

实验目的

1. 掌握虚拟文件系统的实现原理
2. 实践文件、目录、索引节点等概念

实验内容

在 Linux 0.11 上实现 `procfs`（`proc` 文件系统）内的 `psinfo` 结点。当读取此结点的内容时，可得到系统当前所有进程的状态信息。例如，用 `cat` 命令显示 `/proc/psinfo` 的内容，可得到：

```
# cat /proc/psinfo
pid    state  father  counter  start_time
0      1      -1      0        0
1      1      0       28       1
4      1      1       1       73
3      1      1       27       63
6      0      4       12      817
# cat /proc/hdinfo
total_blocks:62000;
free_blocks:39037;
used_blocks:22963;
total_inodes:20666;
...
```

`procfs` 及其结点要在内核启动时自动创建。相关功能实现在 `fs/proc.c` 文件内。

实验报告

完成实验后，在实验报告中回答如下问题：

1. 如果要求你在 `psinfo` 之外再实现另一个结点，具体内容自选，那么你会实现一个给出什么信息的结点？为什么？
2. 一次 `read()` 未必能读出所有的数据，需要继续 `read()`，直到把数据读空为止。而数次 `read()` 之间，进程的状态可能会发生变化。你认为后几次 `read()` 传给用户的数据，应该是变化后的，还是变化前的？
3. 如果是变化后的，那么用户得到的数据衔接部分是否会有混乱？如何防止混乱？
4. 如果是变化前的，那么该在什么样的情况下更新 `psinfo` 的内容？
5. 删除文件以后，`/proc/inodeinfo` 那个 `inode` 号的 `inode`，你发现了什么，为什么会这样？

评分标准

- 自动创建 `/proc`, `/proc/psinfo`; `/proc/hdinfo`; `/proc/inodeinfo`, 20%
- `psinfo` 内容可读, 内容符合题目要求, 40%
- `hdinfo` 内容可读, 符合题目要求, 30%
- 实验报告, 10%

实验提示

procfs 简介

正式的 Linux 内核实现了 `procfs`, 它是一个虚拟文件系统, 通常被 `mount` 到 `/proc` 目录上, 通过虚拟文件和虚拟目录的方式提供访问系统参数的机会, 所以有人称它为“了解系统信息的一个窗口”。这些虚拟的文件和目录并没有真实地存在在磁盘上, 而是内核中各种数据的一种直观表示。虽然是虚拟的, 但它们都可以通过标准的系统调用 (`open()`、`read()`等) 访问。

例如, `/proc/meminfo` 中包含内存使用的信息, 可以用 `cat` 命令显示其内容:

```
$ cat /proc/meminfo
MemTotal:      384780 kB
MemFree:       13636 kB
Buffers:       13928 kB
Cached:        101680 kB
SwapCached:    132 kB
Active:        207764 kB
Inactive:      45720 kB
SwapTotal:     329324 kB
SwapFree:      329192 kB
Dirty:         0 kB
Writeback:     0 kB
.....
```

其实, Linux 的很多系统命令就是通过读取 `/proc` 实现的。例如 `uname -a` 的部分信息就来自 `/proc/version`, 而 `uptime` 的部分信息来自 `/proc/uptime` 和 `/proc/loadavg`。

关于 `procfs` 更多的信息请访问: <http://en.wikipedia.org/wiki/Procfs>

基本思路

Linux 是通过文件系统接口实现 **procfs**，并在启动时自动将其 **mount** 到 **/proc** 目录上。此目录下的所有内容都是随着系统的运行自动建立、删除和更新的，而且它们完全存在于内存中，不占用任何外存空间。

Linux 0.11 还没有实现虚拟文件系统，也就是，还没有提供增加新文件系统支持的接口。所以本实验只能在现有文件系统的基础上，通过打补丁的方式模拟一个 **procfs**。

Linux 0.11 使用的是 **Minix** 的文件系统，这是一个典型的基于 **inode** 的文件系统，《注释》一书对它有详细描述。它的每个文件都要对应至少一个 **inode**，而 **inode** 中记录着文件的各种属性，包括文件类型。文件类型有普通文件、目录、字符设备文件和块设备文件等。在内核中，每种类型的文件都有不同的处理函数与之对应。我们可以增加一种新的文件类型——**proc** 文件，并在相应的处理函数内实现 **procfs** 要实现的功能。

