

# 华东师范大学软件学院实验报告

课程名称：操作系统	年级：2023 级本科	上机实践成绩：
指导教师：张民	姓名：张梓卫	
上机实践名称：Pintos Userprog Part 1	学号：10235101526	上机实践日期：2024/12/16
上机实践编号：（5）	组号：	上机实践时间：2 学时

目录

一 实验目的	1	4 处理内存非法访问 . . . . .	10
二 内容与设计思想	1	5 完善系统调用 . . . . .	12
三 使用环境	2	5.1 halt . . . . .	13
四 实验过程与分析	2	6 syscall_exit . . . . .	13
1 进入实验背景 . . . . .	2	7 syscall_exec . . . . .	14
2 开始分析代码 . . . . .	3	8 syscall_wait . . . . .	14
2.1 分析 start_process() . . . . .	5	8.1 syscall_write . . . . .	15
2.2 分析 load() 函数 . . . . .	5		
3 实现系统调用 . . . . .	9	五 实验结果	16
		六 附录	17

## 一 实验目的

本实验的目标是完成参数传递和部分系统调用（exit 和 write），使得 make check 通过 args 相关的 5 个测试。

本项目的主要任务：实现用户程序和 OS 之间的系统调用。本次实验作出修改的代码同时上传到了 Github 之上，仓库地址为：<https://github.com/Shichien/ECNU-23-SEI-Homework>

请在上传的 PDF 文件中直接点击粉色链接即可。

## 二 内容与设计思想

- 参数传递，系统调用

Notice

- 编译生成 build 文件夹
- 创建用户磁盘
- 修改 process\_execute() 实现文件名与参数分离
- 修改 load() 实现文件名与参数分离并将可执行文件存到 ESP 中
- 修改 setup\_stack() 实现参数堆栈
- 注：上述方法均位于 userprog/process.c 中

该开始着手研究允许运行用户程序的系统部分了。基本代码已经支持加载和运行用户程序，但是无法进行 I / O 或交互。在此项目中，将使程序能够通过系统调用与 OS 进行交互，并为用户进程提供系统调用。在 Project 1 中，执行的操作都是在内核模式下运行的，而在这个 Project 中要求我们对非内核进行修改。

- src/userprog
- process.c/ process.h: 进程代码（需修改）
- pagedir.c/pagedir.h: 物理内存管理（不需修改，可以调用实现的）
- syscall.c/syscall.h: 系统调用，目前只搭了骨架（需修改）
- exception.c/exception.h: 异常处理情况（需修改）
- gdt.c/gdt.h: 全局描述表（不需修改）
- tss.c/tss.h: 任务状态段（不需修改）

## 三 使用环境

使用 Docker v27.1.1 进行 Pintos 的安装实验，基于 Windows 11 操作系统使用 WSL2。  
实验报告使用 L<sup>A</sup>T<sub>E</sub>X 进行撰写，使用 VSCode + Vim 编辑器进行文本编辑。

## 四 实验过程与分析

### 1 进入实验背景

进入 `cd src/pintos/userprog` 中使用 `make` 命令编译  
然后使用 `cd build` 进入到 build 文件夹中建立一个新的用户磁盘。

```
root@e60a96920998:~/pintos/src/userprog/build# pintos-mkdisk filesys.dsk --fileysys-size=2
root@e60a96920998:~/pintos/src/userprog/build# pintos -f -q
Unknown option: f
Unknown option: q
root@e60a96920998:~/pintos/src/userprog/build# pintos -- -f -q
qemu-system-i386 -device isa-debug-exit -drive format=raw,media=disk,index=0,file=/tmp/YWZgP
media=disk,index=1,file=filesys.dsk -m 4 -net none -nographic -monitor null
Pintos hda1
Loading.....
Kernel command line: -f -q
Pintos booting with 3,968 kB RAM...
367 pages available in kernel pool.
367 pages available in user pool.
Calibrating timer... 157,081,600 loops/s.
hda: 1,008 sectors (504 kB), model "QM00001", serial "QEMU HARDDISK"
hda1: 189 sectors (94 kB), Pintos OS kernel (20)
hdb: 5,040 sectors (2 MB), model "QM00002", serial "QEMU HARDDISK"
hdb1: 4,096 sectors (2 MB), Pintos file system (21)
fileysys: using hdb1
Formatting file system...done.
Boot complete.
Timer: 52 ticks
Thread: 0 idle ticks, 53 kernel ticks, 0 user ticks
hdb1 (fileysys): 3 reads, 6 writes
Console: 581 characters output
```

图 1: 编译

注意到这里直接使用 `pintos -f -q` 会提示未知的参数，要修改为: `pintos -- -f -q`，才能编译成功。  
接下来，执行下一条命令，将 `echo.c` 复制到磁盘中，且命名为 `echo`。

```

root@e60a96920998:~/pintos/src/userprog/build# pintos -p ../../examples/echo.c -a echo -- -q
Copying ../../examples/echo.c to scratch partition...
qemu-system-i386 -device isa-debug-exit -drive format=raw,media=disk,index=0,file=/tmp/8CQHn1edQ.dsk -drive format=raw,
media=disk,index=1,file=filesys.dsk -m 4 -net none -nographic -monitor null
Pintos hda1
Loading.....
Kernel command line: -q extract
Pintos booting with 3,968 kB RAM...
367 pages available in kernel pool.
367 pages available in user pool.
Calibrating timer... 157,081,600 loops/s.
hda: 1,008 sectors (504 kB), model "QM00001", serial "QEMU HARDDISK"
hda1: 189 sectors (94 kB), Pintos OS kernel (20)
hda2: 4 sectors (2 kB), Pintos scratch (22)
hdb: 5,040 sectors (2 MB), model "QM00002", serial "QEMU HARDDISK"
hdb1: 4,096 sectors (2 MB), Pintos file system (21)
filesys: using hdb1
scratch: using hda2
Boot complete.

