

xv6 中文文档

鲜染

Published
with GitBook



目錄

1. [介紹](#)
2. [封面](#)
3. [前言](#)
4. [第零章 操作系统接口](#)
5. [第一章 第一个进程](#)
6. [第二章 页表](#)
7. [第三章 陷入，中断和驱动程序](#)
8. [第四章 锁](#)
9. [第五章 调度](#)
10. [第六章 文件系统](#)
11. [附录A PC 硬件](#)
12. [附录B 引导加载器](#)
13. [术语](#)
14. [勘误](#)

xv6 中文文档

xv6 是 MIT 开发的一个教学用的完整的类 Unix 操作系统，并且在 MIT 的操作系统课程 [6.828](#) 中使用。通过阅读并理解 xv6 的代码，可以清楚地了解操作系统中众多核心的概念，对操作系统感兴趣的同学十分推荐一读！这份文档是中文翻译的 MIT xv6 文档，是阅读代码过程中非常好的参考资料。

[原文在此](#)

[文中引用的 xv6 源代码](#)

强烈推荐 xv6 源代码同本书一同阅读！原作和翻译中遇到的括号内的数字，都是指上面链接中文件的源代码行号。

同时，我们的翻译文档也可以通过 [gitbook](#) 阅读

译者

- 鲜染 北京大学 信息科学技术学院 计算机系
- 赵天雨 北京大学 信息科学技术学院 计算机系
- 胡树伟 北京大学 信息科学技术学院 计算机系 (I guess)
- 胡文涛 KAUST CS
- 曹扬 上海交通大学 电子信息与电气工程学院 计算机系

如果你愿意贡献，你的名字也会出现在这里！

翻译状况

| 章节 | 初稿 | 审校 | 二审 |
|----|----|----|----|
| 封面 | + | + | + |
| 前言 | + | + | + |
| 零 | + | + | + |
| 一 | + | + | |
| 二 | + | + | |
| 三 | + | + | |
| 四 | + | + | |
| 五 | + | + | |
| 六 | + | + | + |
| 附A | + | + | + |
| 附B | + | + | + |

参与审校

热情欢迎大家参与到审校工作中！请访问 <https://github.com/Th0ar/xv6-chinese>

1. Fork

2. 审校并修改，保证修改后 markdown 解析正确
3. 发送 Pull Request
4. 等待你的名字（不久后）出现在译者列表中！

许可证（License）

文档中涉及到的 xv6 源代码使用 [MIT](#) 许可证。中文翻译使用 [GNU GPL V3.0](#) 许可证，在 GNU GPL V3.0 之上，转载和引用须注明本项目 Github 地址。

xv6

一个简单，类 **Unix** 的教学操作系统

Russ Cox

Frans Kasshoek

Robert Morris

xv6-book@pdos.csail.mit.edu

翻译：

鲜染

xianran@pku.edu.cn

赵天雨

zhaoty.ting@gmail.com

2013年国庆

前言和致谢

这是一份为操作系统课编写的教学草案。它通过研究一个名为 xv6 的操作系统内核来解释操作系统中的主要概念。xv6 是 Dennis Ritchie 和 Ken Thompson 合著的 Unix Version 6 (v6) 操作系统的重新实现。xv6 在一定程度上遵守 v6 的结构和风格，但它是用 ANSI C 实现的，并且是基于 x86 多核处理器的。

这本教材应该和 xv6 源代码一起阅读，这是 John Lion 在 Unix 6th Edition (Peer to Peer Communications ; ISBN : 1-57398-013-7 ; 第一版 (2000年7月14日) 的评注中推荐的学习方式。参见<http://pdos.csail.mit.edu/6.828> 上有关于 v6 和 xv6 的资料。

我们已经在 6.828 —— MIT 的操作系统课程中使用了这本教材。我们向参与 6.828 的教职员工、助教和学生表示感谢，他们都直接或间接向 xv6 做出了贡献。此处我们要特别感谢 Austin Clements 和 Nicholai Zeldovich。

第0章

操作系统接口

操作系统的工作是(1)将计算机的资源在多个程序间共享，并且给程序提供一系列比硬件本身更有用的服务。(2)管理并抽象底层硬件，举例来说，一个文字处理软件（比如 word）不用去关心自己使用的是何种硬盘。(3)多路复用硬件，使得多个程序可以(至少看起来是)同时运行的。(4)最后，给程序间提供一种受控的交互方式，使得程序之间可以共享数据、共同工作。

操作系统通过接口向用户程序提供服务。设计一个好的接口实际是很难的。一方面我们希望接口设计得简单和精准，使其易于正确地实现；另一方面，我们可能忍不住想为应用提供一些更加复杂的功能。解决这种矛盾的办法是让接口的设计依赖于少量的机制（*mechanism*），而通过这些机制的组合提供强大、通用的功能。

本书通过 xv6 操作系统来阐述操作系统的概念，它提供 Unix 操作系统中的基本接口（由 Ken Thompson 和 Dennis Ritchie 引入），同时模仿 Unix 的内部设计。Unix 里机制结合良好的窄接口提供了令人吃惊的通用性。这样的接口设计非常成功，使得包括 BSD, Linux, Mac OS X, Solaris（甚至 Microsoft Windows 在某种程度上）都有类似 Unix 的接口。理解 xv6 是理解这些操作系统的一个良好起点。

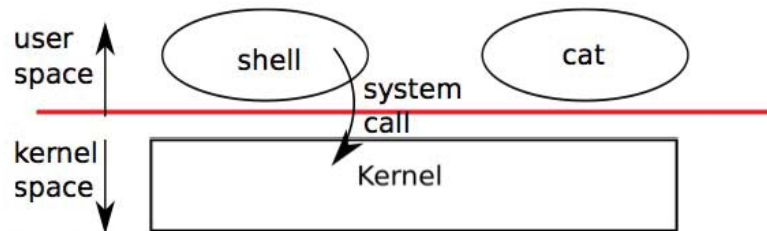


Figure 0-1. A kernel and two user processes.

如图0-1所示，xv6 使用了传统的内核概念 - 一个向其他运行中程序提供服务的特殊程序。每一个运行中程序（称之为进程）都拥有包含指令、数据、栈的内存空间。指令实现了程序的运算，数据是用于运算过程的变量，栈管理了程序的过程调用。

进程通过系统调用使用内核服务。系统调用会进入内核，让内核执行服务然后返回。所以进程总是在用户空间和内核空间之间交替运行。

内核使用了 CPU 的硬件保护机制来保证用户进程只能访问自己的内存空间。内核拥有实现保护机制所需的硬件权限 (hardware privileges)，而用户程序没有这些权限。当一个用户程序进行一次系统调用时，硬件会提升特权级并且开始执行一些内核中预定义的功能。

内核提供的一系列系统调用就是用户程序可见的操作系统接口，xv6 内核提供了 Unix 传统系统调用的一部分，它们是：

| 系统调用 | 描述 |
|-----------------------|----------------|
| fork() | 创建进程 |
| exit() | 结束当前进程 |
| wait() | 等待子进程结束 |
| kill(pid) | 结束 pid 所指进程 |
| getpid() | 获得当前进程 pid |
| sleep(n) | 睡眠 n 秒 |
| exec(filename, *argv) | 加载并执行一个文件 |
| sbrk(n) | 为进程内存空间增加 n 字节 |
| | |

| | |
|---------------------------|----------------------|
| open(filename, flags) | 打开文件，flags 指定读/写模式 |
| read(fd, buf, n) | 从文件中读 n 个字节到 buf |
| write(fd, buf, n) | 从 buf 中写 n 个字节到文件 |
| close(fd) | 关闭打开的 fd |
| dup(fd) | 复制 fd |
| pipe(p) | 创建管道， 并把读和写的 fd 返回到p |
| chdir(dirname) | 改变当前目录 |
| mkdir(dirname) | 创建新的目录 |
| mknod(name, major, minor) | 创建设备文件 |
| fstat(fd) | 返回文件信息 |
| link(f1, f2) | 给 f1 创建一个新名字(f2) |
| unlink(filename) | 删除文件 |

这一章剩下的部分将说明 xv6 系统服务的概貌 —— 进程，内存，文件描述符，管道和文件系统，为了描述他们，我们给出了代码和一些讨论。这些系统调用在 shell 上的应用阐述了他们的设计是多么独具匠心。

shell 是一个普通的程序，它接受用户输入的命令并且执行它们，它也是传统 Unix 系统中最基本的用户界面。shell 作为一个普通程序，而不是内核的一部分，充分说明了系统调用接口的强大：shell 并不是一个特别的用户程序。这也意味着 shell 是很容易被替代的，实际上这导致了现代 Unix 系统有着各种各样的 shell，每一个都有着自己的用户界面和脚本特性。xv6 shell 本质上是一个 Unix Bourne shell 的简单实现。它的实现在第 7850 行。

进程和内存

一个 xv6 进程由两部分组成，一部分是用户内存空间（指令，数据，栈），另一部分是仅对内核可见的进程状态。xv6 提供了分时特性：它在可用 CPU 之间不断切换，决定哪一个等待中的进程被执行。当一个进程不在执行时，xv6 保存它的 CPU 寄存器，当他们再次被执行时恢复这些寄存器的值。内核将每个进程和一个 **pid** (process identifier) 关联起来。

一个进程可以通过系统调用 `fork` 来创建一个新的进程。`fork` 创建的新进程被称为子进程，子进程的内存内容同创建它的进程（父进程）一样。`fork` 函数在父进程、子进程中都返回（一次调用两次返回）。对于父进程它返回子进程的 pid，对于子进程它返回 0。考虑下面这段代码：

```
int pid;
pid = fork();
if(pid > 0){
    printf("parent: child=%d\n", pid);
    pid = wait();
    printf("child %d is done\n", pid);
} else if(pid == 0){
    printf("child: exiting\n");
    exit();
} else {
    printf("fork error\n");
}
```

系统调用 `exit` 会导致调用它的进程停止运行，并且释放诸如内存和打开文件在内的资源。系统调用 `wait` 会返回一个当前进程已退出的子进程，如果没有子进程退出，`wait` 会等候直到有一个子进程退出。在上面的例子中，下面的两行输出

```
parent: child=1234
child: exiting
```


可能以任意顺序被打印，这种顺序由父进程或子进程谁先结束 `printf` 决定。当子进程退出时，父进程的 `wait` 也就返回了，于是父进程打印：

```
parent: child 1234 is done
```

需要留意的是父子进程拥有不同的内存空间和寄存器，改变一个进程中的变量不会影响另一个进程。

系统调用 `exec` 将从某个文件（通常是可执行文件）里读取内存镜像，并将其替换到调用它的进程的内存空间。这份文件必须符合特定的格式，规定文件的哪一部分是指令，哪一部分是数据，哪里是指令的开始等等。xv6 使用 ELF 文件格式，第2章将详细介绍它。当 `exec` 执行成功后，它并不返回到原来的调用进程，而是从 ELF 头中声明的入口开始，执行从文件中加载的指令。`exec` 接受两个参数：可执行文件名和一个字符串参数数组。举例来说：

```
char *argv[3];
argv[0] = "echo";
argv[1] = "hello";
argv[2] = 0;
exec("/bin/echo", argv);
printf("exec error\n");
```

这段代码将调用程序替换为 `/bin/echo` 这个程序，这个程序的参数列表为 `echo hello`。大部分的程序都忽略第一个参数，这个参数惯例上是程序的名字（此例是 `echo`）。

xv6 shell 用以上调用为用户执行程序。shell 的主要结构很简单，详见 `main` 的代码（8001）。主循环通过 `getcmd` 读取命令行的输入，然后它调用 `fork` 生成一个 shell 进程的副本。父 shell 调用 `wait`，而子进程执行用户命令。举例来说，用户在命令行输入“`echo hello`”，`getcmd` 会以 `echo hello` 为参数调用 `runcmd`（7906），由 `runcmd` 执行实际的命令。对于 `echo hello`，`runcmd` 将调用 `exec`。如果 `exec` 成功被调用，子进程就会转而去执行 `echo` 程序里的指令。在某个时刻 `echo` 会调用 `exit`，这会使得其父进程从 `wait` 返回。你可能会疑惑为什么 `fork` 和 `exec` 为什么没有被合并成一个调用，我们之后将发现，将创建进程——加载程序分为两个过程是一个非常机智的设计。

xv6 通常隐式地分配用户的内存空间。`fork` 在子进程需要装入父进程的内存拷贝时分配空间，`exec` 在需要装入可执行文件时分配空间。一个进程在需要额外内存时可以通过调用 `sbrk(n)` 来增加 `n` 字节的数据内存。`sbrk` 返回新的内存的地址。

xv6 没有用户这个概念当然更没有不同用户间的保护隔离措施。按照 Unix 的术语来说，所有的 xv6 进程都以 `root` 用户执行。

I/O 和文件描述符

文件描述符是一个整数，它代表了一个进程可以读写的被内核管理的对象。进程可以通过多种方式获得一个文件描述符，如打开文件、目录、设备，或者创建一个管道（pipe），或者复制已经存在的文件描述符。简单起见，我们常常把文件描述符指向的对象称为“文件”。文件描述符的接口是对文件、管道、设备等的抽象，这种抽象使得它们看上去就是字节流。

每个进程都有一张表，而 xv6 内核就以文件描述符作为这张表的索引，所以每个进程都有一个从 0 开始的文件描述符空间。按照惯例，进程从文件描述符 0 读入（标准输入），从文件描述符 1 输出（标准输出），从文件描述符 2 输出错误（标准错误输出）。我们会看到 shell 正是利用了这种惯例来实现 I/O 重定向。shell 保证在任何时候都有 3 个打开的文件描述符（8007），他们是控制台（console）的默认文件描述符。

系统调用 `read` 和 `write` 从文件描述符所指的文件中读或者写 `n` 个字节。`read(fd, buf, n)` 从 `fd` 读最多 `n` 个字节（`fd` 可能没有 `n` 个字节），将它们拷贝到 `buf` 中，然后返回读出的字节数。每一个指向文件的文件描述符都和一个偏移关联。`read` 从当前文件偏移处读取数据，然后把偏移增加读出字节数。紧随其后的 `read` 会从新的起点开始读数据。当没有数据可读时，`read` 就会返回 0，这就表示文件结束了。

`write(fd, buf, n)` 写 `buf` 中的 `n` 个字节到 `fd` 并且返回实际写出的字节数。如果返回值小于 `n` 那么只可能是发生了错误。就像 `read` 一样，`write` 也从当前文件的偏移处开始写，在写的过程中增加这个偏移。

下面这段程序（实际上就是 `cat` 的本质实现）将数据从标准输入复制到标准输出，如果遇到了错误，它会在标准错误输出输出一条信息。

```
char buf[512];
int n;

for(;;){
    n = read(0, buf, sizeof buf);
    if(n == 0)
        break;
    if(n < 0){
        fprintf(2, "read error\n");
        exit();
    }
    if(write(1, buf, n) != n){
        fprintf(2, "write error\n");
        exit();
    }
}
```

这段代码中值得一提的是 `cat` 并不知道它是从文件、控制台或者管道中读取数据的。同样地 `cat` 也不知道它是写到文件、控制台或者别的什么地方。文件描述符的使用和一些惯例（如0是标准输入，1是标准输出）使得我们可以轻松实现 `cat`。

系统调用 `close` 会释放一个文件描述符，使得它未来可以被 `open`，`pipe`，`dup` 等调用重用。一个新分配的文件描述符永远都是当前进程的最小的未被使用的文件描述符。

文件描述符和 `fork` 的交叉使用使得 I/O 重定向能够轻易实现。`fork` 会复制父进程的文件描述符和内存，所以子进程和父进程的文件描述符一模一样。`exec` 会替换调用它的进程的内存但是会保留它的文件描述符表。这种行为使得 shell 可以这样实现重定向：`fork` 一个进程，重新打开指定文件的文件描述符，然后执行新的程序。下面是一个简化版的 shell 执行 `cat<input.txt` 的代码：

```
char *argv[2];
argv[0] = "cat";
### argv[1] = 0;
if(fork() == 0) {
    close(0);
    open("input.txt", O_RDONLY);
    exec("cat", argv);
}
```

子进程关闭文件描述符0后，我们可以保证 `open` 会使用0作为新打开的文件 `input.txt` 的文件描述符（因为0是 `open` 执行时的最小可用文件描述符）。之后 `cat` 就会在标准输入指向 `input.txt` 的情况下运行。

xv6 的 shell 正是这样实现 I/O 重定位的（7930）。在 shell 的代码中，记得这时 `fork` 出了子进程，在子进程中 `runcmd` 会调用 `exec` 加载新的程序。现在你应该很清楚为何 `fork` 和 `exec` 是单独的两种系统调用了吧。这种区分使得 shell 可以在子进程执行指定程序之前对子进程进行修改。

虽然 `fork` 复制了文件描述符，但每一个文件当前的偏移仍然是在父子进程之间共享的，考虑下面这个例子：

```
if(fork() == 0) {
    write(1, "hello ", 6);
    exit();
} else {
    wait();
    write(1, "world\n", 6);
}
```

在这段代码的结尾，绑定在文件描述符1上的文件有数据"hello world"，父进程的 `write` 会从子进程 `write` 结束的地方继续写（因为 `wait`，父进程只在子进程结束之后才运行 `write`）。这种行为有利于顺序执行的 shell 命令的顺序输出，例如 `(echo`

```
hello; echo world)>output.txt。
```