增加新文件类型

在 `include/sys/stat.h` 文件中定义了几种文件类型和相应的测试宏：

```
#define S_IFMT 00170000
#define S_IFREG 0100000 //普通文件
#define S_IFBLK 0060000 //块设备
#define S_IFDIR 0040000 //目录
#define S_IFCHR 0020000 //字符设备
#define S_IFIFO 0010000
.....

#define S_ISREG(m)      (((m) & S_IFMT) == S_IFREG) //测试 m 是否是普通文件
#define S_ISDIR(m)      (((m) & S_IFMT) == S_IFDIR) //测试 m 是否是目录
#define S_ISCHR(m)      (((m) & S_IFMT) == S_IFCHR) //测试 m 是否是字符设备
#define S_ISBLK(m)      (((m) & S_IFMT) == S_IFBLK) //测试 m 是否是块设备
#define S_ISFIFO(m)     (((m) & S_IFMT) == S_IFIFO)
```

增加新的类型的方法分两步：

1. 定义一个类型宏 **S_IFPROC**，其值应在 00100000 到 01000000 之间，但后四位八进制数必须是 0（这是 **S_IFMT** 的限制，分析测试宏可知原因），而且不能和已有的任意一个 **S_IFXXX** 相同；
2. 定义一个测试宏 **S_ISPROC(m)**，形式仿照其它的 **S_ISXXX(m)**

注意，C 语言中以“0”直接接数字的常数是八进制数。

让 **mknod()** 支持新的文件类型

psinfo 结点要通过 `mknod()` 系统调用建立，所以要让它支持新的文件类型。直接修改 `fs/namei.c` 文件中的 `sys_mknod()` 函数中的一行代码，如下：

```
if (S_ISBLK(mode) || S_ISCHR(mode) || S_ISPROC(mode))
    inode->i_zone[0] = dev;
```

文件系统初始化

内核初始化的全部工作是在 `main()` 中完成，而 `main()` 在最后从内核态切换到用户态，并调用 `init()`。`init()` 做的第一件事情就是挂载根文件系统：

```
void init(void)
{
    .....
    setup((void *) &drive_info);
    .....
}
```

`procfs` 的初始化工作应该在根文件系统挂载之后开始。它包括两个步骤：

1. 建立 `/proc` 目录；
2. 建立 `/proc` 目录下的各个结点。本实验只建立 `/proc/psinfo`。

建立目录和结点分别需要调用 `mkdir()` 和 `mknod()` 系统调用。因为初始化时已经在用户态，所以不能直接调用 `sys_mkdir()` 和 `sys_mknod()`。必须在初始化代码所在文件中实现这两个系统调用的用户态接口，即 API：

```
#include
#define __LIBRARY__
#include

_syscall2(int,mkdir,const char*,name,mode_t,mode)
_syscall3(int,mknod,const char*,filename,mode_t,mode,dev_t,dev)
```

`mkdir()` 时 `mode` 参数的值可以是“0755”（`rwxr-xr-x`），表示只允许 `root` 用户改写此目录，其它人只能进入和读取此目录。

`procfs` 是一个只读文件系统，所以用 `mknod()` 建立 `psinfo` 结点时，必须通过 `mode` 参数将其设为只读。建议使用“`S_IFPROC|0444`”做为 `mode` 值，表示这是一个 `proc` 文件，权限为 0444（`r--r--r--`），对所有用户只读。

`mknod()` 的第三个参数 `dev` 用来说明结点所代表的设备编号。对于 `procfs` 来说，此编号可以完全自定义。`proc` 文件的处理函数将通过这个编号决定对应文件包含的信息是什么。例如，可以把 0 对应 `psinfo`，1 对应 `meminfo`，2 对应 `cpuinfo`。

如此项工作完成得没有问题，那么编译、运行 0.11 内核后，用“ll /proc”可以看到：

```
# ll /proc
total 0
?r--r--r--  1 root    root          0 ??? ??  ???? psinfo
```

此时可以试着读一下此文件：

```
# cat /proc/psinfo
(Read)inode->i_mode=XXX444
cat: /proc/psinfo: EINVAL
```

inode->i_mode 就是通过 `mknod()` 设置的 mode。信息中的 XXX 和你设置的 `S_IFPROC` 有关。通过此值可以了解 `mknod()` 工作是否正常。这些信息说明内核在对 `psinfo` 进行读操作时不能正确处理，向 `cat` 返回了 `EINVAL` 错误。因为还没有实现处理函数，所以这是很正常的。

