考虑加入系统调用的必要性。

我们这里以一个统计系统中进程个数的系统调用 `ps_counter(int *num)` 作为演示，num是返回值对应的地址。

因为系统调用的返回值表示系统调用的完成状态（0是正常，其他值是异常，代表错误编号），所以不建议用返回值传递信息。

本节的各步是没有严格顺序的，但在编译前都要完成：

3.3.1 注册系统调用

内核的汇编代码会在 `arch/x86/include/generated/asm/syscalls_64.h` 中查找调用号。为便于添加系统调用，x86平台提供了一个专门用来注册系统调用的文件 `/arch/x86/entry/syscalls/syscall_64.tbl`。在编译时，脚本 `arch/x86/entry/syscalls/syscalltbl.sh` 会被运行，将上述 `syscall_64.tbl` 文件中登记过的系统调用都生成到前面的 `syscalls_64.h` 文件中。因此我们需要修改 `syscall_64.tbl`。

打开 `linux源码路径/arch/x86/entry/syscalls/syscall_64.tbl`。

```
# linux-4.9.263/arch/x86/entry/syscalls/syscall_64.tbl
# ... 前面还有，但没展示
329 common pkey_mprotect      sys_pkey_mprotect
330 common pkey_alloc          sys_pkey_alloc
331 common pkey_free            sys_pkey_free
# 以上是系统自带的，不是我写的，下面这个是
332 common ps_counter          sys_ps_counter
```

这个文件是x86平台下的系统调用注册表，记录了很多数字与函数名的映射关系。我们在这个文件中注册系统调用。

每一行从左到右各字段的含义，及填写方法是：

- 系统调用号：请为你的系统调用找一个没有用过的数字作为调用号（记住调用号，后面要用）
- 调用类型：本次实验请选择common（含义为：x86_64和x32都适用）
- 系统调用名xxx
- 内核函数名sys_xxx

提示：为保持代码美观，每一列用制表符排得很整齐。但经本人测试，无论你是用制表符还是用空格分隔，无论你用多少空格或制表符，只要你把每一列分开了，无论是否整齐，对这一步的完成都没有影响。

3.3.2 声明内核函数原型

打开 `linux-4.9.263/include/linux/syscalls.h`，里面是对于系统调用函数原型的定义，在最后面加上我们创建的新的系统调用函数原型，格式为 `asmlinkage long sys_xxx(...)`。注意，如果传入了用户空间的地址，需要加入 `__user` 宏来说明。

```
asmlinkage long sys_mlock2(unsigned long start, size_t len, int flags);
asmlinkage long sys_pkey_mprotect(unsigned long start, size_t len,
    unsigned long prot, int pkey);
asmlinkage long sys_pkey_alloc(unsigned long flags, unsigned long init_val);
asmlinkage long sys_pkey_free(int pkey);
//这里是我们新增的系统调用
asmlinkage long sys_ps_counter(int __user * num);

#endif
```

3.3.3 实现内核函数

你可以在 `linux-4.9.264/kernel/sys.c` 代码的最后添加你自己的函数，这个文件中有很多已经实现的系统调用的函数作为参考。我们给出了一段示例代码：

```
SYSCALL_DEFINE1(ps_counter, int __user *, num){
    struct task_struct* task;
    int counter = 0;
    printk("[Syscall] ps_counter\n");
    for_each_process(task){
        counter ++;
    }
    copy_to_user(num, &counter, sizeof(int));
    return 0;
}
```

以下是这段代码的解释：

1. 使用宏 `SYSCALL_DEFINEx` 来简化实现的过程，其中x代表参数的个数。传入宏的参数为：调用名(不带sys_)、以及每个参数的类型、名称（注意这里输入时为两项）。

以我们现在的 `sys_ps_counter` 为例。因为使用了一个参数，所以定义语句是 `SYSCALL_DEFINE1(ps_counter, int *, num)`。如果该系统调用有两个参数，那就应该是 `SYSCALL_DEFINE2(ps_counter, 参数1的类型, 参数1的名称, 参数2的类型, 参数2的名称)`，以此类推。

无法通过 `SYSCALL_DEFINEx` 定义二维数组（如 `char (*p)[50]`）为参数。你可以尝试定义一个，然后结合编译器的报错信息分析该宏的实现原理，并分析不能定义的原因。