`dup` 复制一个已有的文件描述符，返回一个指向同一个输入/输出对象的新描述符。这两个描述符共享一个文件偏移，正如被 `fork` 复制的文件描述符一样。这里有另一种打印“hello world”的办法：

```
fd = dup(1);
write(1, "hello", 6);
write(fd, "world\n", 6);
```

从同一个原初文件描述符通过一系列 `fork` 和 `dup` 调用产生的文件描述符都共享同一个文件偏移，而其他情况下产生的文件描述符就不是这样了，即使他们打开的都是同一份文件。`dup` 允许 shell 像这样实现命令：`ls existing-file non-exsiting-file > tmp1 2>&1 . 2>&1` 告诉 shell 给这条命令一个复制描述符1的描述符2。这样 `existing-file` 的名字和 `non-exsiting-file` 的错误输出都将出现在 `tmp1` 中。xv6 shell 并未实现标准错误输出的重定向，但现在你知道该怎么去实现它。

文件描述符是一个强大的抽象，因为他们将他们所连接的细节隐藏起来了：一个进程向描述符1写出，它有可能是写到一份文件，一个设备（如控制台），或一个管道。

管道

管道是一个小的内核缓冲区，它以文件描述符对的形式提供给进程，一个用于写操作，一个用于读操作。从管道的一端写的数据可以从管道的另一端读取。管道提供了一种进程间交互的方式。

接下来的示例代码运行了程序 `wc`，它的标准输出绑定到了一个管道的读端口。

```
int p[2];
char *argv[2];
argv[0] = "wc";
argv[1] = 0;
pipe(p);
if(fork() == 0) {
    close(0);
    dup(p[0]);
    close(p[0]);
    close(p[1]);
    exec("/bin/wc", argv);
} else {
    write(p[1], "hello world\n", 12);
    close(p[0]);
    close(p[1]);
}
```

这段程序调用 `pipe`，创建一个新的管道并且将读写描述符记录在数组 `p` 中。在 `fork` 之后，父进程和子进程都有了指向管道的文件描述符。子进程将管道的读端口拷贝在描述符0上，关闭 `p` 中的描述符，然后执行 `wc`。当 `wc` 从标准输入读取时，它实际上是从管道读取的。父进程向管道的写端口写入然后关闭它的两个文件描述符。

如果数据没有准备好，那么对管道执行的 `read` 会一直等待，直到有数据了或者其他绑定在这个管道写端口的描述符都已经关闭了。在后一种情况中，`read` 会返回 0，就像是一份文件读到了最后。读操作会一直阻塞直到不可能再有新数据到来了，这就是为什么我们在执行 `wc` 之前要关闭子进程的写端口。如果 `wc` 指向了一个管道的写端口，那么 `wc` 就永远看不到 eof 了。

xv6 shell 对管道的实现（比如 `fork sh.c | wc -l`）和上面的描述是类似的（7950行）。子进程创建一个管道连接管道的左右两端。然后它为管道左右两端都调用 `runcmd`，然后通过两次 `wait` 等待左右两端结束。管道右端可能也是一个带有管道的指令，如 `a | b | c`，它 `fork` 两个新的子进程（一个 `b` 一个 `c`），因此，shell 可能创建出一颗进程树。树的叶子节点是命令，中间节点是进程，它们会等待左子和右子执行结束。理论上，你可以让中间节点都运行在管道的左端，但做的如此精确会使得实现变得复杂。

`pipe` 可能看上去和临时文件没有什么两样：命令

```
echo hello world | wc
```

可以用无管道的方式实现：

```
echo hello world > /tmp/xyz; wc < /tmp/xyz
```

但管道和临时文件起码有三个关键的不同点。首先，管道会进行自我清扫，如果是 shell 重定向的话，我们必须要在任务完成后删除 `/tmp/xyz`。第二，管道可以传输任意长度的数据。第三，管道允许同步：两个进程可以使用一对管道来进行二者之间的信息传递，每一个读操作都阻塞调用进程，直到另一个进程用 `write` 完成数据的发送。

文件系统

xv6 文件系统提供文件和目录，文件就是一个简单的字节数组，而目录包含指向文件和其他目录的引用。xv6 把目录实现为一种特殊的文件。目录是一棵树，它的根节点是一个特殊的目录 `root`。`/a/b/c` 指向一个在目录 `b` 中的文件 `c`，而 `b` 本身又是在目录 `a` 中的，`a` 又是处在 `root` 目录下的。不从 `/` 开始的目录表示的是相对调用进程当前目录的目录，调用进程的当前目录可以通过 `chdir` 这个系统调用进行改变。下面的这些代码都打开同一个文件（假设所有涉及到的目录都是存在的）。

```
chdir("/a");
chdir("b");
open("c", O_RDONLY);

open("/a/b/c", O_RDONLY);
```

第一个代码段将当前目录切换到 `/a/b`；第二个代码片段则对当前目录不做任何改变。

有很多的系统调用可以创建一个新的文件或者目录：`mkdir` 创建一个新的目录，`open` 加上 `O_CREATE` 标志打开一个新的文件，`mknod` 创建一个新的设备文件。下面这个例子说明了这3种调用：

```
mkdir("/dir");
fd = open("/dir/file", O_CREATE|O_WRONLY);
close(fd);
mknod("/console", 1, 1);
```

`mknod` 在文件系统中创建一个文件，但是这个文件没有任何内容。相反，这个文件的元信息标志它是一个设备文件，并且记录主设备号和辅设备号（`mknod` 的两个参数），这两个设备号唯一确定一个内核设备。当一个进程之后打开这个文件的时候，内核将读、写的系统调用转发到内核设备的实现上，而不是传递给文件系统。

`fstat` 可以获取一个文件描述符指向的文件的信息。它填充一个名为 `stat` 的结构体，它在 `stat.h` 中定义为：

```
#define T_DIR 1
#define T_FILE 2
#define T_DEV 3
// Directory
// File
// Device
struct stat {
    short type; // Type of file
    int dev; // File system's disk device
    uint ino; // Inode number
    short nlink; // Number of links to file
    uint size; // Size of file in bytes
};
```

文件名和这个文件本身是有很大的区别。同一个文件（称为 `inode`）可能有多个名字，称为连接（`links`）。系统调用 `link` 创建另一个文件系统的名称，它指向同一个 `inode`。下面的代码创建了一个既叫做 `a` 又叫做 `b` 的新文件。

```
open("a", O_CREATE|O_WRONLY);
link("a", "b");
```

读写 `a` 就相当于读写 `b`。每一个 inode 都由一个唯一的 inode 号直接确定。在上面这段代码中，我们可以通过 `fstat` 知道 `a` 和 `b` 都指向同样的内容：`a` 和 `b` 都会返回同样的 inode 号（`ino`），并且 `nlink` 数会设置为 2。

系统调用 `unlink` 从文件系统移除一个文件名。一个文件的 inode 和磁盘空间只有当它的链接数变为 0 的时候才会被清空，也就是没有一个文件再指向它。因此在上面的代码最后加上

```
unlink("a"),
```

我们同样可以通过 `b` 访问到它。另外，

```
fd = open("/tmp/xyz", O_CREATE|O_RDWR);
unlink("/tmp/xyz");
```

是创建一个临时 inode 的最佳方式，这个 inode 会在进程关闭 `fd` 或者退出的时候被清空。

xv6 关于文件系统的操作都被实现为用户程序，诸如 `mkdir`，`ln`，`rm` 等等。这种设计允许任何人都可以通过用户命令拓展 shell。现在看起来这种设计是很显然的，但是 Unix 时代的其他系统的设计都将这样的命令内置在了 shell 中，而 shell 又是内置在内核中的。

有一个例外，那就是 `cd`，它是在 shell 中实现的（8016）。`cd` 必须改变 shell 自身的当前工作目录。如果 `cd` 作为一个普通命令执行，那么 shell 就会 `fork` 一个子进程，而子进程会运行 `cd`，`cd` 只会改变子进程的当前工作目录。父进程的工作目录保持原样。

现实情况

UNIX 将“标准”的文件描述符，管道，和便于操作它们的 shell 命令整合在一起，这是编写通用、可重用程序的重大进步。这个想法激发了 UNIX 强大和流行的“软件工具”文化，而且 shell 也是首个所谓的“脚本语言”。UNIX 的系统调用接口在今天仍然存在许多操作系统中，诸如 BSD，Linux，以及 Mac OS X。

现代内核提供了比 xv6 要多的多的系统调用和内核服务。最重要的一点，现代基于 Unix 的操作系统并不遵循早期 Unix 将设备暴露为特殊文件的设计，比如刚才所说的控制台文件。Unix 的作者继续打造 Plan 9 项目，它将“资源是文件”的概念应用到现代设备上，将网络、图形和其他资源都视作文件或者文件树。

文件系统抽象是一个强大的想法，它被以万维网的形式广泛的应用在互联网资源上。即使如此，也存在着其他的操作系统接口的模型。Multics，一个 Unix 的前辈，将文件抽象为一种类似内存的概念，产生了十分不同的系统接口。Multics 的设计的复杂性对 Unix 的设计者们产生了直接的影响，他们因此想把文件系统的设计做的更简单。

这本书考察 xv6 是如何实现类似 Unix 的接口的，但涉及的想法和概念可以运用到 Unix 之外很多地方上。任何一个操作系统都需要让多个进程复用硬件，实现进程之间的相互隔离，并提供进程间通讯的机制。在学习 xv6 之后，你应该了解一些其他的更加复杂的操作系统，看一下他们当中蕴含的 xv6 的概念。

第1章

第一个进程

本章通过第一个进程的创建来解释 xv6 是如何开始运行的，让我们得以一窥 xv6 提供的各个抽象是如何实现和交互的。xv6 尽量复用了普通操作的代码来建立第一个进程，避免单独为其撰写代码。接下来的各小节中，我们将详细探索其中的奥秘。

xv6 可以运行在搭载 Intel 80386 及其之后（即“x86”）处理器的 PC 上，因而许多底层功能（例如虚存的实现）是 x86 处理器专有的。本书假设读者已有些许在一些体系结构上进行机器级编程的经验。我们将在有关 x86 专有概念出现时，对其进行介绍。另外，附录 A 中简要地描述了 PC 平台的整体架构。

进程概览

进程是一个抽象概念，它让一个程序可以假设它独占一台机器。进程向程序提供“看上去”私有的，其他进程无法读写的内存系统（或地址空间），以及一颗“看上去”仅执行该程序的CPU。

xv6 使用页表（由硬件实现）来为每个进程提供其独有的地址空间。页表将虚拟地址（x86 指令所使用的地址）翻译（或说“映射”）为物理地址（处理器芯片向主存发送的地址）。

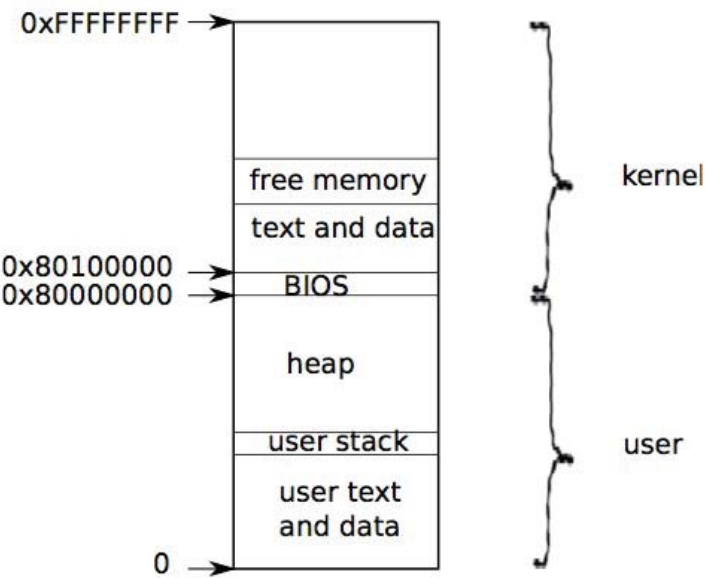


Figure 1-1. Layout of a virtual address space

xv6 为每个进程维护了不同的页表，这样就能够合理地定义进程的地址空间了。如图表1-1所示，一片地址空间包含了从虚拟地址0开始的用户内存。它的地址最低处放置进程的指令，接下来则是全局变量，栈区，以及一个用户可按需拓展的“堆”区（malloc 用）。

和上面提到的用户内存一样，内核的指令和数据也会被进程映射到每个进程的地址空间中。当进程使用系统调用时，系统调用实际上会在进程地址空间中的内核区域执行。这种设计使得内核的系统调用代码可以直接指向用户内存。为了给用户留下足够的内存空间，xv6 将内核映射到了地址空间的高地址处，即从 0x80100000 开始。

xv6 使用结构体 `struct proc` 来维护一个进程的状态，其中最为重要的状态是进程的页表，内核栈，当前运行状态。我们接下来会用 `p->xxx` 来指代 `proc` 结构中的元素。

每个进程都有一个运行线程（或简称为线程）来执行进程的指令。线程可以被暂时挂起，稍后再恢复运行。系统在进程之间切换实际上就是挂起当前运行的线程，恢复另一个进程的线程。线程的大多数状态（局部变量和函数调用的返回地址）都保

存在线程的栈上。

每个进程都有用户栈和内核栈（`p->kstack`）。当进程运行用户指令时，只有其用户栈被使用，其内核栈则是空的。然而当进程（通过系统调用或中断）进入内核时，内核代码就在进程的内核栈中执行；进程处于内核中时，其用户栈仍然保存着数据，只是暂时处于不活跃状态。进程的线程交替地使用着用户栈和内核栈。要注意内核栈是用户代码无法使用的，这样即使一个进程破坏了自己的用户栈，内核也能保持运行。

当进程使用系统调用时，处理器转入内核栈中，提升硬件的特权级，然后运行系统调用对应的内核代码。当系统调用完成时，又从内核空间回到用户空间：降低硬件特权级，转入用户栈，恢复执行系统调用指令后面的那条用户指令。线程可以在内核中“阻塞”，等待 I/O，在 I/O 结束后再恢复运行。

`p->state` 指示了进程的状态：新建、准备运行、运行、等待 I/O 或退出状态中。

`p->pgdir` 以 x86 硬件要求的格式保存了进程的页表。xv6 让分页硬件在进程运行时使用 `p->pgdir`。进程的页表还记录了保存进程内存的物理页的地址。

代码：第一个地址空间

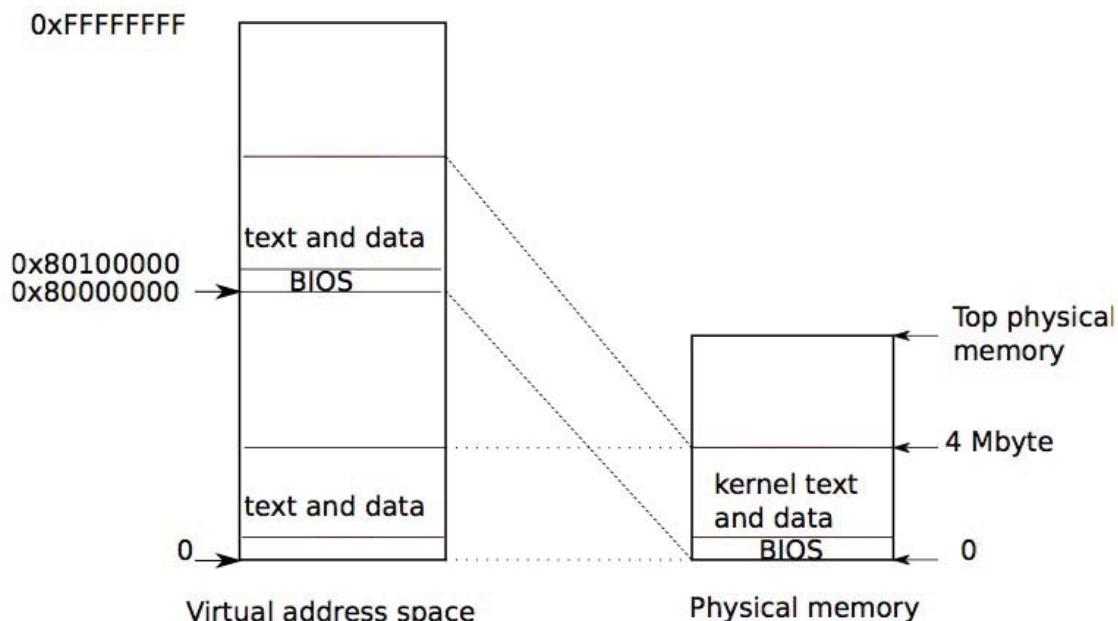


Figure 1-2. Layout of a virtual address space

当 PC 开机时，它会初始化自己然后从磁盘中载入 *boot loader* 到内存并运行。附录 B 介绍了其具体细节。然后，*boot loader* 把 xv6 内核从磁盘中载入并从 `entry`（1040）开始运行。x86 的分页硬件在此时还没有开始工作；所以这时的虚拟地址是直接映射到物理地址上的。

boot loader 把 xv6 内核装载到物理地址 0x100000 处。之所以没有装载到内核指令和内核数据应该出现的 0x80100000，是因为小型机器上很可能没有这么大的物理内存。而之所以在 0x100000 而不是 0x0 则是因为地址 0xa0000 到 0x100000 是属于 I/O 设备的。