这些信息至少说明，`psinfo` 被正确 `open()` 了。所以我们不需要对 `sys_open()` 动任何手脚，唯一要打补丁的，是 `sys_read()`。

让 proc 文件可读

`open()` 没有变化，那么需要修改的就是 `sys_read()` 了。首先分析 `sys_read`（在文件 `fs/read_write.c` 中）：

```
int sys_read(unsigned int fd, char * buf, int count)
{
    struct file * file;
    struct m_inode * inode;
    .....
    inode = file->f_inode;
    if (inode->i_pipe)
        return (file->f_mode&1)?read_pipe(inode, buf, count):-EIO;
    if (S_ISCHR(inode->i_mode))
        return rw_char(READ, inode->i_zone[0], buf, count, &file->f_pos);
    if (S_ISBLK(inode->i_mode))
        return block_read(inode->i_zone[0], &file->f_pos, buf, count);
    if (S_ISDIR(inode->i_mode) || S_ISREG(inode->i_mode)) {
        if (count+file->f_pos > inode->i_size)
            count = inode->i_size - file->f_pos;
        if (count<=0)
            return 0;
        return file_read(inode, file, buf, count);
    }
}
```

```
    printk("(Read)inode->i_mode=%06o\n\r",inode->i_mode);    //这条信息很面
善吧?
    return -EINVAL;
}
```

显然，要在这里一群 if 的排比中，加上 **S_IFPROC()** 的分支，进入对 **proc** 文件的处理函数。需要传给处理函数的参数包括：

1. **inode->i_zone[0]**，这就是 **mknod()** 时指定的 **dev**——设备编号
2. **buf**，指向用户空间，就是 **read()** 的第二个参数，用来接收数据
3. **count**，就是 **read()** 的第三个参数，说明 **buf** 指向的缓冲区大小
4. **&file->f_pos**，**f_pos** 是上一次读文件结束时“文件位置指针”的指向。这里必须传指针，因为处理函数需要根据传给 **buf** 的数据量修改 **f_pos** 的值。

proc 文件的处理函数

proc 文件的处理函数的功能是根据设备编号，把不同的内容写入到用户空间的 **buf**。写入的数据要从 **f_pos** 指向的位置开始，每次最多写 **count** 个字节，并根据实际写入的字节数调整 **f_pos** 的值，最后返回实际写入的字节数。当设备编号表明要读的是 **psinfo** 的内容时，就要按照 **psinfo** 的形式组织数据。

实现此函数可能要用到如下几个函数：

malloc()和 free()

包含 **linux/kernel.h** 头文件后，就可以使用 **malloc()** 和 **free()** 函数。它们是可以被核心态代码调用的，唯一的限制是一次申请的内存大小不能超过一个页面。

sprintf()

Linux 0.11 没有 **sprintf()**，可以参考 **printf()** 自己实现一个，如下：

```
#include <stdarg.h>
.....
int sprintf(char *buf, const char *fmt, ...)
{
    va_list args; int i;
    va_start(args, fmt);
    i=vsprintf(buf, fmt, args);
    va_end(args);
    return i;
}
```

cat 命令

cat 是 Linux 下的一个常用命令，功能是将文件的内容打印到标准输出。它核心实现大体如下：

```
#include <stdio.h>
#include <unistd.h>
int main(int argc, char* argv[])
{
    char buf[513] = {'\0'};
    int nread;

    int fd = open(argv[1], O_RDONLY, 0);
    while(nread = read(fd, buf, 512))
    {
        buf[nread] = '\0';
        puts(buf);
    }

    return 0;
}
```

psinfo 的内容

进程的信息就来源于内核全局结构数组 `struct task_struct * task[NR_TASKS]` 中，具体读取细节可参照 `sched.c` 中的函数 `schedule()`>

```
for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
    if (*p)
        (*p)->counter = ((*p)->counter >> 1)+...;
```

hdinfo 的内容

硬盘总共有多少块，多少块空闲，有多少 `inode` 等信息都放在 `super` 块中，`super` 块可以通过 `get_super()` 函数获得，其中的信息可以借鉴如下代码。

```
struct super_block * sb;
sb=get_super(inode->i_dev);
struct buffer_head * bh;
total_blocks = sb->s_nzones;
for(i=0; is_zmap_blocks; i++)
{
    bh = sb->s_zmap[i];
```

```
p=(char *)bh->b_data;
```