2. `printk()` 是内核中实现的print函数，和 `printf()` 的使用方式相同，`printk()` 会将信息写入到内核日志中。尤其地，在qemu下调试内核时，系统日志会直接打印到屏幕上，所以我们可以直接在屏幕上看到 `printk` 打印出的内容。`printf()` 是使用了C的标准库函数的时候才能使用的，而内核中无法使用标准库函数。你不能 `#include<stdio.h>`，自然不能用 `printf()`。

在使用 `printk` 时，行首输出的时间等信息是去不掉的。我们会用此方法检查ps的输出是否是用户态。

3. 为了获取当前的所有进程，我们使用了宏函数 `for_each_process(p)`（定义在 `include/linux/sched.h` 中），遍历当前所有任务的信息结构体 `task_struct`（同样定义在 `include/linux/sched.h` 中），并将地址赋值给参数 `p`。后面实验内容中我们会进一步用到 `task_struct` 中的成员变量来获取更多信息，注意我们暂时不讨论多线程的场景，每一个 `task_struct` 对应的可以认为就是一个进程。
4. 前面说到，系统调用是在内核空间中执行，所以如果我们要将在内核空间获取的值传递给用户空间，需要使用函数 `copy_to_user(void __user *to, const void *from, unsigned long n)` 来完成从内核空间到用户空间的复制。其中，`from` 和 `to` 是复制的来源和目标地址，`n` 是复制的大小，我们这里就是 `sizeof(int)`。

3.4 测试

若想使用VS Code调试Linux内核，可以参考此文章：<https://zhuanlan.zhihu.com/p/105069730>，直接从“VS Code配置”一章阅读即可。实际操作中，你可能还需要在VS Code中安装"C/C++"和"C/C++ Extension Pack"插件。

3.4.1 编写测试代码

在你的Linux环境下（不是打开qemu后弹出的那个）编写测试代码。我们在这里命名其为 `get_ps_num.c`。新增的系统调用可以使用 `long int syscall (long int sysno, ...)` 来触发，`sysno` 就是我们前面添加的调用号，后面跟着要传入的参数，下面是一个简单的测试代码样例：

```
#include<stdio.h>
#include<unistd.h>
#include<sys/syscall.h>
int main(void)
{
    int result;
    syscall(332, &result);
    printf("process number is %d\n",result);
    return 0;
}
```

提示：

- 这里的测试代码是用户态代码，不是内核态代码，所以可以使用 `printf`。
- 为使用系统调用号调用系统调用，需要使用 `sys/syscall.h` 内的 `syscall` 函数。
- 我们推荐使用C语言而不是C++，因为C++在下一步(3.4.2)的静态编译时容易出现問題。使用C语言时需要注意C和C++之间的语法差异。

3.4.2 编译

使用gcc编译器编译 `get_ps_num.c`。本次实验需要使用静态编译（`-static` 选项）。要用到的gcc指令原型是：`gcc [-static] [-o outfile] infile`。

3.4.3 运行测试程序

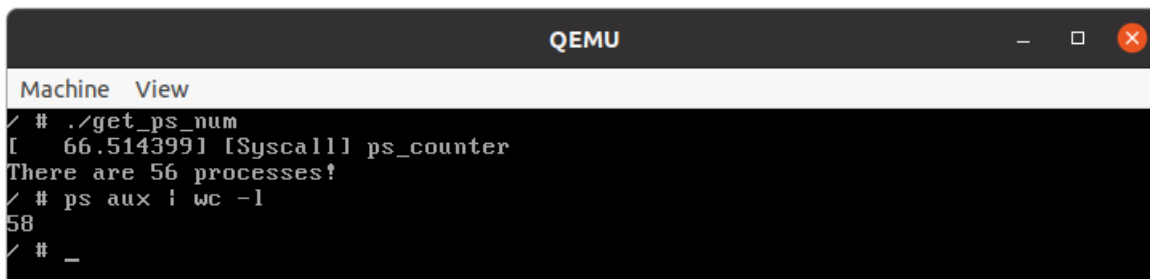
1. 将3.4.2编译出的可执行文件 `get_ps_num` 放到 `busybox-1.32.1/_install` 下面，重新制作initramfs文件（重新执行Lab1的3.2.6.6的find操作），这样我们才能在qemu中看见编译好的 `get_ps_num` 可执行文件。

注意：因为`_install`文件夹只有root用户才有修改权限，所以复制文件应该在命令行下使用 `sudo`。

2. 重新编译linux源码（重新执行Lab1的3.2.4.3，无需重新make menuconfig）
3. 运行qemu（重新执行Lab1的3.2.7，我们在这里不要求使用gdb调试）