为了让内核的剩余部分能够运行，`entry` 的代码设置了页表，将 0x80000000（称为 `KERNBASE`（0207））开始的虚拟地址映射到物理地址 0x0 处。注意，页表经常会这样把两段不同的虚拟内存映射到相同的一段物理内存，我们将会看到更多类似的例子。

`entry` 中的页表的定义在 `main.c`（1311）中。我们将在第 2 章讨论页表的细节，这里简单地说明一下，页表项 0 将虚拟地址 0:0x400000 映射到物理地址 0:0x400000。只要 `entry` 的代码还运行在内存的低地址处，我们就必须这样设置，但最后这个页表项是会被移除的。页表项 512（译注：原文中似乎误写为 960）将虚拟地址的 `KERNBASE:KERNBASE+0x400000`

映射到物理地址 0:0x400000。这个页表项将在 `entry` 的代码结束后被使用；它将内核指令和内核数据应该出现的高虚拟地址处映射到了 `boot loader` 实际将它们载入的低物理地址处。这个映射就限制内核的指令+代码必须在 4mb 以内。

让我们回到 `entry` 中继续页表的设置工作，它将 `entrypgdir` 的物理地址载入到控制寄存器 `%cr3` 中。分页硬件必须知道 `entrypgdir` 的物理地址，因为此时它还不知道如何翻译虚拟地址；它也还没有页表。`entrypgdir` 这个符号指向内存的高地址处，但只要用宏 `V2P_W0` (0220) 减去 `KERNBASE` 便可以找到其物理地址。为了让分页硬件运行起来，xv6 会设置控制寄存器 `%cr0` 中的标志位 `CR0_PG`。

现在 `entry` 就要跳转到内核的 C 代码，并在内存的高地址中执行它了。首先它将栈指针 `%esp` 指向被用作栈的一段内存 (1054)。所有的符号包括 `stack` 都在高地址，所以当低地址的映射被移除时，栈仍然是可用的。最后 `entry` 跳转到高地址的 `main` 代码中。我们必须使用间接跳转，否则编译器会生成 PC 相关的直接跳转 (PC-relative direct jump)，而该跳转会运行在内存低地址处的 `main`。`main` 不会返回，因为栈上并没有返回 PC 值。好了，现在内核已经运行在高地址处的函数 `main` (1217) 中了。

代码：创建第一个进程

在 `main` 初始化了一些设备和子系统后，它通过调用 `userinit` (1239) 建立了第一个进程。`userinit` 首先调用 `allocproc`。`allocproc` (2205) 的工作是在页表中分配一个槽 (即结构体 `struct proc`)，并初始化进程的状态，为其内核线程的运行做准备。注意一点：`userinit` 仅仅在创建第一个进程时被调用，而 `allocproc` 创建每个进程时都会被调用。`allocproc` 会在 `proc` 的表中找到一个标记为 `UNUSED` (2211-2213) 的槽位。当它找到这样一个未被使用的槽位后，`allocproc` 将其状态设置为 `EMBRYO`，使其被标记为被使用的并给这个进程一个独有的 `pid` (2201-2219)。接下来，它尝试为进程的内核线程分配内核栈。如果分配失败了，`allocproc` 会把这个槽位的状态恢复为 `UNUSED` 并返回 0 以标记失败。

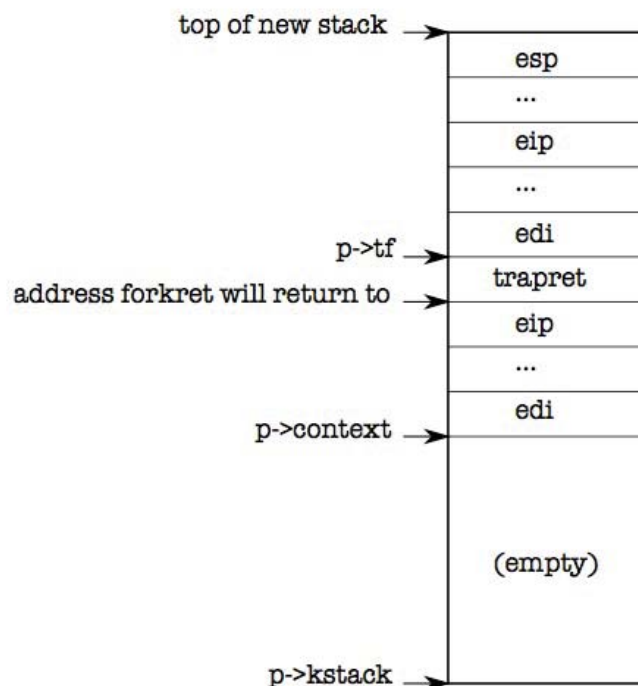


Figure 1-3. A new kernel stack.

现在 `allocproc` 必须设置新进程的内核栈，`allocproc` 以巧妙的方式，使其既能在创建第一个进程时被使用，又能在 `fork` 操作时被使用。`allocproc` 为新进程设置好一个特别准备的内核栈和一系列内核寄存器，使得进程第一次运行时“返回”到用户空间。准备好的内核栈就像图表 1-3 展示的那样。`allocproc` 通过设置返回程序计数器的值，使得新进程的内核线程首先运行在 `forkret` 的代码中，然后返回到 `trapret` (2236-2241) 中运行。

内核线程会从 `p->context` 中拷贝的内容开始运行。所以我们可以通过将 `p->context->eip` 指向 `forkret` 从而让内核线程从 `forkret` (2533) 的开头开始运行。这个函数会返回到那个时刻栈底的地址。`context switch` (2708) 的代码把栈指针指向 `p->context` 结尾。`allocproc` 又将 `p->context` 放在栈上，并在其上方放一个指向 `trapret` 的指针；这样运行完的 `forkret` 就会返回到 `trapret` 中了。`trapret` 接着从栈顶恢复用户寄存器然后跳转到 `process` (3027) 的代码。

这样的设置对于普通的 `fork` 和建立第一个进程都是适用的，虽然一种情况进程会从用户空间的地址0处开始执行而非真正的从 `fork` 返回。

我们将会在第3章看到，将控制权从用户转到内核是通过中断机制实现的，具体地说是系统调用、中断和异常。每当进程运行中要将控制权交给内核时，硬件和 xv6 的 `trap entry` 代码就会在进程的栈上保存用户寄存器。`userinit` 把值写在新建的栈的顶部，使之就像进程是通过中断进入内核的一样（2264-2270）。所以用于从内核返回到用户代码区的通用代码也能适用于第一个进程。这些保存的值就构成了一个结构体 `struct trapframe`，其中保存的是用户寄存器。现在如图表1-3所示，进程的栈已经准备好了。

第一个进程会先运行一个小程序（`initcode.S`（7700）），于是进程需要找到物理内存来保存这段程序。程序不仅需要被拷贝到内存中，还需要页表来指向那段内存。

`userinit` 调用 `setupkvm`（1737）来为进程创建一个（最初）只映射内核区的页表。我们将在第2章学习该函数的具体细节，概括地说，`setupkvm` 和 `userinit` 创建了图表1-1所示的地址空间。

第一个进程内存中的初始内容是由 `initcode.S` 汇编得到的；作为建立内核的进程的一部分，链接器将这段二进制代码嵌入内核中并定义两个特殊的符号：`_binary_initcode_start` 和 `_binary_initcode_size`，表示了这段二进制码的位置和大小。`userinit` 调用 `inituvm`，分配一页物理内存，将虚拟地址0映射到那一段内存，并把二进制码拷贝到那一页中（1803）。

接下来，`userinit` 把 `trap frame`（0602）设置为初始的用户模式状态：`%cs` 寄存器保存着一个段选择器，指向段 `SEG_UCODE` 并处于特权级 `DPL_USER`（即在用户模式而非内核模式）。类似的，`%ds`，`%es`，`%ss` 的段选择器指向段 `SEG_UDATA` 并处于特权级 `DPL_USER`。`%eflags` 的 `FL_IF` 位被设置为允许硬件中断；我们将在第3章回头看这段代码。

栈指针 `%esp` 被设为了进程的最大有效虚拟内存，即 `p->sz`。指令指针则指向初始化代码的入口点，即地址0。

函数 `userinit` 把 `p->name` 设置为 `initcode`，这主要是为了方便调试。还要将 `p->cwd` 设置为进程当前的工作目录；我们将在第6章回过头来看看 `namei` 的细节。

一旦进程初始化完毕，`userinit` 将 `p->state` 设置为 `RUNNABLE`，使进程能够被调度。

运行第一个进程

现在第一个进程的状态已经被设置好了，让我们来运行它。在 `main` 调用了 `userinit` 之后，`mpmain` 调用 `scheduler` 开始运行进程（1267）。`scheduler`（2458）会找到一个 `p->state` 为 `RUNNABLE` 的进程 `initproc`，然后将 `per-cpu` 的变量 `proc` 为该进程，接着调用 `switchuvm` 通知硬件开始使用目标进程的页表（1768）。注意，由于 `setupkvm` 使得所有的进程的页表都有一份相同的映射，指向内核的代码和数据，所以当内核运行时我们改变页表是没有问题的。`switchuvm` 同时还设置好任务状态段 `SEG_TSS`，让硬件在进程的栈中执行系统调用与中断。我们将在第3章研究任务状态段。

`scheduler` 接着把进程的 `p->state` 设置为 `RUNNING`，调用 `swtch`（2708），切换上下文到目标进程的内核线程中。`swtch` 会保存当前的寄存器，并把目标内核线程中保存的寄存器（`proc->context`）载入到 x86 的硬件寄存器中，其中也包括栈指针和指令指针。当前的上下文并非是进程的，而是一个特殊的 `per-cpu` 调度器的上下文。所以 `scheduler` 会让 `swtch` 把当前的硬件寄存器保存在 `per-cpu` 的存储（`cpu->scheduler`）中，而非进程的内核线程上下文中。我们将在第5章讨论 `swtch` 的细节。最后的 `ret`（2727）指令从栈中弹出目标进程的 `%eip`，从而结束上下文切换工作。现在处理器就运行在进程 `p` 的内核栈上了。

`allocproc` 通过把 `initproc` 的 `p->context->eip` 设置为 `forkret` 使得 `ret` 开始执行 `forkret` 的代码。第一次被使用（就是这一次）时，`forkret`（2533）会调用一些初始化函数。注意，我们不能在 `main` 中调用它们，因为它们必须在一个拥有自己的内核栈的普通进程中运行。接下来 `forkret` 返回。由于 `allocproc` 的设计，目前栈上在 `p->context` 之后即将被弹出的字是 `trapret`，因而接下来会运行 `trapret`，此时 `%esp` 保存着 `p->tf`。`trapret`（3027）用弹出指令从 `trap frame`（0602行）中恢复寄存器，就像 `swtch` 对内核上下文的操作一样：`popal` 恢复通用寄存器，`popl` 恢复 `%gs`，`%fs`，`%es`，`%ds`。`addl` 跳过 `trapno` 和 `errcode` 两个数据，最后 `iret` 弹出 `%cs`，`%eip`，`%flags`，`%esp`，`%ss`。`trap frame` 的内容已经转移到 CPU 状态中，所以处理器会从 `trap frame` 中 `%eip` 的值继续执行。对于 `initproc` 来说，这个值就是虚拟地址0，即 `initcode.S` 的第一个指令。

这时 `%eip` 和 `%esp` 的值为0和4096，这是进程地址空间中的虚拟地址。处理器的分页硬件会把它们翻译为物理地址。`allocvm` 为进程建立了页表，所以现在虚拟地址0会指向为该进程分配的物理地址处。`allocvm` 还会设置标志位 `PTE_U` 来让分页硬件允许用户代码访问内存。`userinit` 设置了 `%cs` 的低位，使得进程的用户代码运行在 `CPL = 3` 的情况下，这意味着用户代码只能使用带有 `PTE_U` 设置的页，而且无法修改像 `%cr3` 这样的敏感硬件寄存器。这样，处理器就受限只能使用自己的内存了。

第一个系统调用：`exec`

`initcode.S` 干的第一件事是触发 `exec` 系统调用。就像我们在第0章看到的一样，`exec` 用一个新的程序来代替当前进程的内存和寄存器，但是其文件描述符、进程 `id` 和父进程都是不变的。

`initcode.S` (7708) 刚开始会将 `$argv`, `$init`, `$0` 三个值推入栈中，接下来把 `%eax` 设置为 `SYS_exec` 然后执行 `int T_SYSCALL`：这样做是告诉内核运行 `exec` 这个系统调用。如果运行正常的话，`exec` 不会返回：它会运行名为 `$init` 的程序，`$init` 是一个以空字符结尾的字符串，即 `/init` (7721-7723)。如果 `exec` 失败并且返回了，`initcode` 会不断调用一个不会返回的系统调用 `exit`。

系统调用 `exec` 的参数是 `$init`, `$argv`。最后的 `0` 让这个手动构建的系统调用看起来就像普通的系统调用一样，我们会在第3章详细讨论这个问题。和之前的代码一样，xv6 努力避免为第一个进程的运行单独写一段代码，而是尽量使用通用于普通操作的代码。

第2章讲了 `exec` 的具体实现，概括地讲，它会用从文件系统中获取的 `/init` 的二进制代码代替 `initcode` 的代码。现在 `initcode` 已经执行完了，进程将要运行 `/init`。`init` (7810行) 会在需要的情况下创建一个新的控制台设备文件，然后把它作为描述符0, 1, 2打开。接下来它将不断循环，开启控制台 shell，处理没有父进程的僵尸进程，直到 shell 退出，然后再反复。系统就这样建立起来了。

现实情况

大多操作系统都采用了进程这个概念，而大多的进程都和 xv6 的进程类似。但是真正的操作系统会利用一个显式的链表在常数时间内找到空闲的 `proc`，而不像 `allocproc` 中那样花费线性时间；xv6 使用的是朴素的线性搜索（找第一个空闲的 `proc`）。

xv6 的地址空间结构有一个缺点，即无法使用超过 2GB 的物理 RAM。当然我们可以解决这个问题，不过最好的解决方法还是使用64位的机器。

练习

1. 在 `swtch` 中设断点。用 `gdb` 的 `stepi` 单步调试返回到 `forkret` 的代码，然后使用 `gdb` 的 `finish` 继续执行到 `trapret`，然后再用 `stepi` 直到你进入虚拟地址0处的 `initcode`。
2. `KERNBASE` 会限制一个进程能使用的内存量，在一台有着 4GB 内存的机器上，这可能会让人感到不悦。那么提高 `KERNBASE` 的值是否能让进程使用更多的内存呢？

第2章

页表

操作系统通过页表机制实现了对内存空间的控制。页表使得 xv6 能够让不同进程各自的地址空间映射到相同的物理内存上，还能够为不同进程的内存提供保护。除此之外，我们还能够通过使用页表来间接地实现一些特殊功能。xv6 主要利用页表来区分多个地址空间，保护内存。另外，它也使用了一些简单的技巧，即把不同地址空间的多段内存映射到同一段物理内存（内核部分），在同一地址空间中多次映射同一段物理内存（用户部分的每一页都会映射到内核部分），以及通过一个没有映射的页保护用户栈。本章的其余部分将详细地探讨 x86 硬件提供的页表以及 xv6 对页表的使用。

分页硬件

回顾一下，x86 的指令（用户和内核均是如此）计算的都是虚拟地址。机器的 RAM，或者物理内存，则是用物理地址来作标记的。x86 的页表硬件通过映射机制将虚拟地址和物理地址联系起来。

一个 x86 页表就是一个包含 2^{20} (1,048,576) 条页表条目 (PTE) 的数组。每条 PTE 包含了一个 20 位的物理页号 (PPN) 及一些标志位。分页硬件要找到一个虚拟地址对应的 PTE，只需使用其高 20 位来找到该虚拟地址在页表中的索引，然后将其高 20 位替换为对应 PTE 的 PPN。而低 12 位是会被分页硬件原样复制的。因此在虚拟地址-物理地址的翻译机制下，页表可以为操作系统提供对一块块大小为 4096 (2^{12}) 字节的内存片，这样的一个内存片就是一页。

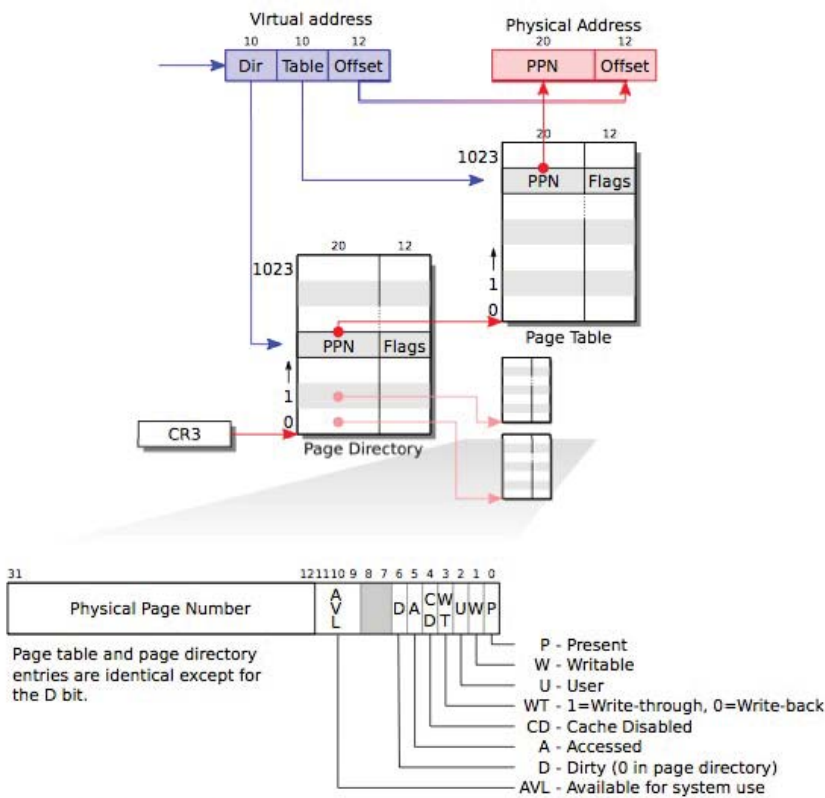


Figure 2-1. x86 page table hardware.