```

图 2: 复制文件

在未实现参数传递时，我们直接执行 `pintos -- -q run 'echo x'` 这样的命令其实是很不合理的。这个输入的含义是：执行 `echo` 文件，传递的参数是 `x`。实际上，`echo x` 被当成了一个整体来处理。

```

root@e60a96920998:~/pintos/src/userprog/build# pintos -- -q run 'echo x'
qemu-system-i386 -device isa-debug-exit -drive format=raw,media=disk,index=0,
media=disk,index=1,file=filesys.dsk -m 4 -net none -nographic -monitor null
Pintos hda1
Loading.....
Kernel command line: -q run 'echo x'
Pintos booting with 3,968 kB RAM...
367 pages available in kernel pool.
367 pages available in user pool.
Calibrating timer... 209,510,400 loops/s.
hda: 1,008 sectors (504 kB), model "QM00001", serial "QEMU HARDDISK"
hda1: 189 sectors (94 kB), Pintos OS kernel (20)
hdb: 5,040 sectors (2 MB), model "QM00002", serial "QEMU HARDDISK"
hdb1: 4,096 sectors (2 MB), Pintos file system (21)
filesys: using hdb1
Boot complete.
Executing 'echo x':
Execution of 'echo x' complete.
Timer: 51 ticks
Thread: 0 idle ticks, 51 kernel ticks, 0 user ticks
hdb1 (filesys): 2 reads, 0 writes
Console: 612 characters output
Keyboard: 0 keys pressed
Exception: 0 page faults
Powering off...

```

图 3: Echo

## 2 开始分析代码

先查看官方文档中对目前清空的介绍：

### 3.3.3 Argument Passing

Currently, `process_execute()` does not support passing arguments to new processes. Implement this functionality, by extending `process_execute()` so that instead of simply taking a program file name as its argument, it divides it into words at spaces. The first word is the program name, the second word is the first argument, and so on. That is, `process_execute("grep foo bar")` should run `grep` passing two arguments `foo` and `bar`.

Within a command line, multiple spaces are equivalent to a single space, so that `process_execute("grep foo bar")` is equivalent to our original example. You can

图 4: Pintos References

翻译如下：

目前，`process_execute()` 不支持向新进程传递参数。通过扩展 `process_execute()` 来实现此功能，而不是简单的 `task`。将程序 `filename` 作为其参数，它在空格处将其拆分为单词。第一个词是程序名，第二个单词是第一个参数，以此类推。也就是 `process execute ("grep-foo-bar")` 应该运行 `grep`，传递两个参数 `foo` 和 `bar`。在命令行中多个空格相当于一个空格，以便进程执行 ("`grep-foo-bar`") 与我们最初的示例等效。

查看相关的代码处，我们可以看到以下内容：

```

1  tid_t process_execute (const char *file_name) {
2      char *fn_copy;
3      tid_t tid;
4
5      /* Make a copy of FILE_NAME.
6       * Otherwise there's a race between the caller and load(). */
7      fn_copy = palloc_get_page (0);
8      if (fn_copy == NULL)
9          return TID_ERROR;
10     strcpy (fn_copy, file_name, PGSIZE);
11
12     /* Create a new thread to execute FILE_NAME. */
13     tid = thread_create (file_name, PRI_DEFAULT, start_process, fn_copy);
14     if (tid == TID_ERROR)
15         palloc_free_page (fn_copy);
16     return tid;
17 }

```

process

### Notice

在传递参数时，传递 `process_execute` 函数的参数 `file_name` 既包括可执行文件的名称，也包含了可执行文件的参数。所以，要做的第一件事就是换把可执行文件名称和参数互相分开。我们可以使用 `string.h` 中的 `strtok_r()` 来分离参数。

```

1  tid_t process_execute (const char *file_name) {
2      char *fn_copy;
3      char *thread_name;
4      char *save_ptr; // 用来保存strtok_r的状态
5
6      thread_name = malloc(strlen(file_name) + 1);
7      strcpy(thread_name, file_name, strlen(file_name) + 1);
8      // 拷贝 file_name 到动态分配的 thread_name，以便后续通过 strtok_r 截取线程名称
9      thread_name = strtok_r(thread_name, " ", &save_ptr); // 创建一个线程用来运行"文件名"
10     tid_t tid;
11
12     /* Make a copy of FILE_NAME.
13      * Otherwise there's a race between the caller and load(). */
14     fn_copy = palloc_get_page (0);
15     if (fn_copy == NULL)
16         return TID_ERROR;
17     strcpy (fn_copy, file_name, PGSIZE);
18
19     /* Create a new thread to execute FILE_NAME. */
20     tid = thread_create (file_name, PRI_DEFAULT, start_process, fn_copy);
21     free(thread_name); // 释放由 malloc 创建的线程名
22     if (tid == TID_ERROR)
23         palloc_free_page (fn_copy);
24     return tid;
25 }

```

process

### Notice

这里我们首先在最开始添加两个 `char *` 类型变量，`*save_ptr` 指向 `file_name` 的一个拷贝，是为了实现保存从命令语句中分离出来的真正的可执行文件名，给 `fn` 分配与命令语句容量相等的内存空间。