- 在qemu中运行 `get_ps_num` 程序，得到当前的进程数量（这个进程数量是包括 `get_ps_num` 程序本身的，下图中，除了 `get_ps_num` 本身外，还有55个进程）。
- 验证：使用 `ps aux | wc -l` 获取当前进程数。但这个输出的结果与上一步中得到的不同。请自己思考数字不同的原因。

在实际操作中，系统中的进程数量可能不是55。只要与 `ps aux | wc -l` 得到的结果匹配即可。



```
QEMU
Machine View
/ # ./get_ps_num
[ 66.514399] [Syscall] ps_counter
There are 56 processes!
/ # ps aux | wc -l
58
/ # _
```

3.5 任务目标

- 在linux4.9下创建适当的（可以是一个或多个）系统调用。利用新实现的系统调用，实现一个linux4.9下的进程状态信息统计程序，参考命令 `top`：

PID	PR	S	%CPU	COMMAND
142420	20	S	6.0	netease-cloud-m
146353	20	S	4.0	chrome
2100	20	S	3.3	pulseaudio
2542	20	S	3.3	deepin-system-m
786	20	S	2.7	Xorg
1742	20	S	1.3	kwin_x11
1822	20	S	1.3	dde-dock

- 输出的信息需要包括：
 - 进程的PID
 - 进程的COMMAND（进程名）
 - 进程是否处于running状态
 - 进程的CPU占用率
- 每一秒刷新一次信息。输出的信息需要按CPU占用率排序。输出前20行即可，无需输出在两次刷新间隔之间新建/消失的进程。程序需要一直运行，无需考虑你写的 `top` 程序如何关闭。
- 添加的系统调用在被调用时需要 `printf` 出自己的系统调用名（如 `ps_counter`）
- 具体评分规则见4.3。我们不限制大家创建系统调用的数量、功能，也不限制测试程序的实现。能完成实验内容即可。
- 下面是示例输出。进程信息是随时间变化的，状态值、排序并不固定。

```
QEMU
Machine View
[ 77.990693] [Syscall] ps_counter
[ 77.990943] [Syscall] ps_info
PID      COMM      CPU  ISRUNNING
122      get_info  0.90%  1
7        rcu_sched  0.14%  0
121      kworker/0:2  0.04%  0
1        sh      0.00%  0
2        kthreadd  0.00%  0
3        ksoftirqd/0  0.00%  0
4        kworker/0:0  0.00%  0
5        kworker/0:0H  0.00%  0
6        kworker/u2:0  0.00%  0
8        rcu_bh    0.00%  0
9        migration/0  0.00%  0
10       lru-add-drain  0.00%  0
11       watchdog/0  0.00%  0
12       cpuhp/0    0.00%  0
13       kdevtmpfs  0.00%  0
14       netns     0.00%  0
15       khungtaskd  0.00%  0
16       oom_reaper  0.00%  0
17       writeback  0.00%  0
18       kcompactd0  0.00%  0
```

3.6 任务提示

- 前面示例中 `for_each_process` 获取到的 `task_struct` 结构体保存有对应进程的很多信息，如：

成员变量	成员含义	其他说明
<code>char</code> <code>comm[TASK_COMM_LEN]</code>	进程名	
<code>pid_t pid</code>	pid	
<code>volatile long</code> <code>state</code>	进程的状态，如运行、停止、僵尸进程等	进程的状态定义于 <code>include/linux/sched.h</code> 的197 ~ 244行。Linux使用位掩码(bitmask)实现多种状态的叠加，若想了解请自行搜索。
<code>struct sched_entity</code> <code>se</code>	该进程的调度实体，里面记录了很多与进程调度相关的信息	<code>sched_entity</code> 结构体中的 <code>sum_exec_runtime</code> 成员定义进程从进程开始时算起的实际运行时间。它的单位是纳秒(ns)。

- 内核空间的成员变量（如 `task_struct` 结构体及其成员）要借助 `copy_to_user(void __user *to, const void *from, unsigned long n)` 函数复制到用户空间的变量中，才可以在用户空间访问内容。
- 如果需要将用户空间的变量复制到内核空间，用 `copy_from_user(void *to, const void __user *from, unsigned long n)`
- 在你的top程序（不是内核代码）中，使用 `<unistd.h>` 中的 `sleep` 函数可以使程序等待，使用 `<stdlib.h>` 中的 `system("clear")` 可以清屏。
- 进程在一段物理时间内的CPU占用率的计算方法：进程在这段时间内的实际运行时间/物理时间长度。
- 如果访问数组越界，会出现很多奇怪的错误（如程序运行时报malloc错误、段错误、输出其他数组的值或一段乱码（注意，烫烫烫是Visual C的特性，Linux下没有烫烫烫）、其他数组的值被改变等）。提问前请先排查是否出现了此类问题。
- 一些常见错误：
 - 拼写错误，如 `asmlkage`、`commom`、`__uesr`、中文逗号
 - 语法错误，如 `sizeof(16 * int)`
 - C标准问题，如C语言声明结构体、算 `sizeof` 时应该用 `struct xxx` 而不是直接 `xxx`