如图 2-1 所示，实际上，地址的翻译有两个步骤。一个页表在物理内存中像一棵两层的树。树的根是一个 4096 字节的项目录，其中包含了 1024 个类似 PTE 的条目，但其实每个条目是指向一个页表页的引用。而每个页表页又是包含 1024 个 32 位 PTE 的数组。分页硬件使用虚拟地址的高 10 位来决定对应项目录条目。如果想要的条目已经放在了项目录中，分页硬件就会继续使用接下来的 10 位来从页表页中选择出对应的 PTE。否则，分页硬件就会抛出错误。通常情况下，大部分虚拟地址不会进行映射，而这样的二级结构就使得项目录可以忽略那些没有任何映射的页表页。

每个 PTE 都包含一些标志位，说明分页硬件对应的虚拟地址的使用权限。PTE_P 表示 PTE 是否陈列在页表中：如果不是，那么一个对该页的引用会引发错误（也就是：不允许被使用）。PTE_W 控制着能否对页执行写操作；如果不能，则只允许对其进行读操作和取指令。PTE_U 控制着用户程序能否使用该页；如果不能，则只有内核能够使用该页。图 2-1 对此进行了说明。这些的标志位和页表硬件相关的结构体都在 `mmu.h`（0200）定义。

下面对一些名词作出解释。物理内存是指 DRAM 中的储存单元。每个字节的物理内存都有一个地址，称为物理地址。而虚拟地址则是程序所使用的。分页硬件会将程序发出的虚拟地址翻译为物理地址，然后发送给 DRAM 硬件以读写存储器。这一层面的讨论中我们仅仅考虑虚拟地址，暂不考虑虚拟内存。

进程地址空间

`entry` 中建立的页表已经产生了足够多的映射来让内核的 C 代码正常运行。但是 `main` 还是调用了 `kvmalloc`（1757）立即转换到新的页表中，这是因为内核建立的页表更加精巧地映射了内存空间。

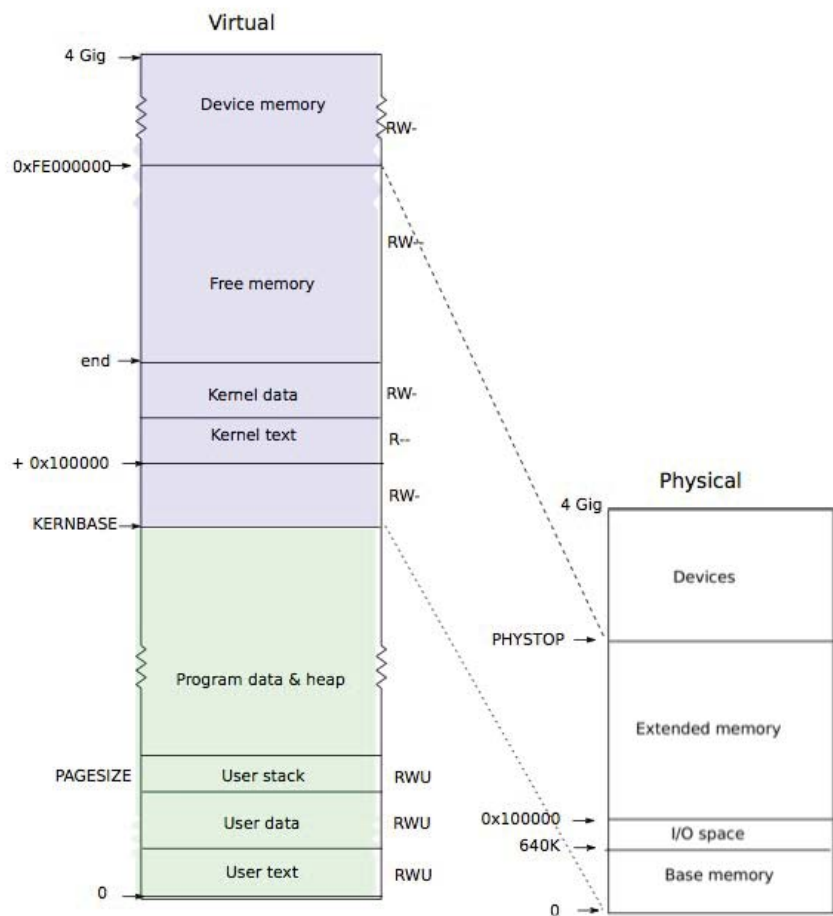


Figure 2-2. Layout of a virtual address space and the physical address space.

每个进程都有自己的页表，xv6 会在进程切换时通知分页硬件切换页表。如图表 2-2 所示，进程的用户内存从 0 开始，最多能够增长到 `KERNBASE`，这使得一个进程最多只能使用 2GB 的内存。当进程向 xv6 要求更多的内存时，xv6 首先要找到空闲的物理页，然后把这些页对应的 PTE 加入该进程的页表中，并让 PTE 指向对应的物理页。xv6 设置了 PTE 中的 `PTE_U`、`PTE_W`、`PTE_P` 标志位。大多数进程是用不完整个内存空间的；xv6 会把没有被使用的 PTE 的 `PTE_P` 标志位设为 0。不同进程的页表将其用户内存映射到不同的物理内存中，因此每个进程就拥有了私有的用户内存。

xv6 在每个进程的页表中都包含了内核运行所需要的所有映射，而这些映射都出现在 `KERNBASE` 之上。它将虚拟地址 `KERNBASE:KERNBASE+PHYSTOP` 映射到 `0:PHYSTOP`。这样映射的原因之一是内核可以使用自己的指令和数据；原因之二是内核有时需要对物理页进行写操作，譬如在创建页表页的时候，而使得每一个物理页都在对应的虚拟地址上被映射就让这些操作变得很方便。这样的安排有一个缺点，即 xv6 无法使用超过 2GB 的物理内存。有一些使用内存映射的 I/O 设备的物理内存存在 `0xFE000000` 之上，对于这些设备 xv6 页表采用了直接映射。`KERNBASE` 之上的页对应的 PTE 中，`PTE_U` 位均被置 0，因而

只有内核能够使用这些页。

每个进程的页表同时包括用户内存和内核内存的映射，这样当用户通过中断或者系统调用转入内核时就不需要进行页表的转换了。大多数情况下，内核都没有自己的页表，所以内核几乎都是在借用用户进程的页表。

现在来回顾一下，xv6 保证了每个进程只能使用其自己的内存，并且每个进程所看到的内存都是从虚拟地址 0 开始的一段连续内存。对于一个进程，xv6 只把该进程所使用的内存对应的 PTE 的 `PTE_U` 设为 1，其他 PTE 则不然，这样就可以实现前者。对于后者，则是让页表把连续的虚拟页映射到实际分配的物理页。

代码：建立一个地址空间

`main` 调用 `kvmalloc` (1757)，创建并切换到一个拥有内核运行所需的 `KERNBASE` 以上映射的页表。这里的大多数工作都是由 `setupkvm` (1737) 完成的。首先，它会分配一页内存来放置页目录，然后调用 `mappages` 来建立内核需要的映射，这些映射可以在 `kmap` (1728) 数组中找到。这里的映射包括内核的指令和数据，`PHYSTOP` 以下的物理内存，以及 I/O 设备所占的内存。`setupkvm` 不会建立任何用户内存的映射，这些映射稍后会建立。

`mappages` (1679) 做的工作是在页表中建立一段虚拟内存到一段物理内存的映射。它是在页的级别，即一页一页地建立映射的。对于每一个待映射虚拟地址，`mappages` 调用 `walkpgdir` 来找到该地址对应的 PTE 地址。然后初始化该 PTE 以保存对应物理页号、许可级别 (`PTE_W` 和/或 `PTE_U`) 以及 `PTE_P` 位来标记该 PTE 是否是有效的 (1691)。

`walkpgdir` (1654) 模仿 x86 的分页硬件为一个虚拟地址寻找 PTE 的过程 (见图表2-1)。`walkpgdir` 通过虚拟地址的前 10 位来找到在页目录中的对应条目 (1659)，如果该条目不存在，说明要找的页表页尚未分配；如果 `alloc` 参数被设置了，`walkpgdir` 会分配页表页并将其物理地址放到页目录中。最后用虚拟地址的接下来 10 位来找到其在页表中的 PTE 地址 (1672)。

物理内存的分配

在运行时，内核需要为页表、进程的用户内存、内核栈及管道缓冲区分配空闲的物理内存。

xv6 使用从内核结尾到 `PHYSTOP` 之间的物理内存为运行时分配提供内存资源。每次分配，它会将整块 4096 字节大小的页分配出去。xv6 还会通过维护一个物理页组成的链表来寻找空闲页。所以，分配内存需要将页移出该链表，而释放内存需要将页加入该链表。

这里我们遇到了一个自举的问题：为了让分配器能够初始化该空闲链表，所有的物理内存都必须建立起映射，但是建立包含这些映射的页表又必须要分配存放页表的页。xv6 通过在 `entry` 中使用一个特别的页分配器来解决这个问题。该分配器会在内核数据部分的后面分配内存。该分配器不支持释放内存，并受限于 `entrypgdir` 中规定的 4MB 分配大小。即便如此，该分配器还是足够为内核的第一个页表分配出内存。

代码：物理内存分配器

分配器中的数据结构是一个由可分配物理内存页构成的空闲链表。这个空闲页的链表的元素是结构体 `struct run` (2764)。那么分配器从哪里获得内存来存放这些数据呢？实际上，分配器将每个空闲页的 `run` 结构体保存在该空闲页本身中，因为空闲页中没有其他数据。分配器还用了一个 spin lock (2764-2766) 来保护空闲链表。链表和这个锁都封装在一个结构体中，这样逻辑就比较明晰：锁保护了该结构体中的域。不过现在让我们先忽略这个锁，以及对 `acquire` 和 `release` 的调用；我们会在第 4 章了解其细节。

`main` 函数调用了 `kinit1` 和 `kinit2` 两个函数对分配器进行初始化 (2780)。这样做是由于 `main` 中的大部分代码都不能使用锁以及 4MB 以上的内存。`kinit1` 在前 4MB 进行了不需要锁的内存分配。而 `kinit2` 允许了锁的使用，并使得更多的内存可用于分配。原本应该由 `main` 决定有多少物理内存可用于分配，但在 x86 上很难实现。所以它假设机器中有 240MB (`PHYSTOP`) 物理内存，并将内核末尾和 `PHYSTOP` 之间的内存都作为一个初始的空闲内存池。`kinit1` 和 `kinit2` 调用 `freerange` 将内存加入空闲链表中，`freerange` 则是通过对每一页调用 `kfree` 实现该功能。一个 PTE 只能指向一个 4096 字节对齐的物理地址 (即是 4096 的倍数)，因此 `freerange` 用 `PGROUNDUP` 来保证分配器只会释放对齐的物理地址。分配器原本一开始没有内存可用，正是对 `kfree` 的调用将可用内存交给了分配器来管理。

分配器用映射到高内存区域的虚拟地址找到对应的物理页，而非物理地址。所以 `kinit` 会使用 `p2v(PHYSTOP)` 来将 `PHYSTOP`（一个物理地址）翻译为虚拟地址。分配器有时将地址看作是整型，这是为了对其进行运算（譬如在 `kinit` 中遍历所有页）；而有时将地址看作读写内存用的指针（譬如操作每个页中的 `run` 结构体）；对地址的双重使用导致分配器代码中充满了类型转换。另外一个原因是，释放和分配内存隐性地改变了内存的类型。

函数 `kfree`（2815）首先将被释放内存的每一字节设为 1。这使得访问已被释放内存的代码所读到的不是原有数据，而是垃圾数据；这样做能让这种错误的代码尽早崩溃。接下来 `kfree` 把 `v` 转换为一个指向结构体 `struct run` 的指针，在 `r->next` 中保存原有空闲链表的表头，然后将当前的空闲链表设置为 `r`。`kalloc` 移除并返回空闲链表的表头。

地址空间中的用户部分

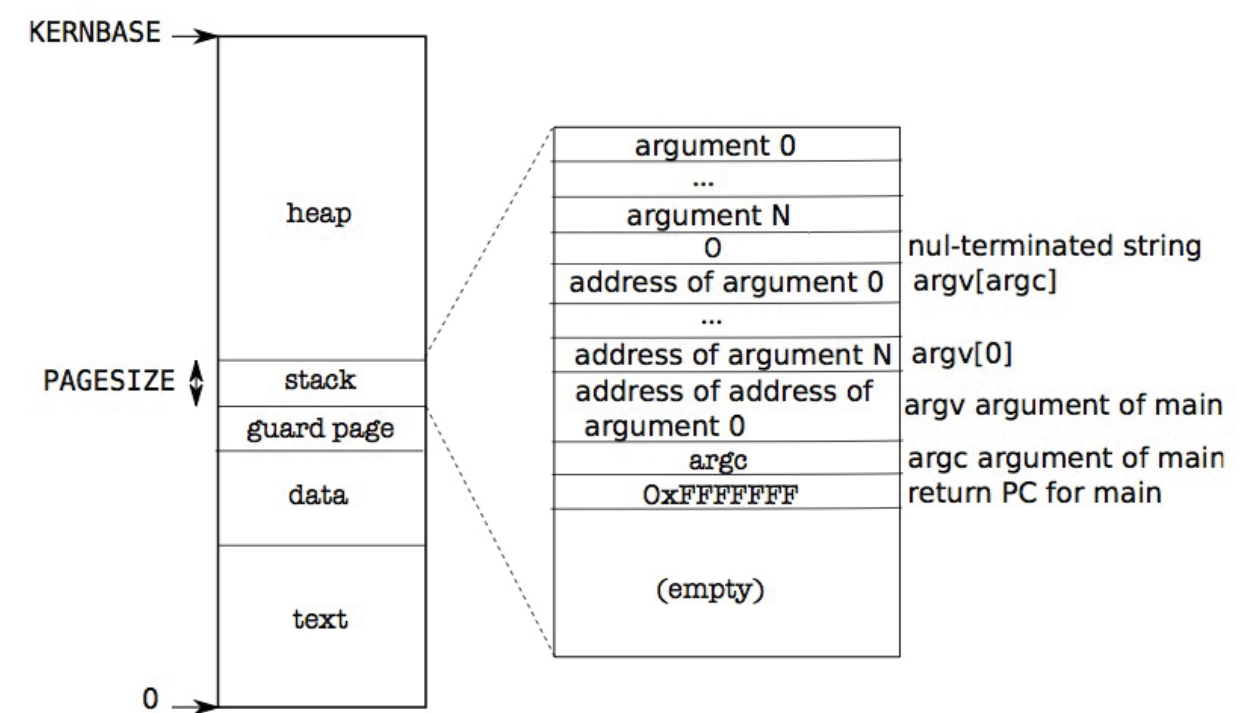


Figure 2-3. Memory layout of a user process with its initial stack.