## 2.1 分析 `start__process()`

我们在函数中将 `file_name` 通过 `load` 存储到栈中。

下面这个是 Pintos 中的原始代码。

```

1 static void
2 start__process (void *file_name_)
3 {
4     char *file_name = file_name_;
5     struct intr_frame if_;
6     bool success;
7
8     /* Initialize interrupt frame and load executable. */
9     memset (&if_, 0, sizeof if_);
10    if_.gs = if_.fs = if_.es = if_.ds = if_.ss = SEL_UDSEG;
11    if_.cs = SEL_UCSEG;
12    if_.eflags = FLAG_IF | FLAG_MBS;
13    success = load (file_name, &if_.eip, &if_.esp);
14
15    /* If load failed, quit. */
16    palloc_free_page (file_name);
17    if (!success)
18        thread_exit ();
19
20    /* Start the user process by simulating a return from an
21       interrupt, implemented by intr_exit (in
22       threads/intr-stubs.S). Because intr_exit takes all of its
23       arguments on the stack in the form of a 'struct intr_frame',
24       we just point the stack pointer (%esp) to our stack frame
25       and jump to it. */
26    asm volatile ("movl %0, %%esp; jmp intr_exit" : : "g" (&if_) : "memory");
27    NOT_REACHED ();
28 }

```

start\_\_process

修改后的代码如下所示：

```

1 char *token, *save_ptr;
2 file_name = strtok_r(file_name, " ", &save_ptr); // 获取文件名
3 success = load (file_name, &if_.eip, &if_.esp);

```

process

## 2.2 分析 `load()` 函数

主要是将从 `FILE_NAME` 获取的可执行的程序放入到当前线程中。

针对传参，我们需要第二次参数分离：传递过来的是 `file_name="args-single onearg"`。

```

1 /* Separate file name and arguments. */
2 char *exec_name;
3 char *save_ptr;
4
5 /* Make a modifiable copy of file_name. */
6 char modifiable_file_name[PGSIZE];
7 strcpy(modifiable_file_name, file_name, PGSIZE);
8
9 /* Extract the executable name (first token). */

```

```

10 exec_name = strtok_r(modifiable_file_name, " ", &save_ptr);
11
12 /* Open the executable file. */
13 file = filesys_open(exec_name);
14 if (file == NULL)
15 {
16     printf ("load: %s: open failed\n", exec_name);
17     goto done;
18 }

```

load

另外，在 load() 函数中，注释中写明了 `setup_stack()` 函数的作用是设置参数堆栈。

```

/* Set up stack. */
if (!setup_stack (esp))
    goto done;

```

图 5: load()

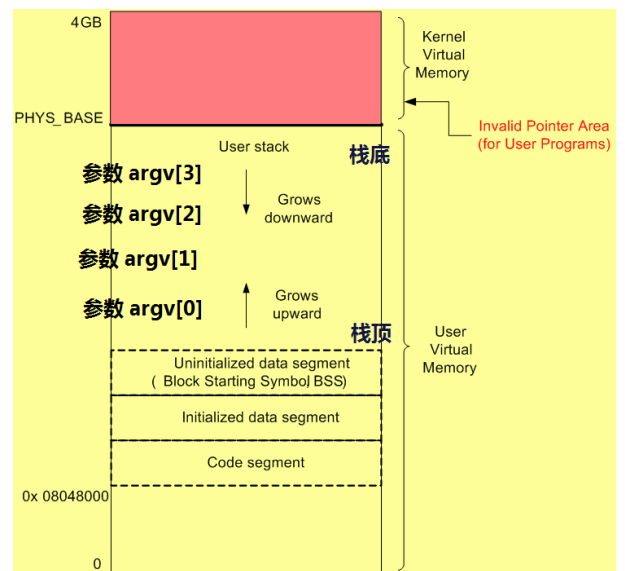
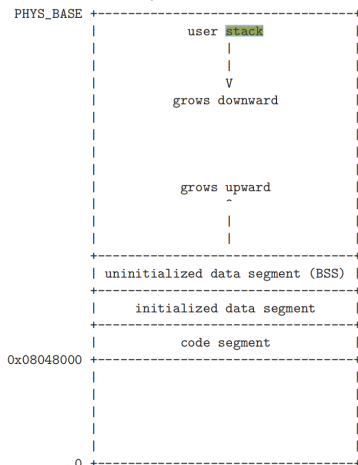
而在 `setup_stack()` 函数中，它仅仅实现了页的初始化，并没有参数传递实现的代码。

因此我们需要将用 `strok` 分离出的 `filename` 以及 `argv` 传递进去。

一个实现思路是，在 `setup_stack()` 函数中添加参数 `file_name` 和 `save_ptr`

#### 3.1.4.1 Typical Memory Layout

Conceptually, each process is free to lay out its own user virtual memory however it chooses. In practice, user virtual memory is laid out like this:



#### Notice

`if _esp` 原来指向 `esp` 的内核和用户区边界为避免内存访问越界（访问到内核区）  
 所以需将 `setup_stack()` 函数中的 `esp = PHYS_BASE` 修改成 `*esp = PHYS_BASE - 12;`  
 后来，文档中还附加了一条对于 `bin/ls -l foo bar` 命令执行时栈的增长情况：

The table below shows the state of the stack and the relevant registers right before the beginning of the user program, assuming `PHYS_BASE` is `0xc0000000`:

Address	Name	Data	Type
0xbfffffff	<code>argv[3][...]</code>	<code>'bar\0'</code>	<code>char[4]</code>
0xbffffff8	<code>argv[2][...]</code>	<code>'foo\0'</code>	<code>char[4]</code>
0xbffffff5	<code>argv[1][...]</code>	<code>'-l\0'</code>	<code>char[3]</code>
0xbffffffd	<code>argv[0][...]</code>	<code>'/bin/ls\0'</code>	<code>char[8]</code>
0xbffffec	<code>word-align</code>	0	<code>uint8_t</code>
0xbffffe8	<code>argv[4]</code>	0	<code>char *</code>
0xbffffe4	<code>argv[3]</code>	0xbfffffff	<code>char *</code>
0xbffffe0	<code>argv[2]</code>	0xbffffff8	<code>char *</code>
0xbffffdc	<code>argv[1]</code>	0xbffffff5	<code>char *</code>
0xbffffd8	<code>argv[0]</code>	0xbffffffd	<code>char *</code>
0xbffffd4	<code>argv</code>	0xbffffd8	<code>char **</code>
0xbffffd0	<code>argc</code>	4	<code>int</code>
0xbffffcc	<code>return address</code>	0	<code>void (*)()</code>

In this example, the stack pointer would be initialized to `0xbffffcc`.