上述错误均可以通过阅读编译器报错信息查出。

第四部分 实验评分标准

本次实验共10分，无实验报告。评分方式为按点给分，某个点做出来即可拿到对应的分数，不同的得分点之间没有约束关系，比如本实验中最难的按cpu使用率排序和多管道、重定向这些都没做，但是其他的基础内容都做出来了，也是可以拿到8分的。

4.1 知识问答 (2')

我们会使用随机数发生器随机地从下面题库中抽三道题，每回答正确一题得一分，满分2分。也就是说，允许答错一题。请按照自己的理解答题，答题时不准念实验文档或念提前准备的答案稿。

下面是题库：

1. 解释 `wc` 和 `grep` 指令的含义。
2. 解释 `ps aux | grep firefox | wc -l` 的含义。
3. `echo aaa | echo bbb | echo ccc` 是否适合做shell实验中管道符的检查用例？说明原因。
4. 对于匿名管道，如果写端不关闭，并且不写，读端会怎样？
5. 对于匿名管道，如果读端关闭，但写端仍尝试写入，写端会怎样？
6. 假如使用匿名管道从父进程向子进程传输数据，这时子进程不写数据，为什么子进程要关闭管道的写端？
7. `fork`之后，是管道从一分为二，变成两根管道了吗？如果不是，复制的是什么？
8. 解释系统调用 `dup2` 的作用。
9. 什么是shell内置指令，为什么不能fork一个子进程然后 `exec cd` ？
10. 为什么 `ps aux | wc -l` 得出的结果比 `get_ps_num` 多2？
11. 进程名的最大长度是多少？这个长度在哪定义？
12. `task_struct` 在Linux源码的哪个文件中定义？
13. 为什么无法通过 `SYSCALL_DEFINE` 定义二维数组（如 `char (*p)[50]` ）为参数？
14. 在修改内核代码的时候，能用 `printf` 调试吗？如果不能，应该用什么调试？
15. `read()`、`write()`、`dup2()` 都能直接调用。现在我们已经写好了一个名为`ps_counter`的系统调用。为什么我们不能在测试代码中直接调 `ps_counter()` 来调用系统调用？

大家可以自己思考或互相讨论这些题目的答案（但不得在本课程群内公布答案）。助教不会公布这些题目的答案。

4.2 “实现一个Shell”部分检查标准 (4')

- 支持基本的单条命令运行、支持两条命令间的管道 `|`、内建命令（只要求 `cd` / `exit` ），得1分。
- 选做：
 - 支持多条命令间的管道 `|` 操作，得2分。
 - 支持重定向符 `>`、`>>`、`<` 和分号 `;`，得2分。
- 若两个选做都做，另加1分。
- 你需要流畅地说明你是如何实现实验要求的（即，现场讲解代码）。本部分共1分。若未能完成任何一个功能，则该部分分数无效。

注：实验检查时，请**自行准备**可以验证相应功能支持的测试样例。若我们认为测试样例不能验证程序是否有误，则该部分不计分。

4.3 “编写系统调用”部分检查标准 (4')

- 你的程序需要在【用户态】每秒打印一次各进程的CPU占用统计，并按CPU占用统计排序。本部分共1分。
- 你的程序需要在【用户态】每秒打印一次各进程的PID、进程名、进程是否处于running状态。本部分共1分。
- 你需要在现场阅读Linux源码，展示pid的数据类型是什么（即，透过多层 `typedef` 找到真正的pid数据类型）。本部分共1分。
- 你需要流畅地展示你修改了哪些文件的代码，允许对着实验文档描述（即，现场讲解代码）。本部分共1分。若未能完成任何一个功能，则该部分分数无效。
- 不准在内核态直接 `printk` 进程信息。否则会被视为该部分未完成。
- 不准在用户态直接调用ps/top等工具，或读取procfs获取进程信息，否则会被视为该部分未完成。