图表 2-3 展示了在 `xv6` 中，一个运行中进程的用户内存结构。堆在栈之上，所以它可以增长（通过 `sbrk`）。栈占用了单独的一页内存，其中存有 `exec` 创建的初始数据。栈的最上方放着字符串形式的命令行参数以及指向这些参数的指针数组，其下方放的值使得一个程序可以从 `main` 开始，仿佛刚刚调用了函数 `main(argc, argv)`。为了防止栈使用了它不应该使用的页，栈的下方有一个保护页。保护页没有映射，因此当栈的增长超出其所在页时就会产生异常，因为无法翻译这个错误的地址。

代码：`exec`

`exec` 是创建地址空间中用户部分的系统调用。它根据文件系统中保存的某个文件来初始化用户部分。`exec`（5910）通过 `namei`（5920）打开二进制文件，这一点将在第 6 章进行解释。然后，它读取 ELF 头。`xv6` 应用程序以通行的 ELF 格式来描述，该格式在 `elf.h` 中定义。一个 ELF 二进制文件包括了一个 ELF 头，即结构体 `struct elfhdr`（0955），然后是连续几个程序段的头，即结构体 `struct proghdr`（0974）。每个 `proghdr` 都描述了需要载入到内存中的程序段。`xv6` 中的程序只有一个程序段的头，但其他操作系统中可能有多。

`exec` 第一步是检查文件是否包含 ELF 二进制代码。一个 ELF 二进制文件是以 4 个“魔法数字”开头的，即 `0x7F`，“E”，“L”，“F”，或者写为宏 `ELF_MAGIC`（0952）。如果 ELF 头中包含正确的魔法数字，`exec` 就会认为该二进制文件的结构是正确的。

`exec` 通过 `setupkvm`（5931）分配了一个没有用户部分映射的页表，再通过 `allocuvmm`（5943）为每个 ELF 段分配内存，然后通过 `loaduvmm`（5945）把段的内容载入内存中。`allocuvmm` 会检查请求分配的虚拟地址是否是在 `KERNBASE` 之下。

`loaduvm` (1818) 通过 `walkpgdir` 来找到写入 ELF 段的内存的物理地址；通过 `readi` 来将段的内容从文件中读出。

`exec` 创建的第一个用户程序 `/init` 程序段的头是这样的：

```
#objdump -p _init

_init:      file format elf32-i386

Program Header:
  LOAD off   0x00000054 vaddr 0x00000000 paddr 0x00000000 align 2**2
        filesz 0x000008c0 memsz 0x000008cc flags
```

程序段头中的 `filesz` 可能比 `memsz` 小，这表示中间相差的地方应该用 0 填充（对于 C 的全局变量）而不是继续从文件中读取数据。对于 `/init`，`filesz` 是 2240 字节而 `memsz` 是 2252 字节。所以 `allocuvm` 会分配足够的内存来装 2252 字节的内容，但只从文件 `/init` 中读取 2240 字节的内容。

现在 `exec` 要分配以及初始化用户栈了。它只为栈分配一页内存。`exec` 一次性把参数字符串拷贝到栈顶，然后把指向它们的指针保存在 `ustack` 中。它还会在 `main` 参数列表 `argv` 的最后放一个空指针。这样，`ustack` 中的前三项就是伪造的返回 PC，`argc` 和 `argv` 指针了。

`exec` 会在栈的页下方放一个无法访问的页，这样当程序尝试使用超过一个页的栈时就会出错。另外，这个无法访问的页也让 `exec` 能够处理那些过于庞大的参数；当参数过于庞大时，`exec` 用于将参数拷贝到栈上的函数 `copyout` 会发现目标页无法访问，并且返回 -1。

在创建新的内存映像时，如果 `exec` 发现了错误，比如一个无效的程序段，它就会跳转到标记 `bad` 处，释放这段内存映像，然后返回 -1。`exec` 必须在确认该调用可以成功后才能释放原来的内存映像，否则，若原来的内存映像被释放，`exec` 甚至都无法向它返回 -1 了。`exec` 中的错误只可能发生在新建内存映像时。一旦新的内存映像建立完成，`exec` 就能装载新映像（5989）而把旧映像释放（5990）。最后，`exec` 成功地返回 0。

现实情况

和大多数操作系统一样，xv6 使用分页硬件来保护和映射内存。但是很多操作系统的实现更加精巧；例如，xv6 不能向磁盘中请求页，没有实现 copy-on-write 的 fork 操作、共享内存和惰性分配页（lazily-allocated page），也不能自动扩展栈。x86 支持段式内存转换（见附录 B），但 xv6 仅用它来实现 `proc` 这种有固定地址，但在不同 CPU 上有不同值的 per-CPU 变量（见 `seginit`）。对于不支持段式内存的体系结构而言，想要实现这种 per-CPU（或 per-thread）的变量，就必须额外用一个寄存器来保存指向 per-CPU 数据区的指针。由于 x86 的寄存器实在太少，所以花费额外代价用段式内存来实现 per-CPU 变量是值得的。

在内存较多的机器上使用 x86 的 4MB 大小的“超级页”还是很划算的，能够减少页表的工作负担。而内存较小时，就比较适合用比较小的页，使得分配和向磁盘换出页时都拥有较细的粒度。譬如，当一个程序只需 8KB 的内存时，分配 4MB 的页就太浪费了。xv6 只在初始页表（1311）中使用了“超级页”。数组的初始化设置了 1024 条 PDE 中的 2 条，即 0 号和 512 号（`KERBASE >> PDXSHIFT`），而其他的 PDE 均为 0。xv6 设置了这两条 PDE 中的 `PTE_PS` 位，标记它们为“超级页”。内核还通过设置 `%cr4` 中的 `CP_PSE`（Page Size Extension）位来通知分页硬件允许使用超级页。

xv6 本来应该确定实际 RAM 的配置，而不是假设有 240MB 内存。在 x86 上，至少有三个通用算法：第一种是探测物理地址空间，寻找像内存一样能够维持被写入数据的区域；第二种是从 PC 非易失性 RAM 中某个已知的 16-bit 位置处读取内存大小；第三种是在 BIOS 中查看作为多处理器表一部分的内存布局表。读取内存布局表是一项比较复杂的工作。

内存分配曾经是一个热门话题，其主要问题就是如何对有限内存进行高效使用，以及如何为可能出现的各种内存请求做好准备；相关资料请搜索 Knuth。今天人们更加关心速度，而非空间利用率。另外，精巧的内核往往会分配不同大小的内存块，而不是像 xv6 一样固定分配 4096 字节；实际使用的内存分配器必须做到，对小块内存和大块内存的分配请求都能很好地处理。

练习

1. 查看一下真实的操作系统的内存大小。
2. 如果 xv6 没有使用超级页，我们应该如何声明 `entrypgdir` ？
3. Unix 在 `exec` 的实现中考虑了对 shell 脚本的特殊处理。执行一个以文本 `#!` 开头的文件时，第一行内容会被系统理解为执行该文件的翻译器。例如，如果用 `exec` 执行 `myprog arg1`，而 `myprog` 的第一行又是 `#!/interp`，那么 `exec` 会运行命令 `/interp myprog arg1`。如何 xv6 中实现该功能？

第3章

陷入，中断和驱动程序

运行进程时，cpu 一直处于一个大循环中：取指，更新 PC，执行，取指……。但有些情况下用户程序需要进入内核，而不是执行下一条用户指令。这些情况包括设备信号的发出、用户程序的非法操作（例如引用一个找不到页表项的虚拟地址）。处理这些情况面临三大挑战：1) 内核必须使处理器能够从用户态转换到内核态（并且再转换回用户态）2) 内核和设备必须协调好他们并行的活动。3) 内核必须知道硬件接口的细节。解决这三个问题需要对硬件的深入理解和小心翼翼的编程，并且有可能导致难以理解的内核代码。这一章告诉你 xv6 是如何解决这些问题的。

系统调用，异常和中断

正如我们上一章最后所见，用户程序通过系统调用请求系统服务。术语 `exception` 指产生中断的非法程序操作，例如除以0，尝试访问 PTE 不存在的内存等等。术语 `interrupt` 指硬件产生的希望引起操作系统注意的信号，例如时钟芯片可能每100毫秒产生一个中断，以此来实现分时。再举一个例子，当硬盘读完一个数据块时，它会产生一个中断来提醒操作系统这个块已经准备好被获取了。

所有的中断都由内核管理，而不是进程。因为在大多数情况下只有内核拥有处理中断所需的特权和状态。例如为了使进程响应时钟中断而在进程间实现时间分片，就必须在内核中执行这些操作，因为我们有可能强迫进程服从处理器的调度。

在所有三种情况下，操作系统的设计必须保证下面这些事情。系统必须保存寄存器以备将来的状态恢复。系统必须准备好在内核中执行，必须选择一个内核开始执行的地方。内核必须能够获得关于这个事件的信息，例如系统调用的参数。同时还必须保证安全性；系统必须保持用户进程和系统进程的隔离。

为了达成这个目标操作系统必须知道硬件是如何处理系统调用、异常和中断的。在大多数处理器中这三种事件都用同样的硬件机制处理。比如说，在 x86 中，一个程序可以通过 `int` 指令产生一个中断来进行系统调用。同样的，异常也会产生一个中断。因此，如果操作系统能够处理中断，那么操作系统也可以处理系统调用以及异常。

我们的计划是这样的。中断终止正常的处理器循环然后开始执行中断处理程序中的代码。在开始中断处理程序之前，处理器保存寄存器，这样在操作系统从中断中返回时就可以恢复他们。切换到中断服务程序面临的问题是处理器需要在用户模式和内核模式之间切换。

咱们说说术语：虽然官方的 x86 术语是中断，xv6 都用陷入来代表他们，很大程度上是因为这个术语被 PDP11/40 使用，从而也是传统的 Unix 术语。这一章交替使用陷入和中断这两个术语，但一定要记住陷入是由在 cpu 上运行的当前进程导致的，而中断是由设备导致的，可能与当前进程毫无关系。比如说，磁盘可能在接受了一个进程的数据块之后发出一个中断，但是在中断的时候可能运行的是其他进程。中断的这一特性使得思考中断的相关问题比陷入要难，因为中断和其它活动是并行的。然而正如我们马上就要讨论的，他们都依赖相同的硬件机制在用户模式和内核模式之间进行切换。

X86 的保护机制

x86 有四个特权级，从 0（特权最高）编号到 3（特权最低）。在实际使用中，大多数的操作系统都使用两个特权级，0 和 3，他们被称为内核模式和用户模式。当前执行指令的特权级存在于 `%cs` 寄存器中的 CPL 域中。

在 x86 中，中断处理程序的入口在中断描述符表（IDT）中被定义。这个表有256个表项，每一个都提供了相应的 `%cs` 和 `%eip`。

一个程序要在 x86 上进行一个系统调用，它需要调用 `int n` 指令，这里 `n` 就是 IDT 的索引。`int` 指令进行下面一些步骤：

- 从 IDT 中获得第 `n` 个描述符，`n` 就是 `int` 的参数。
- 检查 `%cs` 的域 `CPL ≤ DPL`，`DPL` 是描述符中记录的特权级。
- 如果目标段选择符的 `PL < CPL`，就在 CPU 内部的寄存器中保存 `%esp` 和 `%ss` 的值。
- 从一个任务段描述符中加载 `%ss` 和 `%esp`。

- 将 %ss 压栈。
- 将 %esp 压栈。
- 将 %eflags 压栈。
- 将 %cs 压栈。
- 将 %eip 压栈。
- 清除 %eflags 的一些位。
- 设置 %cs 和 %eip 为描述符中的值。

int 指令是一个非常复杂的指令，可能有人会问是不是所有的这些操作都是必要的。检查 $CPL \leq DPL$ 使得内核可以禁止一些特权级系统调用。例如，如果用户成功执行了 int 指令，那么 DPL 必须是 3。如果用户程序没有合适的特权级，那么 int 指令就会触发 int 13，这是一个通用保护错误。再举一个例子，int 指令不能使用用户栈来保存值，因为用户可能还没有建立一个合适的栈，因此硬件会使用任务段中指定的栈（这个栈在内核模式中建立）。

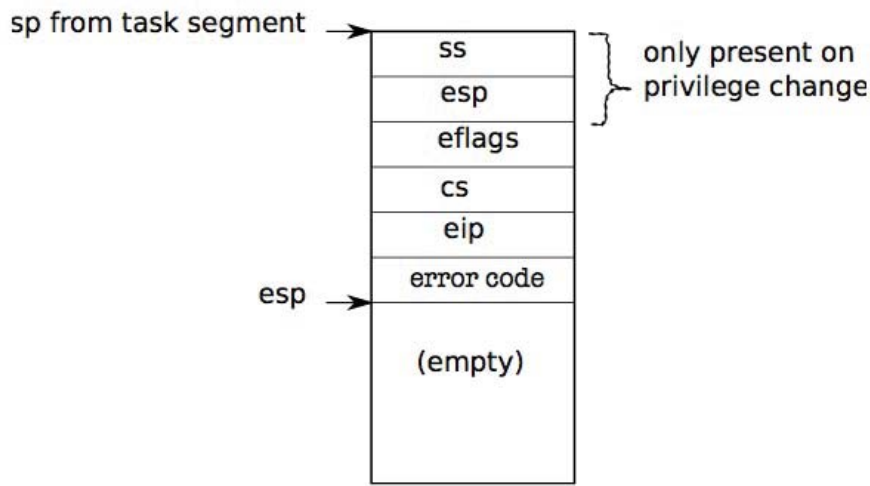


Figure 3-1. Kernel stack after an int instruction.

图 3-1 展示了一个 int 指令之后的栈的情况，注意这是发生了特权级转换（即描述符中的特权级比 CPL 中的特权级低的时候）栈的情况。如果这条指令没有导致特权级转换，x86 就不会保存 %ss 和 %esp。在任何一种情况下，%eip 都指向中断描述符表中指定的地址，这个地址的第一条指令就是要执行的下一条指令，也是 int n 的中断处理程序的第一条指令。操作系统应该实现这些中断处理程序，之后我们会看到 xv6 干了些什么。

操作系统可以使用 `iret` 指令来从一个 int 指令中返回。它从栈中弹出 int 指令保存的值，然后通过恢复保存的 %eip 的值来继续用户程序的执行。

代码：第一个系统调用

第一章的最后在 `initcode.S` 中调用了系统调用。让我们再看一遍 (7713)。这个进程将 `exec` 所需的参数压栈，然后把系统调用号存在 %eax 中。这个系统调用号和 `syscalls` 数组中的条目匹配，（`syscall` 是一个函数指针的数组）(3350)。我们需要设法使得 int 指令将处理器的状态从用户模式切换到内核模式，调用适当的内核函数（例如在这里是 `sys_exec`），并且使内核可以取出 `sys_exec` 的参数。接下来的几个小节将描述 xv6 是如何做到这一点的，你会发现我们可以用同样的代码来实现中断和异常。

代码：汇编陷入处理程序

xv6 必须设置硬件在遇到 int 指令时进行一些特殊的操作，这些操作会使处理器产生一个中断。x86 允许 256 个不同的中断。中断 0-31 被定义为软件异常，比如除 0 错误和访问非法的内存页。xv6 将中断号 32-63 映射给硬件中断，并且用 64 作为系统调用的中断号。

`tvinit` (3067) 在 `main` 中被调用，它设置了 `idt` 表中的 256 个表项。中断 `i` 被位于 `vectors[i]` 的代码处理。每一个中断处理程序的入口点都是不同的，因为 x86 并未把中断号传递给中断处理程序，使用 256 个不同的处理程序是区分这 256 种

情况的唯一办法。

`Tvinit` 处理 `T_SYSCALL`，用户系统会调用 `trap`，特别地：它通过传递第二个参数值为 1 来指定这是一个陷阱门。陷阱门不会清除 `FL` 位，这使得在处理系统调用的时候也接受其他中断。

同时也设置系统调用门的权限为 `DPL_USER`，这使得用户程序可以通过 `int` 指令产生一个内陷。xv6 不允许进程用 `int` 来产生其他中断（比如设备中断）；如果它们这么做了，就会抛出通用保护异常，也就是发出 13 号中断。

当特权级从用户模式向内核模式转换时，内核不能使用用户的栈，因为它可能不是有效的。用户进程可能是恶意的或者包含了一些错误，使得用户的 `%esp` 指向一个不是用户内存的地方。xv6 会使得在内陷发生的时候进行一个栈切换，栈切换的方法是让硬件从一个任务段描述符中读出新的栈选择符和一个新的 `%esp` 的值。函数 `switchuvm` (1773) 把用户进程的内核栈顶地址存入任务段描述符中。

当内陷发生时，处理器会做下面一些事。如果处理器在用户模式下运行，它会从任务段描述符中加载 `%esp` 和 `%ss`，把老的 `%ss` 和 `%esp` 压入新的栈中。如果处理器在内核模式下运行，上面的事件就不会发生。处理器接下来会把 `%eflags`，`%cs`，`%eip` 压栈。对于某些内陷来说，处理器会压入一个错误字。而后，处理器从相应 IDT 表项中加载新的 `%eip` 和 `%cs`。

xv6 使用一个 perl 脚本 (2950) 来产生 IDT 表项指向的中断处理函数入口点。每一个入口都会压入一个错误码（如果 CPU 没有压入的话），压入中断号，然后跳转到 `alltraps`。

`Alltraps` (3004) 继续保存处理器的寄存器：它压入 `%ds`，`%es`，`%fs`，`%gs`，以及通用寄存器 (3005-3010)。这么做使得内核栈上压入一个 `trapframe`（中断帧）结构体，这个结构体包含了中断发生时处理器的寄存器状态（参见图 3-2）。处理器负责压入 `%ss`，`%esp`，`%eflags`，`%cs` 和 `%eip`。处理器或者中断入口会压入一个错误码，而 `alltraps` 负责压入剩余的。中断帧包含了所有处理器从当前进程的内核态恢复到用户态需要的信息，所以处理器可以恰如中断开始时那样继续执行。回顾一下第二章，`userinit` 通过手动建立中断帧来达到这个目标（参见图 1-3）。