As shown above, your code should start the stack at the very top of the user virtual address space, in the page just below virtual address `PHYS_BASE` (defined in `threads/vaddr.h`).

图 7: 示例栈增长

根据 PPT 上的提示, 我们编写以下代码来实现初始化栈:

另外, 官方文档中还提到了我们可以使用 `hex_dump()` 函数来打印此时栈的信息, 不妨添加进去。

As shown above, your code should start the stack at the very top of the user virtual address space, in the page just below virtual address `PHYS_BASE` (defined in `threads/vaddr.h`).

如上所示, 你的代码应该在虚拟地址 `PHYS_BASE` (定义在 `threads/vaddr.h`) 下面的页面启动堆栈。

You may find the non-standard `hex_dump()` function, declared in `<stdio.h>`, useful for debugging your argument passing code. Here's what it would show in the above example:

你可能会发现非标准的 `hex_dump()` 函数, 声明在 `<stdio.h>` 中, 它对调试传递参数的代码很有用。这是它在上面的例子中显示的内容:

```
bffffffc 00 00 00 00 | .....|
bffffffd 04 00 00 00 d8 ff ff bf-ed ff ff bf f5 ff ff bf |.....|
bffffffe f8 ff ff bf fc ff ff bf-00 00 00 00 2f 62 69 |...../bi|
bfffffff 6e 2f 6c 73 00 2d 6c 00-66 6f 6f 00 62 61 72 00 |n/ls.-l.foo.bar.|
```

图 8: `hex_dump()`

在函数中, 写好了为什么要用这样的方式实现的具体注释:

```
1 static bool
2 setup_stack (void **esp, char* file_name)
3 {
4     uint8_t *kpage;
5     bool success = false;
6
7     // 分配一页用户内存并将其安装到用户栈的顶部
8     kpage = pallocc_get_page (PAL_USER | PAL_ZERO);
9     if (kpage != NULL)
10     {
11         success = install_page (((uint8_t *) PHYS_BASE) - PGSIZE, kpage, true);
12         if (success)
13             *esp = PHYS_BASE - 12; // 初始化栈指针到栈顶
14         else
15             pallocc_free_page (kpage);
16     }
17
18     // 创建一个临时拷贝的文件名字符串
19     char *token, *temp_ptr;
20     char *filename_cp = malloc(strlen(file_name) + 1);
21     strcpy (filename_cp, file_name, strlen(file_name) + 1);
22
23     // 计算参数数量 (argc)
24     enum intr_level old_level = intr_disable();
25     int argc = 1; // 初始值为 1, 因为至少有一个参数
```

```

26 bool is_lastone_space = false; // 记录最后一个字符是否为空格
27 for (int j = 0; j != strlen(file_name); j++) {
28     if (file_name[j] == ' ') {
29         if (!is_lastone_space)
30             argc++; // 如果遇到新的空格且不是连续空格, 参数计数加一
31         is_lastone_space = true;
32     } else {
33         is_lastone_space = false;
34     }
35 }
36 intr_set_level (old_level);
37
38 // 为参数指针 (argv) 分配内存
39 int *argv = calloc(argc, sizeof(int));
40
41 // 将文件名分割为各个参数, 并依次压入栈中
42 int i;
43 token = strtok_r (file_name, " ", &temp_ptr);
44 for (i = 0; ; i++) {
45     if (token) {
46         *esp -= strlen(token) + 1; // 为参数字符串分配空间
47         memcpy(*esp, token, strlen(token) + 1); // 将参数字符串复制到栈
48         argv[i] = *esp; // 保存参数地址
49         token = strtok_r (NULL, " ", &temp_ptr); // 获取下一个参数
50     } else {
51         break; // 没有更多参数时退出循环
52     }
53 }
54
55 // 对齐到 4 字节边界 (Word Align)
56 *esp -= ((unsigned)*esp % 4);
57
58 // 空指针标志 (null pointer sentinel), 表示 argv[argc] 为空
59 *esp -= sizeof(int);
60
61 // 将参数地址依次压入栈中
62 for (i = argc - 1; i >= 0; i--)
63 {
64     *esp -= sizeof(int);
65     memcpy(*esp, &argv[i], sizeof(int));
66 }
67
68 // 压入 argv 的地址
69 int tmp = *esp;
70 *esp -= sizeof(int);
71 memcpy(*esp, &tmp, sizeof(int));
72
73 // 压入 argc (参数数量)
74 *esp -= sizeof(int);
75 memcpy(*esp, &argc, sizeof(int));
76
77 // 压入返回地址 (return address)
78 *esp -= sizeof(int);
79 memcpy(*esp, &argv[argc], sizeof(int));
80
81 // 打印栈的状态, 便于调试
82 printf("STACK SET. ESP: %p\n", esp);
83 hex_dump((uintptr_t)esp, esp, 100, true); // 打印栈中前 100 字节的数据, 便于调试
84
85 // 释放临时分配的内存
86 free(filename_cp);
87 free(argv);
88
89 return success;

```



```

}

```

```

setup_stack

```

### Notice

在官方文档中提到，栈的增长是反方向的：栈地址是向下增长的，而字符串的地址是向上增长的，所以在放置参数的时候需要把参数倒序放置。

直到这里，即使我们使用 **make check** 命令，仍然无法看到输出，只会出现：

Run didn't produce any output”。

因为用户程序的打印（输出到 stdout）也是通过写的系统调用实现的。

- The `write` system call for writing to fd 1, the system console. All of our test programs write to the console (the user process version of `printf()` is implemented this way), so they will all malfunction until `write` is available.
- `write` 系统调用，用于写入 fd1，即系统控制台。我们所有的测试程序都是写控制台的（用户进程版本的 `printf()` 是这样实现的），所以它们都会出现故障，直到 `write` 可用。

图 9: write()

## 3 实现系统调用

除了我们要实现的 `userprog/syscall.c`，在 `lib/user` 文件夹下，还有一对 `syscall.h` 和 `syscall.c` 文件

### Notice

这组文件，是提供给用户的系统调用接口，其通过 `int $0x30` 中断指令以唤起一个系统调用，最终会落入 `userprog/syscall.c` 的 `syscall_handler()` 函数中。

其 `enum` 中定义了不同的系统调用的序号：

```

1  enum {
2      /* Projects 2 and later. */
3      SYS_HALT,           /**< Halt the operating system. */
4      SYS_EXIT,           /**< Terminate this process. */
5      SYS_EXEC,           /**< Start another process. */
6      SYS_WAIT,           /**< Wait for a child process to die. */
7      SYS_CREATE,         /**< Create a file. */
8      SYS_REMOVE,         /**< Delete a file. */
9      SYS_OPEN,           /**< Open a file. */
10     SYS_FILESIZE,        /**< Obtain a file's size. */
11     SYS_READ,            /**< Read from a file. */
12     SYS_WRITE,           /**< Write to a file. */
13     SYS_SEEK,            /**< Change position in a file. */
14     SYS_TELL,            /**< Report current position in a file. */
15     SYS_CLOSE,           /**< Close a file. */
16
17     /* Project 3 and optionally project 4. */
18     SYS_MMAP,             /**< Map a file into memory. */
19     SYS_MUNMAP,           /**< Remove a memory mapping. */
20
21     /* Project 4 only. */
22     SYS_CHDIR,            /**< Change the current directory. */
23     SYS_MKDIR,            /**< Create a directory. */
24     SYS_READDIR,          /**< Reads a directory entry. */
25     SYS_ISDIR,            /**< Tests if a fd represents a directory. */
26     SYS_INUMBER,          /**< Returns the inode number for a fd. */
27 };

```

```

syscall.h

```

故我们需要在 `syscall_handler()` 中实现分发：根据类型分发给各个具体处理函数，逐个实现。返回值要放在 EAX 中传回调用者。

### Notice

在 Pintos 中，`struct intr_frame` 是用来保存中断处理器状态的结构体。在系统调用返回时，通常会将结果值存储到 `f->eax` 中，因为这是 x86 架构的约定——返回值通过 `eax` 寄存器传递。

对于系统调用 `SYS_EXIT`，需要返回一个退出码到用户进程，返回值就需要存放到 `f->eax` 中。类似地，`SYS_WRITE` 也会将写入的字节数返回到 `f->eax`。

```

1 static void syscall_handler (struct intr_frame *f UNUSED) {
2     int *p = f->esp;
3     is_valid_addr(p);
4
5     int system_call = *p;
6     switch (system_call) {
7         case SYS_HALT: syscall_halt(); break;
8         case SYS_EXIT: syscall_exit(f); break; // 直接关机，不需要返回值，故是 void function
9         case SYS_EXEC: f->eax = syscall_exec(f); break;
10        case SYS_WAIT: f->eax = syscall_wait(f); break;
11        case SYS_CREATE: f->eax = syscall_creat(f); break;
12        case SYS_REMOVE: f->eax = syscall_remove(f); break;
13        case SYS_OPEN: f->eax = syscall_open(f); break;
14        case SYS_FILESIZE: f->eax = syscall_filesize(f); break;
15        case SYS_READ: f->eax = syscall_read(f); break;
16        case SYS_WRITE: f->eax = syscall_write(f); break;
17        case SYS_SEEK: syscall_seek(f); break; // 调整文件指针的位置，也不需要返回值
18        case SYS_TELL: f->eax = syscall_tell(f); break;
19        case SYS_CLOSE: syscall_close(f); break;
20
21        default:
22            printf("Default %d\n", *p);
23    }
24 }
```

syscall.c

在上方的函数中，有一些函数是不需要返回值的，这些 `void function` 在定义的时候要注意一下。

其次，Pintos 中已经使用了 `syscall_init` 将中断处理函数注册到寄存器：

```

1 void syscall_init (void) {
2     intr_register_int (0x30, 3, INTR_ON, syscall_handler, "syscall");
3 }
```

init.c

另外要注意的是，Pintos 中不允许在系统调用中使用 `malloc()` 函数：

Pintos可以运行普通的C程序，只要它们能装入内存并只使用你实现的系统调用。

- Notably, `malloc()` cannot be implemented because none of the system calls required for this project allow for memory allocation.  
值得注意的是，`malloc()` 不能实现，因为该项目所需的系统调用都不允许内存分配。
- Pintos also can't run programs that use floating point operations, since the kernel doesn't save and restore the processor's floating-point unit when switching threads.  
Pintos也不能运行使用浮点操作的程序，因为内核在切换线程时不会保存和恢复处理器的浮点单元。

图 10: `malloc()`

## 4 处理内存非法访问

PPT 中的提示了使用 `syscall_create` 来实现：

## 以syscall\_create为例

```
int
syscall_create(struct intr_frame *f)
{
    int ret;
    off_t initial_size;
    char *name;

    pop_stack(f->esp, &initial_size, 5);
    pop_stack(f->esp, &name, 4);
    if (!is_valid_addr(name))
        ret = -1;

    lock_acquire(&fileys_lock);
    ret = fileys_create(name, initial_size);
    lock_release(&fileys_lock);
    return ret;
}
```