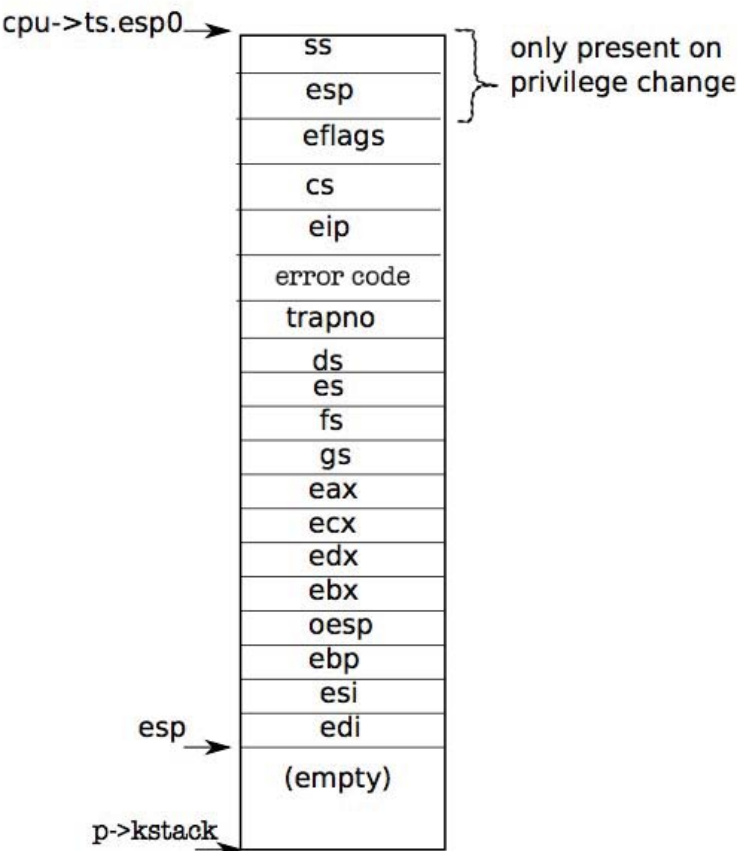


Figure 3-2. The trapframe on the kernel stack

考虑第一个系统调用，被保存的 `%eip` 是 `int` 指令下一条指令的地址。`%cs` 是用户代码段选择符。`%eflags` 是执行 `int` 指令时

的 `eflags` 寄存器，`alltraps` 同时也保存 `%eax`，它存有系统调用号，内核在之后会使用到它。

现在用户态的寄存器都保存了，`alltraps` 可以完成对处理器的设置并开始执行内核的 C 代码。处理器在进入中断处理程序之前设置选择符 `%cs` 和 `%ss`；`alltraps` 设置 `%ds` 和 `%es`（3013-3015）。它设置 `%fs` 和 `%gs` 来指向 `SEG_KCPU`（每个 CPU 数据段选择符）（3016-3018）。

一旦段设置好了，`alltraps` 就可以调用 C 中断处理程序 `trap` 了。它压入 `%esp` 作为 `trap` 的参数，`%esp` 指向刚在栈上建立好的中断帧（3021）。然后它调用 `trap`（3022）。`trap` 返回后，`alltraps` 弹出栈上的参数（3023）然后执行标号为 `trapret` 处的代码。我们在第二章阐述第一个用户进程的时候跟踪分析了这段代码，在那里第一个用户进程通过执行 `trapret` 处的代码来退出到用户空间。同样地事情在这里也发生：弹出中断帧会恢复用户模式下的寄存器，然后执行 `iret` 会跳回到用户空间。

现在我们讨论的是发生在用户模式下的中断，但是中断也可能发生在内核模式下。在那种情况下硬件不需要进行栈转换，也不需要保存栈指针或栈的段选择符；除此之外的别的步骤都和发生在用户模式下的中断一样，执行的 xv6 中断处理程序的代码也是一样的。而 `iret` 会恢复了一个内核模式下的 `%cs`，处理器也会继续在内核模式下执行。

代码：C 中断处理程序

我们在上一节中看到每一个处理程序会建立一个中断帧然后调用 C 函数 `trap`。`trap`（3101）查看硬件中断号 `tf->trapno` 来判断自己为什么被调用以及应该做些什么。如果中断是 `T_SYSCALL`，`trap` 调用系统调用处理程序 `syscall`。我们会在第五章再来讨论这里的两个 `cp->killed` 检查。

当检查完是否是系统调用，`trap` 会继续检查是否是硬件中断（我们会在下面讨论）。中断可能来自硬件设备的正常中断，也可能来自异常的、未预料到的硬件中断。

如果中断不是一个系统调用也不是一个硬件主动引发的中断，`trap` 就认为它是一个发生中断前的一段代码中的错误行为导致的（如除零错误）。如果产生中断的代码来自用户程序，xv6 就打印错误细节并且设置 `cp->killed` 使之待会被清除掉。我们会在第五章看看 xv6 是怎样进行清除的。

如果是内核程序正在执行，那就出现了一个内核错误：`trap` 打印错误细节并且调用 `panic`。

代码：系统调用

对于系统调用，`trap` 调用 `syscall`（3375）。`syscall` 从中断帧中读出系统调用号，中断帧也包括被保存的 `%eax`，以及到系统调用函数表的索引。对第一个系统调用而言，`%eax` 保存的是 `SYS_exec`（3207），并且 `syscall` 会调用第 `SYS_exec` 个系统调用函数表的表项，相应地也就调用了 `sys_exec`。

`syscall` 在 `%eax` 保存系统调用函数的返回值。当 `trap` 返回用户空间时，它从 `cp->tf` 中加载其值到寄存器中。因此，当 `exec` 返回时，它会返回系统调用处理函数返回的返回值（3381）。系统调用按照惯例会在发生错误的时候返回一个小于 0 的数，成功执行时返回正数。如果系统调用号是非法的，`syscall` 会打印错误并且返回 -1。

之后的章节会讲解系统调用的实现。这一章关心的是系统调用的机制。还有一点点的机制没有说到：如何获得系统调用的参数。工具函数 `argint`、`argptr` 和 `argstr` 获得第 `n` 个系统调用参数，他们分别用于获取整数，指针和字符串起始地址。`argint` 利用用户空间的 `%esp` 寄存器定位第 `n` 个参数：`%esp` 指向系统调用结束后的返回地址。参数就恰好在 `%esp` 之上（`%esp+4`）。因此第 `n` 个参数就在 `%esp+4*n`。

`argint` 调用 `fetchint` 从用户内存地址读取值到 `*ip`。`fetchint` 可以简单地将这个地址直接转换成一个指针，因为用户和内核共享同一个页表，但是内核必须检验这个指针的确指向的是用户内存空间的一部分。内核已经设置好了页表来保证本进程无法访问它的私有地址以外的内存：如果一个用户尝试读或者写高于（包含）`p->sz` 的地址，处理器会产生一个段中断，这个中断会杀死此进程，正如我们之前所见。但是现在，我们在内核态中执行，用户提供的任何地址都是有权访问的，因此必须要检查这个地址是在 `p->sz` 之下的。

`argptr` 和 `argint` 的目标是相似的：它解析第 `n` 个系统调用参数。`argptr` 调用 `argint` 来把第 `n` 个参数当做是整数来获取，然后把这个整数看做指针，检查它的确指向的是用户地址空间。注意 `argptr` 的源码中有两次检查。首先，用户的栈指针在获取参数的时候被检查。然后这个获取到得参数作为用户指针又经过了一次检查。

`argstr` 是最后一个用于获取系统调用参数的函数。它将第 `n` 个系统调用参数解析为指针。它确保这个指针是一个 NUL 结尾的字符串并且整个完整的字符串都在用户地址空间中。

系统调用的实现（例如，`sysproc.c` 和 `sysfile.c`）仅仅是封装而已：他们用 `argint`，`argptr` 和 `argstr` 来解析参数，然后调用真正的实现。在第二章，`sys_exec` 利用这些函数来获取参数。

代码：中断

主板上的设备可以产生中断，xv6 必须配置硬件来处理这些中断。没有硬件的支持 xv6 不可能正常使用起来：用户不能够用键盘输入，没有一个能够存储数据的文件系统等等。幸运的是，添加一些简单设备的中断并不会增加太多额外的复杂性。正如我们将会见到的，中断可以使用与系统调用和异常处理相同的代码。

中断和系统调用相似，除了它可以在任何时候产生。主板上的硬件能够在需要的时候向 CPU 发出信号（例如用户在键盘上输入了一个字符）。我们得对设备编程来产生一个中断，然后令 CPU 接受它们的中断。

我们来看一看分时硬件和时钟中断。我们希望分时硬件大约以每秒 100 次的速度产生一个中断，这样内核就可以对进程进行时钟分片。100 次每秒的速度足以提供良好的交互性能并且同时不会使处理器进入不断的中断处理中。

像 x86 处理器一样，PC 主板也在进步，并且提供中断的方式也在进步。早期的主板有一个简单的可编程中断控制器（被称作 PIC），你可以在 `picirq.c` 中找到管理它的代码。

随着多核处理器主板的出现，需要一种新的处理中断的方式，因为每一颗 CPU 都需要一个中断控制器来处理发送给它的中断，而且也得有一个方法来分发中断。这一方式包括两个部分：第一个部分是在 I/O 系统中的（IO APIC，`ioapic.c`），另一部分是关联在每一个处理器上的（局部 APIC，`lapic.c`）。xv6 是为搭载多核处理器的主板设计的，每一个处理器都需要编程接受中断。

为了在单核处理器上也能够正常运行，xv6 也为 PIC 编程（6932）。每一个 PIC 可以处理最多 8 个中断（设备）并且将他们接到处理器的中断引脚上。为了支持多于八个硬件，PIC 可以进行级联，典型的主板至少有两级级联。使用 `inb` 和 `outb` 指令，xv6 配置主 PIC 产生 IRQ 0 到 7，从 PIC 产生 IRQ 8 到 16。最初 xv6 配置 PIC 屏蔽所有中断。`timer.c` 中的代码设置时钟 1 并且使能 PIC 上相应的中断（7574）。这样的说法忽略了编写 PIC 的一些细节。这些 PIC（也包括 IOAPIC 和 LAPIC）的细节对本书来说并不重要，但是感兴趣的读者可以参考 xv6 源码引用的各设备的手册。

在多核处理器上，xv6 必须编写 IOAPIC 和每一个处理器的 LAPIC。IO APIC 维护了一张表，处理器可以通过内存映射 I/O 写这个表的表项，而非使用 `inb` 和 `outb` 指令。在初始化的过程中，xv6 将第 0 号中断映射到 IRQ 0，以此类推，然后把它们都屏蔽掉。不同的设备自己开启自己的中断，并且同时指定哪一个处理器接受这个中断。举例来说，xv6 将键盘中断分发到处理器 0（7516）。将磁盘中断分发到编号最大的处理器，你们将在下面看到。

时钟芯片是在 LAPIC 中的，所以每一个处理器可以独立地接收时钟中断。xv6 在 `lapicinit`（6651）中设置它。关键的一行代码是 `timer`（6664）中的代码，这行代码告诉 LAPIC 周期性地在 `IRQ_TIMER`（也就是 IRQ 0）产生中断。第 6693 行打开 CPU 的 LAPIC 的中断，这使得 LAPIC 能够将中断传递给本地处理器。

处理器可以通过设置 `eflags` 寄存器中的 `IF` 位来控制自己是否想要收到中断。指令 `cli` 通过清除 `IF` 位来屏蔽中断，而 `sti` 又打开一个中断。xv6 在启动主 cpu（8412）和其他 cpu（1126）时屏蔽中断。每个处理器的调度器打开中断（2464）。为了控制一些特殊的代码片段不被中断，xv6 在进入这些代码片段之前关中断（例如 `switchvm`（1773））。

xv6 在 `idtinit`（1265）中设置时钟中断触发中断向量 32（xv6 使用它来处理 IRQ 0）。中断向量 32 和中断向量 64（用于实现系统调用）的唯一区别就是 32 是一个中断门，而 64 是一个陷阱门。中断门会清除 `IF`，所以被中断的处理器在处理当前中断的时候不会接受其他中断。从这儿开始直到 `trap` 为止，中断执行和系统调用或异常处理相同的代码——建立中断帧。

当因时钟中断而调用 `trap` 时，`trap` 只完成两个任务：递增时钟变量的值（3063），并且调用 `wakeup`。我们将在第 5 章看到后者可能会使得中断返回到一个不同的进程。

驱动程序

驱动程序是操作系统中用于管理某个设备的代码：它提供设备相关的中断处理程序，操纵设备完成操作，操纵设备产生中断，等等。驱动程序可能会非常难写，因为它和它管理的设备同时在并发地运行着。另外，驱动程序必须要理解设备的接口（例如，哪一个 I/O 端口是做什么的），而设备的接口又有可能非常复杂并且文档稀缺。

xv6 的硬盘驱动程序给我们提供了一个良好的例子。磁盘驱动程序从磁盘上拷出和拷入数据。磁盘硬件一般将磁盘上的数据表示为一系列的 512 字节的块（亦称扇区）：扇区 0 是最初的 512 字节，扇区 1 是下一个，以此类推。为了表示磁盘扇区，操作系统也有一个数据结构与之对应。这个结构中存储的数据往往和磁盘上的不同步：可能还没有从磁盘中读出（磁盘正在读数据但是还没有完全读出），或者它可能已经被更新但还没有写出到磁盘。磁盘驱动程序必须保证 xv6 的其他部分不会因为不同步的问题而产生错误。

代码：磁盘驱动程序

通过 IDE 设备可以访问连接到 PC 标准 IDE 控制器上的磁盘。IDE 现在不如 SCSI 和 SATA 流行，但是它的接口比较简单使得我们可以专注于驱动程序的整体结构而不是硬件的某个特别部分的细节。

磁盘驱动程序用结构体 `buf`（称为缓冲区）（3500）来表示一个磁盘扇区。每一个缓冲区表示磁盘设备上的一个扇区。域 `dev` 和 `sector` 给出了设备号和扇区号，域 `data` 是该磁盘扇区数据的内存中的拷贝。

域 `flags` 记录了内存和磁盘的联系：`B_VALID` 位代表数据已经被读入，`B_DIRTY` 位代表数据需要被写出。`B_BUSY` 位是一个锁；它代表某个进程正在使用这个缓冲区，其他进程必须等待。当一个缓冲区的 `B_BUSY` 位被设置，我们称这个缓冲区被锁住。

内核在启动时通过调用 `main`（1234）中的 `ideinit`（3851）初始化磁盘驱动程序。`ideinit` 调用 `pickenable` 和 `ioapickenable` 来打开 `IDE_IRQ` 中断（3856-3857）。调用 `pickenable` 打开单处理器的中断；`ioapickenable` 打开多处理器的中断，但只是打开最后一个 CPU 的中断（`ncpu-1`）：在一个双处理器系统上，CPU 1 专门处理磁盘中断。

接下来，`ideinit` 检查磁盘硬件。它最初调用 `idewait`（3858）来等待磁盘接受命令。PC 主板通过 I/O 端口 0x1f7 来表示磁盘硬件的状态位。`idewait`（3833）获取状态位，直到 `busy` 位（`IDE_BSY`）被清除，以及 `ready` 位（`IDE_DRDY`）被设置。

现在磁盘控制器已经就绪，`ideinit` 可以检查有多少磁盘。它假设磁盘 0 是存在的，因为启动加载器和内核都是从磁盘 0 加载的，但它必须检查磁盘 1。它通过写 I/O 端口 0x1f6 来选择磁盘 1 然后等待一段时间，获取状态位来查看磁盘是否就绪（3860-3867）。如果不就绪，`ideinit` 认为磁盘不存在。

`ideinit` 之后，就只能通过块高速缓冲（buffer cache）调用 `iderw`，`iderw` 根据标志位更新一个锁住的缓冲区。如果 `B_DIRTY` 被设置，`iderw` 将缓冲区的内容写到磁盘；如果 `B_VALID` 没有被设置，`iderw` 从磁盘读出数据到缓冲区。

磁盘访问耗时在毫秒级，对于处理器来说是很漫长的。引导加载器发出磁盘读命令并反复读磁盘状态位直到数据就绪。这种轮询或者忙等待的方法对于引导加载器来说是可以接受的，因为没有更好的事儿可做。但是在操作系统中，更有效的方法是让其他进程占有 CPU 并且在磁盘操作完成时接受一个中断。`iderw` 采用的就是后一种方法，维护一个等待中的磁盘请求队列，然后用中断来指明哪一个请求已经完成。虽然 `iderw` 维护了一个请求的队列，简单的 IDE 磁盘控制器每次只能处理一个操作。磁盘驱动程序的原则是：它已将队首的缓冲区送至磁盘硬件；其他的只是在等待他们被处理。

`iderw`（3954）将缓冲区 `b` 送到队列的末尾（3967-3971）。如果这个缓冲区在队首，`iderw` 通过 `idestart` 将它送到磁盘上（3924-3926）；在其他情况下，一个缓冲区被开始处理当且仅当它前面的缓冲区被处理完毕。

`idestart` 发出关于缓冲区所在设备和扇区的读或者写操作，根据标志位的情况不同。如果操作是一个写操作，`idestart` 必须提供数据（3889）而在写出到磁盘完成后会发出一个中断。如果操作是一个读操作，则发出一个代表数据就绪的中断，然后中断处理程序会读出数据。注意 `iderw` 有一些关于 IDE 设备的细节，并且在几个特殊的端口进行读写。如果任何一个 `outb` 语句错误了，IDE 就会做一些我们意料之外的事。保证这些细节正确也是写设备驱动程序的一大挑战。

`iderw` 已经将请求添加到了队列中，并且会在必要的时候开始处理，`iderw` 还必须等待结果。就像我们之前讨论的，轮询并不是有效的利用 CPU 的办法。相反，`iderw` 睡眠，等待中断处理程序在操作完成时更新缓冲区的标志位（3978-3979）。当这个进程睡眠时，xv6 会调度其他进程来保持 CPU 处于工作状态。

最终，磁盘会完成自己的操作并且触发一个中断。`trap` 会调用 `ideintr` 来处理它（3124）。`ideintr`（3902）查询队列中的第一个缓冲区，看正在发生什么操作。如果该缓冲区正在被读入并且磁盘控制器有数据在等待，`ideintr` 就会调用 `insl` 将数据读入缓冲区（3915-3917）。现在缓冲区已经就绪了：`ideintr` 设置 `B_VALID`，清除 `B_DIRTY`，唤醒任何一个睡眠在这个缓冲区上的进程（3919-3922）。最终，`ideintr` 将下一个等待中的缓冲区传递给磁盘（3924-3926）。