栈中有3个参数时：  
后两个参数的地址偏移量加4、5  
栈中有2个参数时：（比如Open）  
后一个参数的地址偏移量加1  
栈中有4个参数时（比如Read/Write）  
后三个参数的地址偏移量加5、6、7

```
void *
is_valid_addr(const void *vaddr)
{
    void *page_ptr = NULL;
    if (!is_user_vaddr(vaddr))
        || !(page_ptr = pagedir_get_page(thread_current()->pagedir, vaddr))
    {
        exit_process(-1);
        return 0;
    }
    return page_ptr;
}
```

```
void pop_stack(int *esp, int *a, int offset){
    int *tmp_esp = esp;
    *a = *((int *)is_valid_addr(tmp_esp + offset));
}
```

1. 出栈，获取文件名和文件大小
2. 判断文件名是否合法
3. 对文件上锁
4. 调用fileys.c中的API
5. 释放锁

图 11: PPT

为了保证内核不受到破坏，需要立即终止发出非法指针请求的进程。查看 [userprog/pagedir.c](#) 和 [threads/vaddr.h](#)，可以发现已经为我们定义了以下两个函数：

```
1  /* Looks up the physical address that corresponds to user virtual address UADDR in PD.
2  Returns the kernel virtual address
3  corresponding to that physical address, or a null pointer if UADDR is unmapped. */
4  void *pagedir_get_page (uint32_t *pd, const void *uaddr) {
5      uint32_t *pte;
6      ASSERT (is_user_vaddr (uaddr));
7      pte = lookup_page (pd, uaddr, false);
8      if (pte != NULL && (*pte & PTE_P) != 0)
9          return pte_get_page (*pte) + pg_ofs (uaddr);
10     else
11         return NULL;
12 }
```

pagedir\_get\_page()

```
1  /* Returns true if VADDR is a user virtual address. */
2  static inline bool is_user_vaddr (const void *vaddr) {
3      return vaddr < PHYS_BASE;
4  }
```

is\_user\_vaddr()

除了空指针以外，这就是最常见的两个异常情况。当出现异常的时候都会返回异常信息。可以借助这两个函数来实现内存访问的安全控制。

我们可以跟着这个思路来编写代码：

```
1  // 检查传入的虚拟地址是否是用户空间的有效地址：
2  // 使用 is_user_vaddr 判断地址是否在用户地址范围内。如果是有效的用户地址，则使用 pagedir_get_page 从当前线程
3  // 的页目录中获取该虚拟地址对应的物理页指针。
4  // 如果地址无效或无法找到对应的物理页，则调用 exit_process(-1) 退出当前进程。
5
6  void *is_valid_addr(const void *vaddr) {
7      void *page_ptr = NULL;
8      if (!is_user_vaddr(vaddr) || !(page_ptr = pagedir_get_page(thread_current()->pagedir, vaddr))) {
9          // 地址无效，退出进程
10         exit_process(-1);
11     }
```

```

10         return 0;
11     }
12     // 地址有效, 返回物理页指针
13     return page_ptr;
14 }

```

syscall\_create()

如果地址无效则以异常状态中止进程, 有效则返回物理地址。

再编写 pop\_stack() 函数, 用于从栈中取得元素, 方便取得参数:

```

1 void pop_stack(int *esp, int *a, int offset){
2     int *tmp_esp = esp;
3     *a = *((int *)is_valid_addr(tmp_esp + offset));
4 }

```

pop\_stack()

这里是因为 Pintos 使用的是 4 字节对齐, 所以 offset 参数可以方便地使用整数来取得参数。

Consider a function `f()` that takes three `int` arguments. This diagram shows a sample stack frame as seen by the callee at the beginning of step 3 above, supposing that `f()` is invoked as `f(1, 2, 3)`. The initial stack address is arbitrary:

考虑一个函数 `f()`, 它接受三个 `int` 参数。这张图显示了上面步骤3开始时被调用者看到的一个样本栈帧, 假设 `f()` 被调用为 `f(1, 2, 3)`。栈的初始地址是任意的:

	0xbffffe7c		3	
	0xbffffe78		2	
	0xbffffe74		1	
stack pointer -->	0xbffffe70		return address	

图 12: 栈指针

## 5 完善系统调用

按照 PPT 要求, 在 `thread.h` 中新增结构体变量:

```

struct thread
{
    /* Owned by thread.c. */
    tid_t tid; /*< Thread identifier. */
    int exit_code; /*< Exit code of the thread. */
    enum thread_status status; /*< Thread state. */
    char name[16]; /*< Name (for debugging purposes). */
    uint8_t *stack; /*< Saved stack pointer. */
    int priority; /*< Priority. */
    struct list_elem allelem; /*< List element for all threads list. */

    /* Shared between thread.c and synch.c. */
    struct list_elem elem; /*< List element. */

    bool load_success; /*< 标志线程加载是否成功(例如加载可执行文件时使用)。 */
    struct semaphore load_sema; /*< 信号量, 用于同步线程加载进程, 例如父线程等待子线程加载完成。 */
    int exit_status; /*< 线程的退出状态码, 表示线程终止时的状态。 */
    struct list children_list; /*< 子线程列表, 存储当前线程的所有子线程, 用于管理父子线程关系。 */
    struct thread *parent; /*< 指向当前线程的父线程, 用于实现父子线程的交互(如退出时通知父线程)。 */
    struct file *self; /*< 当前线程正在执行的可执行文件, 防止文件在执行期间被修改或删除。 */
    struct list opened_files; /*< 打开文件的列表, 存储当前线程打开的所有文件, 用于管理文件资源的生命周期。 */

    int fd_count; /*< 文件描述符计数器, 分配唯一的文件描述符给线程打开的文件。 */
    struct child_process *waiting_child; /*< 指向线程正在等待的子线程, 用于实现父线程等待特定子线程结束的功能。 */
}