实际情况

想要完美的支持所有的设备需要投入大量的工作，这是因为各种各样的设备有各种各样的特性，设备和驱动之间的协议有时会很复杂。在很多操作系统当中，各种驱动合起来的代码数量要比系统内核的数量更多。

实际的设备驱动远比这一章的磁盘驱动要复杂的多，但是他们的基本思想是一样的：设备通常比 CPU 慢，所以硬件必须使用中断来提醒系统它的状态发生了改变。现代磁盘控制器一般在同一时间接受多个未完成的磁盘请求，甚至重排这些请求使得磁盘使用可以得到更高的效率。当磁盘没有这项功能时，操作系统经常负责重排请求队列。

很多操作系统可以驱动固态硬盘，因为固态硬盘提供了更快的数据访问速度。虽然固态硬盘和传统的机械硬盘的工作机制很不一样，但是这两种设备都使用了基于块的接口，在固态硬盘上读写块仍然要比在内存读写成本高的多。

其他硬件和磁盘非常的相似：网络设备缓冲区保存包，音频设备缓冲区保存音频采样，显存保存图像数据和指令序列。高带宽的设备，如硬盘，显卡和网卡在驱动中同样都使用直接内存访问（Direct memory access, DMA）而不是直接用 I/O（`insl`，`outsl`）。DMA 允许磁盘控制器或者其他控制器直接访问物理内存。驱动给予设备缓存数据区域的物理地址，可以让设备直接地从主存中读取或者写入，一旦复制完成就发出中断。使用 DMA 意味着 CPU 不直接参与传输，这样做可以提高 CPU 工作效率，并且 CPU Cache 开销也更小。

在这一章中的绝大多数设备使用了 I/O 指令来进行编程，但这都是针对老设备的了。而所有现代设备都使用内存映射 I/O（Memory-mapped I/O）来进行编程。

有些设备动态地在轮询模式和中断模式之间切换，因为使用中断的成本很高，但是在驱动去处理一个事件之前，使用轮询会导致延迟。举个例子，对于一个收到大量包的网络设备来说，可能会从中断模式到轮询模式之间切换，因为它知道会到来更多的包被处理，使用轮询会降低处理它们的成本。一旦没有更多的包需要处理了，驱动可能就会切换回中断模式，使得当有新的包到来的时候能够被立刻通知。

IDE 硬盘的驱动静态的发送中断到一个特定的处理器上。有些驱动使用了复杂的算法来发送中断，使得处理负载均衡，并且达到良好的局部性。例如，一个网络驱动程序可能为一个网络连接的包向处理这个连接的处理器发送一个中断，而来自其他连接的包的中断发送给另外的处理器。这种分配方式很复杂；例如，如果有某些网络连接的活动时间很短，但是其他的网络连接却很长，这时候操作系统就要保持所有的处理器都工作来获得一个高的吞吐量。

用户在读一个文件的时候，这个文件的数据将会被拷贝两次。第一次是由驱动从硬盘拷贝到内核内存，之后通过 `read` 系统调用，从内核内存拷贝到用户内存。同理当在网络上发送数据的时候，数据也是被拷贝了两次：先是从用户内存到内核空间，然后是从内核空间拷贝到网络设备。对于很多程序来说低延迟是相当重要的（比如说 Web 服务器服务一个静态页面），操作系统使用了一些特别的代码来避免这种多次拷贝。一个在真实世界中的例子就是缓冲区大小通常是符合内存页大小的，这使得只读的数据拷贝可以直接通过分页映射到进程的地址空间，而不用任何的复制。

练习

1. 在 `syscall()` 的第一条指令处设置一个断点来截获第一次系统调用（例如 `br syscall`）。此时栈上有一些什么值？解释这个断点下执行 `x/37x $esp` 的输出，对每一个值说明它的含义（如为 `trap` 保存的 `%ebp`，`trapframe.eip`，临时分配空间等等）。
2. 添加一个新的系统调用
3. 添加一个网络驱动

第4章

锁

xv6 运行在多处理器上，即计算机上有多个单独执行代码的 CPU。这些 CPU 操作同一片地址空间并分享其中的数据结构；xv6 必须建立一种合作机制防止它们互相干扰。即使是在单个处理器上，xv6 也必须使用某些机制来防止中断处理程序与非中断代码之间互相干扰。xv6 为这两种情况使用了相同的低层概念：锁。锁提供了互斥功能，保证某个时间点只有一个 CPU 能持有锁。如果 xv6 只能在持有特定的锁时才能使用数据结构，那么就能保证同一时间只有一个 CPU 能使用这个数据结构。这样，我们就称这个锁保护了数据结构。

本章的其余部分将解释为何 xv6 需要锁，以及 xv6 是如何实现、使用锁的。我们需要重点注意的是在读代码时，你一定要问自己另一个处理器的存在是否会让这行代码无法达到期望的运行结果（因为另一个处理器也可能正在运行该行代码，或者另一行修改这个共享变量的代码），还要考虑如果这里执行一个中断处理程序，又会发生什么情况。与此同时，一定要记住一行 C 代码可能由多条机器指令组成，而另一个处理器或者中断可能在这些指令中间影响之。你不能假设这些代码是顺序执行的，也不能假设一个 C 指令是以原子操作执行的。并发使得考虑代码的正确性变得困难。

竞争条件

下面举一个例子说明为什么我们需要锁，考虑几个共享磁盘的处理器，例如 xv6 中的 IDE 磁盘。磁盘驱动会维护一个未完成磁盘请求的链表（3821），这样处理器可能会并发地向链表中加入新的请求（3954）。如果没有并发请求，你可以这样实现：

```
struct list{
    int data;
    struct list *next;
};

struct list *list = 0;

void
insert(int data)
{
    struct list *l;
    l = malloc(sizeof *l);
    l->data = data;
    l->next = list;
    list = l;
}
```

证明其正确性是数据结构与算法课中的练习。即使可以证明其正确性，实际上这种实现也是错误的，至少不能在多处理器上运行。如果两个不同的 CPU 同时执行 `insert`，可能会两者都运行到15行，而都未开始运行16行（见图表4-1）。这样的话，就会出现两个链表节点，并且 `next` 都被设置为 `list`。当两者都运行了16行的赋值后，后运行的一个会覆盖前运行的一个；于是先赋值的一个进程中添加的节点就丢失了。这种问题就被称为竞争条件。竞争问题在于它们的结果由 CPU 执行时间以及其内存操作的先后决定的，并且这个问题难以重现。例如，在调试 `insert` 时加入输出语句，就足以改变执行时间，使得竞争消失。

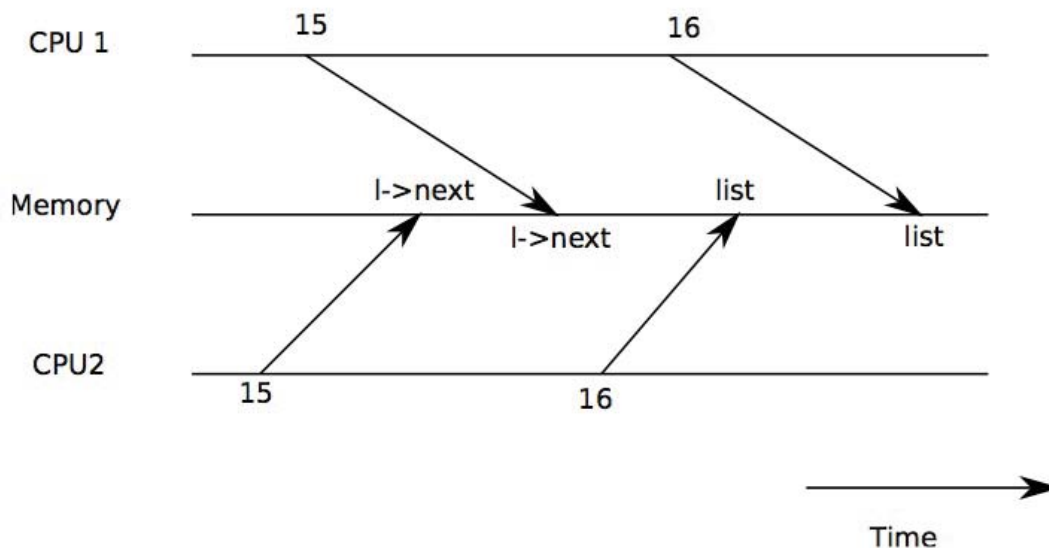


Figure 4-1. Example race

通常我们使用锁来避免竞争。锁提供了互斥，所以一时间只有一个 CPU 可以运行 `insert`；这就让上面的情况不可能发生。只需加入几行代码（未标号的）就能修改为正确的带锁代码：

```
struct list *list = 0;
struct lock listlock;

void
insert(int data)
{
    struct list *l;
    acquire(&listlock);
    l = malloc(sizeof *l);
    l->data = data;
    l->next = list;
    list = l;
    release(&listlock);
}
```

当我们说锁保护了数据时，是指锁保护了数据对应的一组不变量（invariant）。不变量是数据结构在操作中维护的一些状态。一般来说，操作的正确行为会取决于不变量是否为真。操作是有可能暂时破坏不变量的，但在结束操作之前必须恢复不变量。例如，在链表中，不变量即 `list` 指向链表中第一个节点，而每个节点的 `next` 指向下一个节点。`insert` 的实现就暂时破坏了不变量：第13行建立一个新链表元素 `l`，并认为 `l` 是链表中的第一个节点，但 `l` 的 `next` 还没有指向下一个节点（在第15行恢复了该不变量），而 `list` 也还没有指向 `l`（在第16行恢复了该不变量）。上面所说的竞争之所以发生，是因为可能有另一个 CPU 在这些不变量（暂时）没有被恢复的时刻运行了依赖于不变量的代码。恰当地使用锁就能保证一时间只有一个 CPU 操作数据结构，这样在不变量不正确时就不可能有其他 CPU 对数据结构进行操作了。

代码：锁

xv6 用结构体 `struct spinlock`（1401）。结构体中的临界区用 `locked` 表示。这是一个字，在锁可以被获得时值为0，而当锁已经被获得时值为非零。逻辑上讲，xv6 应该用下面的代码来获得锁：

```
void
acquire(struct spinlock *lk)
{
    for(;;) {
        if(!lk->locked) {
            lk->locked = 1;
            break;
        }
    }
}
```

然而这段代码在现代处理器上并不能保证互斥。有可能两个（或多个）CPU 接连执行到第25行，发现 `lk->locked` 为0，然后都执行第26、27行拿到了锁。这时，两个不同的 CPU 持有锁，违反了互斥。这段代码不仅不能帮我们避免竞争条件，它本身就存在竞争。这里的问题主要出在第25、26行是分开执行的。若要保证代码的正确，就必须让第25、26行是原子操作的。

为了让这两行变为原子操作，xv6 采用了386硬件上的一条特殊指令 `xchg`（0569）。在这个原子操作中，`xchg` 交换了内存中的一个字和一个寄存器的值。函数 `acquire`（1474）在循环中反复使用 `xchg`；每一次都读取 `lk->locked` 然后设置为1（1483）。如果锁已经被持有了，`lk->locked` 就已经为1了，故 `xchg` 会返回1然后继续循环。如果 `xchg` 返回0，但是 `acquire` 已经成功获得了锁，即 `locked` 已经从0变为了1，这时循环可以停止了。一旦锁被获得了，`acquire` 会记录获得锁的 CPU 和栈信息，以便调试。当某个进程获得了锁却没有释放时，这些信息可以帮我们找到问题所在。当然这些信息也被锁保护着，只有在持有锁时才能修改。

函数 `release`（1502）则做了相反的事：清除调试信息并释放锁。

模块化与递归锁

系统设计力求简单、模块化的抽象：最好是让调用者不需要了解被调者的具体实现。锁的机制则和这种模块化理念有所冲突。例如，当 CPU 持有锁时，它不能再调用另一个试图获得该锁的函数 `f`：因为调用者在 `f` 返回之前无法释放锁，如果 `f` 试图获得这个锁，就会造成死锁。

现在还没有一种透明方案可以让调用者和被调者可以互相隐藏所使用的锁。我们可以使用递归锁（*recursive locks*）使得被调者能够在此获得调用者已经持有的锁，这种方案虽然是透明通用的，但是十分繁复。还有一个问题就是这种方案不能用来保护不变量。在 `insert` 调用 `acquire(&listlock)` 后，它就可以假设没有其他函数会持有这个锁，也没有其他函数可以操作链表，最重要的是，可以保持链表相关的所有不变量。在使用递归锁的系统中，`insert` 可以假设在它之后 `acquire` 不会再被调用：`acquire` 之所以能成功，只可能是 `insert` 的调用者持有锁，并正在修改链表数据。这时的不变量有可能被破坏了，链表也就不再保护其不变量了。锁不仅要让不同的 CPU 不会互相干扰，还需要让调用者与被调者不会互相干扰；而递归锁就无法保证这一点。

由于没有理想、透明的解决方法，我们不得不在函数的使用规范中加入锁。编程者必须保证一个函数不会在持有锁时调用另一个需要获得该锁的函数 `f`。就这样，锁也成为了我们的抽象中的一员。

代码：使用锁

xv6 非常谨慎地使用锁来避免竞争条件。一个简单的例子就是 IDE 驱动（3800）。就像本章开篇提到的一样，`iderw`（3954）有一个磁盘请求的队列，处理器可能会并发地向队列中加入新请求（3969）。为了保护链表以及驱动中的其他不变量，`iderw` 会请求获得锁 `idelock`（3965）并在函数末尾释放锁。练习1中研究了如何通过把 `acquire` 移动到队列操作之后来触发竞争条件。我们很有必要做一个这些练习，它们会让我们了解到想要触发竞争并不容易，也就是说很难找到竞争条件。并不是说 xv6 的代码中就没有竞争。

使用锁的一个难点在于要决定使用多少个锁，以及每个锁保护哪些数据、不变量。不过有几个基本原则。首先，当一个 CPU 正在写一个变量，而同时另一个 CPU 可能读/写该变量时，需要用锁防止两个操作重叠。第二，当用锁保护不变量时，如果不变量涉及到多个数据结构，通常每个数据结构都需要用一个单独的锁保护起来，这样才能维持不变量。

上面只说了需要锁的原则，那么什么时候不需要锁呢？由于锁会降低并发度，所以我们一定要避免过度使用锁。当效率不是很重要的时候，完全可以使用单处理器计算机，这样就完全不用考虑锁了。当我们要保护内核的数据结构时，使用一个内核锁还是值得的，当进入内核时必须持有该锁，而退出内核时就释放该锁。许多单处理器操作系统就用这种方法运行在了多处理器上，有时这种方法被称为“内核巨锁（giant kernel lock）”，但使用这种方法就牺牲了并发性：即一时间只有一个 CPU 可以运行在内核上。如果我们想要依靠内核做大量的计算，那么使用一组更为精细的锁来让内核可以在多个 CPU 上轮流运行会更有效率。

最后，对于锁的粒度选择是并行编程中的一个重要问题。xv6 只使用了几个简单的锁；例如，xv6 中使用了一个单独的锁来保护进程表及其不变量，我们将在第5章讨论这个问题。更精细的做法是给进程表中的每一个条目都上一个锁，这样在不同条目上运行的线程也能并行了。但是在进程表中维护那么多个不变量就必须使用多个锁，这就让情况变得很复杂了。不过 xv6

中的例子已经足够让我们了解如何使用锁了。

锁的顺序

如果一段代码要使用多个锁，那么必须要注意代码每次运行都要以相同的顺序获得锁，否则就有死锁的危险。假设某段代码的两条执行路径都需要锁 A 和 B，但路径1获得锁的顺序是 A、B，而路径2获得锁的顺序是 B、A。这样就有能路径1获得了锁 A，而在它继续获得锁 B 之前，路径2获得了锁 B，这样就死锁了。这时两个路径都无法继续执行下去了，因为这时路径1需要锁 B，但锁 B 已经在路径2手中了，反之路径2也得不到锁 A。为了避免这种死锁，所有的代码路径获得锁的顺序必须相同。避免死锁也是我们吧锁作为函数使用规范的一部分的原因：调用者必须以固定顺序调用函数，这样函数才能以相同顺序获得锁。

由于 xv6 本身比较简单，它使用的锁也很简单，所以 xv6 几乎没有锁的使用链。最长的锁链也就只有两个锁。例如，`ideintr` 在调用 `wakeup` 时持有 `ide` 锁，而 `wakeup` 又需要获得 `ptable.lock`。还有很多使用 `sleep / wakeup` 的例子，它们要考虑锁的顺序是因为 `sleep` 和 `wakeup` 中有比较复杂的不变量，我们会在第5章讨论。文件系统中有很多两个锁的例子，例如文件系统在删除一个文件时必须持有该文件及其所在文件夹的锁。xv6 总是首先获得文件夹的锁，然后再获得文件的锁。

中断处理程序

xv6 用锁来防止中断处理程序与另一个 CPU 上运行非中断代码使用同一个数据。例如，时钟中断（3114）会增加 `ticks` 但可能有另一个 CPU 正在运行 `sys_sleep`，其中也要使用该变量（3473）。锁 `tickslock` 就能够为该变量实现同步。

即使在单个处理器上，中断也可能导致并发：在允许中断时，内核代码可能在任何时候停下来，然后执行中断处理程序。假设 `iderw` 持有 `idelock`，然后中断发生，开始运行 `ideintr`。`ideintr` 会试图获得 `idelock`，但却发现 `idelock` 已经被获得了，于是就等着它被释放。这样，`idelock` 就永远不会被释放了，只有 `iderw` 能释放它，但又只有让 `ideintr` 返回 `iderw` 才能继续运行，这样处理器、整个系统都会死锁。