```

图 13: thread.h

同时定义全局锁: `struct lock filesystem_lock`, 用于文件系统操作的同步。该结构体 `lock` 位于 `synch.h` 中。

当然，不能忘记了子进程的结构体定义：

```

1 struct child_process {
2     int tid;                /**< 子进程的线程 ID。 */
3     bool if_waited;         /**< 子进程是否已被父进程等待。 */
4     int exit_status;        /**< 子进程的退出状态。 */
5     struct list_elem child_elem; /**< 子进程的元素。 */
6     struct semaphore wait_sema; /**< 等待子进程退出的信号量。 */
7 };

```

child\_process

## 5.1 halt

这里只需要简单地调用 `devices/shutdown.c` 下的函数接口即可：

```

1 void syscall_halt(void){
2     shutdown_power_off();
3 }

```

## 6 syscall\_exit

```

1 static void syscall_exit(struct intr_frame *f) {
2     // exit_code在系统调用之后被解析为参数
3     int exit_code = *(int *)check_read_user_ptr(f->esp + ptr_size, sizeof(int));
4     thread_current()->exit_status = exit_code;
5     thread_exit();
6 }

```

### Notice

`syscall_exit` 函数实现了系统调用，用于终止当前进程。通过遍历查找对应的子进程描述结构（`child_process`），更新其退出状态码（`exit_status`）以及是否已被父线程等待的标志（`if_waited`）。同时，将当前线程的退出状态存储在其自身的 `exit_status` 字段中。为了避免多线程环境下的竞争问题，函数通过禁用中断（`intr_disable`）确保状态更新的原子性。

最后调用 `thread_exit` 终止当前线程，之所以选择 `thread_exit` 而非直接调用 `process_exit`，是因为前者还包含了对信号量的同步操作，并且可以用于终止用户进程和系统线程。

其中，子函数 `exit_process` 的实现如下：

```

1 void exit_process(int status) {
2     struct child_process *cp;
3     struct thread *cur_thread = thread_current();
4     enum intr_level old_level = intr_disable(); // 禁用中断以保证线程状态更新的安全性。
5
6     // 遍历当前线程的父线程的子线程列表，更新当前线程在父线程记录中的状态。
7     for (struct list_elem *e = list_begin(&cur_thread->parent->children_list);
8          e != list_end(&cur_thread->parent->children_list);
9          e = list_next(e)) {
10        cp = list_entry(e, struct child_process, child_elem);
11        if (cp->tid == cur_thread->tid) {
12            cp->if_waited = true; // 标记当前线程已被父线程等待过。
13            cp->exit_status = status; // 更新退出状态。
14        }
15    }
16 }

```

```

16     cur_thread->exit_status = status; // 设置当前线程的退出状态。
17     intr_set_level(old_level); // 恢复中断状态
18     thread_exit(); // 退出当前线程
19 }

```

## 7 syscall\_exec

```

1 // 运行其名称在 cmd_line 中给出的可执行文件，并传递任何给定的参数，返回新进程的进程ID(pid)
2 static void syscall_exec(struct intr_frame *f) {
3     char *cmd_line = *(char **)check_read_user_ptr(f->esp + ptr_size, ptr_size);
4     check_read_user_str(cmd_line);
5     f->eax = process_execute(cmd_line);
6 }

```

### Notice

`syscall_exec` 函数实现了 `exec` 系统调用，用于运行用户程序指定的可执行文件并传递命令行参数。函数首先从用户栈中读取命令行字符串 (`cmd_line`)，这是用户程序传递的参数。

随后，通过 `check_read_user_ptr` 验证命令行指针是否合法，确保指针指向的地址位于用户空间内，并进一步使用 `check_read_user_str` 验证字符串的内容是否安全可读，防止非法内存访问。验证通过后，函数调用 `process_execute(cmd_line)`，用于创建新进程并执行命令行中指定的程序。`process_execute` 返回新进程的进程 ID (pid)，代表子进程的唯一标识符，函数将其存储在中断帧的 `eax` 寄存器中，供用户程序访问。

## 8 syscall\_wait

```

1 int syscall_wait(struct intr_frame *f) {
2     tid_t child_tid;
3     pop_stack(f->esp, &child_tid, 1); // Extract the child process ID from the user stack.
4     return process_wait(child_tid); // Wait for the child process and return its exit status.
5 }

```

`syscall_wait()`

### Notice

`syscall_wait` 实现了 `WAIT` 系统调用，用于当前进程等待一个指定的子进程完成执行。它通过中断帧的栈指针 (`f->esp`) 从用户栈中获取传递的子进程线程 ID (`tid`)。

在未实现 `syscall_write` 之前，会出现这样的情况：

```

pintos -v -k -T 60 --qemu --filesystem-size=2 -p tests/userprog/args-none
2> tests/userprog/args-none.errors > tests/userprog/args-none.output
perl -I../.. ../tests/userprog/args-none.ck tests/userprog/args-none
FAIL tests/userprog/args-none
Test output failed to match any acceptable form.

Acceptable output:
(args) begin
(args) argc = 1
(args) argv[0] = 'args-none'
(args) argv[1] = null
(args) end
args-none: exit(0)

```

```

Differences in 'diff -u' format:
- (args) begin
- (args) argc = 1
- (args) argv[0] = 'args-none'
- (args) argv[1] = null
- (args) end
- args-none: exit(0)
+ Default 9
+ Default 9
+ Default 9
+ Default 9
+ Default 9
pintos -v -k -T 60 --qemu --filesystem-size=2 -p tests/userprog/args-single
rg' < /dev/null 2> tests/userprog/args-single.errors > tests/userprog/args
perl -I../.. ../tests/userprog/args-single.ck tests/userprog/args-sing
FAIL tests/userprog/args-single
Test output failed to match any acceptable form.

```

图 14: Combination of Failed Results

输出的结果显示为: **Default 9** 这在 `user/syscall.c` 中是 9 号, `SYS_WRITE`, 用户程序的打印 (输出到 `stdout`) 也是通过写的系统调用实现的。

## 8.1 syscall\_write

```

1  int syscall_write(struct intr_frame *f) {
2      int ret;                /**< 写入操作的返回值。 */
3      int size;               /**< 要写入的数据大小 (字节)。 */
4      void *buffer;           /**< 数据缓冲区的起始地址。 */
5      int fd;                 /**< 文件描述符。 */
6
7      // 从用户栈中读取系统调用参数, 分别是 size、buffer 和 fd
8      pop_stack(f->esp, &size, 7);
9      pop_stack(f->esp, &buffer, 6);
10     pop_stack(f->esp, &fd, 5);
11
12     // 检查 buffer 是否为有效的用户地址
13     if (!is_valid_addr(buffer))
14         ret = -1;
15
16     // 如果文件描述符是标准输出 (fd == 1), 直接将数据输出到控制台
17     if (fd == 1) {
18         putbuf(buffer, size); /**< 调用系统提供的 putbuf() 函数写入控制台。 */
19         ret = size;           /**< 返回写入的字节数, 即 size。 */
20     } else {
21         // 文件写入操作, 需要查找文件描述符对应的文件指针
22         enum intr_level old_level = intr_disable(); /**< 禁用中断, 避免多线程竞争。 */
23         struct process_file *pf = search_fd(&thread_current()->opened_files, fd);
24         intr_set_level(old_level); /**< 恢复中断状态。 */
25         // 如果文件描述符无效, 返回错误
26         if (pf == NULL) ret = -1;
27         else {
28             // 获取文件系统锁, 确保文件写操作是线程安全的
29             lock_acquire(&fileys_lock);
30             ret = file_write(pf->ptr, buffer, size); /**< 写入数据到文件, 返回写入字节数。 */
31             lock_release(&fileys_lock);
32         }
33     }
34     return ret;
35 }

```

syscall\_write()

### Notice

`syscall_write` 函数实现了系统调用, 用于向文件或标准输出写入数据。函数首先通过从用户栈中读取参数 (`size`、`buffer`、`fd`) 完成参数传递, 确保用户态与内核态的隔离。随后, 检查缓冲区地址是否为有效的用户地址, 以防止非法内存访问, 保护操作系统的安全性。如果文件描述符为标准输出 (`fd == 1`), 直接调用 `putbuf` 将数据输出到控制台, 绕过文件系统以提高效率。

对于其他文件描述符, 函数通过禁用中断和文件系统锁保证线程安全, 查询文件描述符对应的文件指针并执行文件写入操作, 同时确保多线程环境下的资源管理。

## 五 实验结果

### Question

使用 `make check` 命令时，但此时仍然显示：Run didn't produce any output.

查询资料后显示，这是因为残留的 `process_wait` 未实现导致用户程序还没执行整个系统就关机了。

[https://blog.csdn.net/Altair\\_alpha/article/details/126819252](https://blog.csdn.net/Altair_alpha/article/details/126819252)

根据博主写的彩蛋的内容作出相应的修改，重新 `make`，可以看到此时已经能够成功通过参数传递的五个测试点。  
在 `thread.c` 中添加以下函数：

```
1 int thread_dead(tid_t tid) {
2     struct list_elem *e;
3     for (e = list_begin (&all_list); e != list_end (&all_list); e = list_next (e)) {
4         struct thread *t = list_entry (e, struct thread, allelem);
5         if (t->tid == tid) return 0;
6     }
7     return 1;
8 }
```

将 `process_wait` 函数修改为如下所示：

### Notice

之前的问题是系统在用户程序执行完成前就关闭了，导致没有机会输出预期结果。

通过在 `process_wait()` 中添加一个无限循环，主线程可以持续等待子进程完成。这样，用户程序在系统关闭前获得了足够的时间完成执行，并输出结果。

```
1 int process_wait (tid_t child_tid) {
2     while (1) {
3         thread_yield();
4         if (thread_dead(child_tid)) break;
5     }
6     return -1;
7 }
8
9 int thread_dead(tid_t tid) {
10    struct list_elem *e;
11    for (e = list_begin (&all_list); e != list_end (&all_list); e = list_next (e)) {
12        struct thread *t = list_entry (e, struct thread, allelem);
13        if (t->tid == tid) return 0; // 如果找到目标线程，返回 0 (未死亡)
14    }
15    return 1; // 如果目标线程不在 'all_list' 中，返回 1 (已死亡)
16 }
```

- `all_list` 是所有线程的链表。一个线程在死亡后会从 `all_list` 中移除。
- 如果 `tid` 找不到匹配的线程，则认为线程已经死亡。

运行 `make clean` 和 `make check`

得到以下结果：



图 15: make check

```
C:\Windows\system32\cmd.e: × + ▾
pass tests/userprog/args-none
pass tests/userprog/args-single
pass tests/userprog/args-multiple
pass tests/userprog/args-many
pass tests/userprog/args-dbl-space
```

图 16: make check

## 参考资料

- Pintos 单一测试点代码: <https://pastebin.com/tFRfywvv>
- Pintos Lab2: 用户程序 User Programs (上): [https://blog.csdn.net/Altair\\_alpha/article/details/126819252](https://blog.csdn.net/Altair_alpha/article/details/126819252)
- Pintos 指导手册: <https://pkuflyingpig.gitbook.io/pintos/project-description/lab2-user-programs/suggestions>
- OS 实验: pintos project2: <https://zhuanlan.zhihu.com/p/382326973>
- 斯坦福大学 Pintos Project1、2 指南 + 总结: <https://zhuanlan.zhihu.com/p/104497182>