为了避免这种情况，当中断处理程序会使用某个锁时，处理器就不能在允许中断发生时持有锁。xv6 做得更决绝：允许中断时不能持有任何锁。它使用 `pushcli`（1555）和 `popcli`（1566）来屏蔽中断（`cli` 是 x86 屏蔽中断的指令）。`acquire` 在尝试获得锁之前调用了 `pushcli`（1476），`release` 则在释放锁后调用了 `popcli`（1521）。`pushcli`（1555）和 `popcli`（1566）不仅包装了 `cli` 和 `sti`，它们还做了计数工作，这样就需要调用两次 `popcli` 来抵消两次 `pushcli`；这样，如果代码中获得了两个锁，那么只有当两个锁都被释放后中断才会被允许。

`acquire` 一定要在可能获得锁的 `xchg` 之前调用 `pushcli`（1483）。如果两者颠倒了，就可能在几个时钟周期里，中断仍被允许，而锁也被获得了，如果此时不幸地发生了中断，系统就会死锁。类似的，`release` 也一定要在释放锁的 `xchg` 之后调用 `popcli`（1483）。

另外，中断处理程序和非中断代码对彼此的影响也让我们看到了递归锁的缺陷。如果 xv6 使用了递归锁（即如果 CPU 获得了某个锁，那么同一 CPU 上可以再次获得该锁），那么中断处理程序就可能在非中断代码正运行到临界区时运行，这样就非常混乱了。当中断处理程序运行时，它所依赖的不变量可能暂时被破坏了。例如，`ideintr`（3902）会假设未处理请求链表是完好的。若 xv6 使用了递归锁，`ideintr` 就可能在 `iderw` 正在修改链表，这样 `ideintr` 就会使用这个不正确的链表。

内存乱序

在本章中，我们都假设了处理器会按照代码中的顺序执行指令。但是许多处理器会通过指令乱序来提高性能。如果一个指令需要多个周期完成，处理器会希望这条指令尽早开始执行，这样就能与其他指令交叠，避免延误太久。例如，处理器可能会发现一系列 A、B 指令序列彼此并没有关系，在 A 之前执行 B 可以让处理器执行完 A 时也执行完 B。但是并发可能会让这种乱序行为暴露到软件中，导致不正确的结果。

例如，考虑在 `release` 中把0赋给 `lk->locked` 而不是使用 `xchg`。那么结果就不明确了，因为我们难以保证这里的执行顺序。比方说如果 `lk->locked=0` 在乱序后被放到了 `popcli` 之后，可能在锁被释放之前，另一个线程中就允许中断了，`acquire` 就会被打断。为了避免乱序可能造成的不确定性，xv6 决定使用稳妥的 `xchg`，这样就能保证不出现乱序了。

现实情况

由于使用了锁机制的程序编写仍然是个巨大的挑战，所以并发和并行至今还是研究的热点。我们最好以锁为基础来构建高级的同步队列，虽然 xv6 并没有这么做。如果你使用锁进行编程，那么你最好用一些工具来确定竞争条件，否则很容易遗漏掉某些需要锁保护的不变量。

用户级程序也需要锁，但 xv6 的程序只有一个运行线程，进程间也不会共享内存，所以就不需要锁了。

当然我们也有可能用非原子性的操作来实现锁，只不过那非常复杂，而且大多数的操作系统都是使用了原子操作的。

原子操作的代价也不小。如果一个处理器在它的本地缓存中有一个锁，而这时另一个处理器必须获得该锁，那么更新缓存中该行的原子操作就必须把这行从一个处理器的缓存中移到另一个处理器的缓存中，同时还可能需要让这行缓存的其他备份失效。从其他处理器的缓存中取得一行数据要比从本地缓存中取代价大得多。

为了减少使用锁所产生的代价，许多操作系统使用了锁无关的数据结构和算法，并在这些算法中尽量避免原子操作。例如，对于本章开篇提到的链表，我们在查询时不需要获得锁，然后用一个原子操作来添加元素。

练习

1. 若在 `acquire` 中不用 `xchg`，运行 xv6 会发生什么情况？
2. 把 `iderw` 中的 `acquire` 移到 `sleep` 之前会出现竞争吗？你可以通过运行 xv6 并运行 `stressfs` 来观察。用简单的循环扩大临界区看看会发生什么，并对此作出解释。
3. 完成公布的作业。
4. 在缓冲区的 `flags` 中置位并不是原子操作：处理器会在寄存器中拷贝一份 `flags`，修改寄存器然后写回去。所以两个处理器不能同时写 `flags`。XV6 只在持有 `buflock` 的时候修改 `B_BUSY`，但修改 `B_VALID` 和 `B_WRITE` 的时候并没有锁。为什么这么做仍然是安全的呢？

第5章

调度

任何操作系统都可能碰到进程数多于处理器数的情况，这样就需要考虑如何分享处理器资源。理想的做法是让分享机制对进程透明。通常我们对进程造成一个自己独占处理器的假象，然后让操作系统的多路复用机制 (multiplex) 将单独的一个物理处理器模拟为多个虚拟处理器。本章将讲述 xv6 是如何为多个进程模拟出多处理器的。

多路复用

xv6 中多路复用的实现如下：当一个进程等待磁盘请求时，xv6 使之进入睡眠状态，然后调度执行另一个进程。另外，当一个进程耗尽了它在处理器上运行的时间片（100毫秒）后，xv6 使用时钟中断强制它停止运行，这样调度器才能调度运行其他进程。这样的多路复用机制为进程提供了独占处理器的假象，类似于 xv6 使用内存分配器和页表硬件为进程提供了独占内存的假象。

实现多路复用有几个难点。首先，应该如何从运行中的一个进程切换到另一个进程？xv6 采用了普通的上下文切换机制；虽然这里的思想是非常简洁明了的，但是其代码实现是操作系统中最晦涩难懂的一部分。第二，如何让上下文切换透明化？xv6 只是简单地使用时钟中断处理程序来驱动上下文切换。第三，可能出现多个 CPU 同时切换进程的情况，那么我们必须使用一个带锁的方案来避免竞争。第四，进程退出时必须释放其占用内存与资源，但由于它本身在使用自己的资源（譬如其内核栈），所以不能由该进程本身释放其占有的所有资源。xv6 希望能够简洁明了地处理这些难点，不过最后其代码实现还是比较“巧妙”。

xv6 必须为进程提供互相协作的方法。譬如，父进程需要等待子进程结束，以及读取管道数据的进程需要等待其他进程向管道中写入数据。与其让这些等待中的进程消耗 CPU 资源，不如让它们暂时放弃 CPU，进入睡眠状态来等待其他进程发出事件来唤醒它们。但我们必须要小心设计以防睡眠进程遗漏事件通知。本章我们将用管道机制的具体实现来解释上述问题及其解决方法。

代码：上下文切换

如图表5-1所示，xv6 在低层次中实现了两种上下文切换：从进程的内核线程切换到当前 CPU 的调度器线程，从调度器线程到进程的内核线程。xv6 永远不会直接从用户态进程切换到另一个用户态进程；这种切换是通过用户态-内核态切换（系统调用或中断）、切换到调度器、切换到新进程的内核线程、最后这个陷入返回实现的。本节我们将以内核线程与调度器线程的切换作为例子来说明。

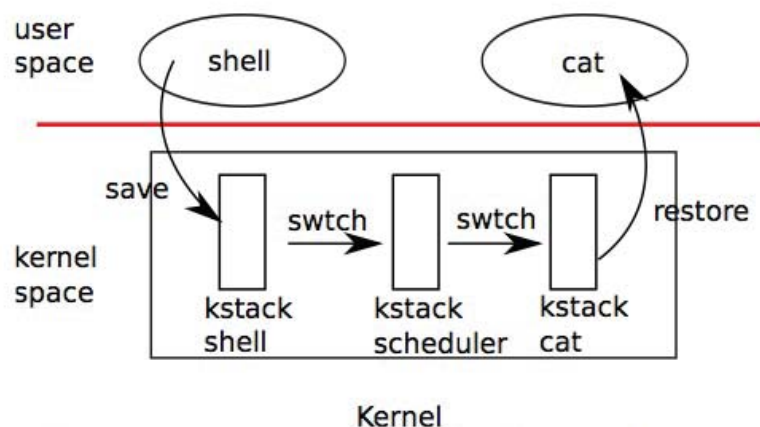


Figure 5-1. Switching from one user process to another. In this example, xv6 runs with one CPU (and thus one scheduler thread).

如我们在第2章中所见，每个 xv6 进程都有自己的内核栈以及寄存器集合。每个 CPU 都有一个单独的调度器线程，这样调度

就不会发生在进程的内核线程中，而是在此调度器线程中。线程的切换涉及到了保存旧线程的 CPU 寄存器，恢复新线程之前保存的寄存器；其中 `%esp` 和 `%eip` 的变换意味着 CPU 会切换运行栈与运行代码。

`swtch` 并不了解线程，它只是简单地保存和恢复寄存器集合，即上下文。当进程让出 CPU 时，进程的内核线程调用 `swtch` 来保存自己的上下文然后返回到调度器的上下文中。每个上下文都是以结构体 `struct context*` 表示的，这实际上是一个保存在内核栈中的指针。`swtch` 有两个参数：`struct context **old`、`struct context *new`。它将当前 CPU 的寄存器压入栈中并将栈指针保存在 `*old` 中。然后 `swtch` 将 `new` 拷贝到 `%esp` 中，弹出之前保存的寄存器，然后返回。

接下来我们先不考察调度器调用 `swtch` 的过程，我们先回到用户进程中看看。在第3章中我们知道，有可能在中断的最后，`trap` 会调用 `yield`。`yield` 又调用 `sched`，其中 `sched` 会调用 `swtch` 来保存当前上下文到 `proc->context` 中然后切换到之前保存的调度器上下文 `cpu->scheduler`（2516）。

`swtch`（2702）一开始从栈中弹出参数，放入寄存器 `%eax` 和 `%edx`（2709-2710）中；`swtch` 必须在改变栈指针以及无法获得 `%esp` 前完成这些事情。然后 `swtch` 压入寄存器，在当前栈上建立一个新的上下文结构。仅有被调用者保存的寄存器此时需要被保存；按照 x86 的惯例即 `%ebp %ebx %esi %ebp %esp`。`swtch` 显式地压入前四个寄存器（2713-2716）；最后一个则是在 `struct context*` 被写入 `old`（2719）时隐式地保存的。要注意，还有一个重要的寄存器，即程序计数器 `%eip`，该寄存器在使用 `call` 调用 `swtch` 时就保存在栈中 `%ebp` 之上的位置上了。保存了旧寄存器后，`swtch` 就准备要恢复新的寄存器了。它将指向新上下文的指针放入栈指针中（2720）。新的栈结构和旧的栈相同，因为新的上下文其实是之前某次的切换中的旧上下文。所以 `swtch` 就能颠倒一下保存旧上下文的顺序来恢复新上下文。它弹出 `%edi %esi %ebx %ebp` 然后返回（2723-2727）。由于 `swtch` 改变了栈指针，所以这时恢复的寄存器就是新上下文中的寄存器值。

在我们的例子中，`sched` 调用 `swtch` 切换到 `cpu->scheduler`，即 per-cpu 的调度器上下文。这个上下文是在之前 `scheduler` 调用 `swtch`（2478）时保存的。当 `swtch` 返回时，它不会返回到 `sched` 中，而是返回到 `scheduler`，其栈指针指向了当前 CPU 的调度器的栈，而非 `initproc` 的内核栈。

代码：调度

上一节中我们查看了 `swtch` 的底层细节；现在让我们将 `swtch` 看做一个既有的功能，来研究从进程到调度器然后再回到进程的切换过程中的一些约定。进程想要让出 CPU 必须要获得进程表的锁 `ptable.lock`，并释放其拥有的其他锁，修改自己的状态（`proc->state`），然后调用 `sched`。`yield`（2522）和 `sleep exit` 都遵循了这个约定，我们稍后将会详细研究。`sched` 检查了两次状态（2507-2512），这里的状态表明由于进程此时持有锁，所以 CPU 应该是在中断关闭的情况下运行的。最后，`sched` 调用 `swtch` 把当前上下文保存在 `proc->context` 中然后切换到调度器上下文即 `cpu->scheduler` 中。`swtch` 返回到调度器栈中，就像是调度器调用的 `swtch` 返回了一样（2478）。调度器继续其 `for` 循环，找到一个进程来运行，切换到该进程，然后继续轮转。

我们看到，在对 `swtch` 的调用的整个过程中，xv6 都持有锁 `ptable.lock`：`swtch` 的调用者必须持有该锁，并将锁的控制权转移给切换代码。锁的这种使用方式很少见，通常来说，持有锁的线程应该负责释放该锁，这样更容易让我们理解其正确性。但对于上下文切换来说，我们必须使用这种方式，因为 `ptable.lock` 会保证进程的 `state` 和 `context` 在运行 `swtch` 时保持不变。如果在 `swtch` 中没有持有 `ptable.lock`，可能引发这样的问题：在 `yield` 将某个进程状态设置为 `RUNNABLE` 之后，但又是在 `swtch` 让它停止在其内核栈上运行之前，有另一个 CPU 要运行该进程。其结果将是两个 CPU 都运行在同一个栈上，这显然是不该发生的。

内核线程只可能在 `sched` 中让出处理器，在 `scheduler` 中切换回对应的地方，当然这里 `scheduler` 也是通过 `sched` 切换到进程中的。所以，如果要写出 xv6 中切换线程的代码行号，我们会发现其执行规律是（2478），（2516），（2516），不断循环。以这种形式在两个线程之间切换的过程有时被称作共行程序（*coroutines*）；在此例中，`sched` 和 `scheduler` 就是彼此的共行程序。

但在一种特殊情况下，调度器调用的切换到新进程的 `swtch` 不会在 `sched` 中结束。我们在第2章学到了这个例子：当一个新进程第一次被调度时，它从 `forkret`（2533）开始运行。之所以要运行 `forkret`，只是为了按照惯例释放 `ptable.lock`；否则，这个新进程是可以就从 `trapret` 开始运行的。

`scheduler`（2458）运行了一个普通的循环：找到一个进程来运行，运行直到其停止，然后继续循环。`scheduler` 大部分时间里都持有 `ptable.lock`，但在每次外层循环中都要释放该锁（并显式地允许中断）。当 CPU 闲置（找不到 `RUNNABLE` 的进程）时这样做十分有必要。如果一个闲置的调度器一直持有锁，那么其他 CPU 就不可能执行上下文切换或任何和进程相关

的系统调用了，也就更不可能将某个进程标记为 `RUNNABLE` 然后让闲置的调度器能够跳出循环了。而之所以周期性地允许中断，则是因为可能进程都在等待 I/O，从而找不到一个 `RUNNABLE` 的进程（例如 `shell`）；如果调度器一直不允许中断，I/O 就永远无法到达了。

`scheduler` 不断循环寻找可运行，即 `p->state == RUNNABLE` 的进程。一旦它找到了这样的进程，就将 `per-cpu` 的当前进程变量 `proc` 设为该进程，用 `switchvm` 切换到该进程的页表，标记该进程为 `RUNNING`，然后调用 `swtch` 切换到该进程中运行（2472-2478）。

下面我们来从另一个层面研究这段调度代码。对于每个进程，调度维护了进程的一系列固定状态，并且保证当状态变化时必须持有锁 `ptable.lock`。第一个固定状态是，如果进程为 `RUNNING` 的，那么必须确保使用时钟中断的 `yield` 时，能够无误地切换到其他进程；这就意味着 CPU 寄存器必须保存着进程的寄存器值（这些寄存器值并非在 `context` 中），`%cr3` 必须指向进程的页表，`%esp` 则要指向进程的内核栈，这样 `swtch` 才能正确地向栈中压入寄存器值，另外 `proc` 必须指向进程的 `proc[]` 槽中。另一个固定状态是，如果进程是 `RUNNABLE`，必须保证调度器能够无误地调度执行该进程；这意味着 `p->context` 必须保存着进程的内核线程变量，并且没有任何 CPU 此时正在其内核栈上运行，没有任何 CPU 的 `%cr3` 寄存器指向进程的页表，也没有任何 CPU 的 `proc` 指向该进程。

正是由于要坚持以上两个原则，所以 `xv6` 必须在一个线程中获得 `ptable.lock`（通常是在 `yield` 中），然后在另一个线程中释放这个锁（在调度器线程或者其他内核线程中）。如果一段代码想要将运行中进程的状态修改为 `RUNNABLE`，那么在恢复到固定状态中之前持有锁；最早的可以释放锁的时机是在 `scheduler` 停止使用该进程页表并清空 `proc` 时。类似地，如果 `scheduler` 想把一个可运行进程的状态修改为 `RUNNING`，在该进程的内核线程完全运行起来（`swtch` 之后，例如在 `yield` 中）之前必须持有锁。

除此之外，`ptable.lock` 也保护了一些其他的状态：进程 ID 的分配，进程表槽的释放，`exit` 和 `wait` 之间的互动，保证对进程的唤醒不会被丢失等等。我们应该思考一下 `ptable.lock` 有哪些不同的功能可以分离，使之更为简洁高效。

睡眠与唤醒

锁的机制使得 CPU 之间，进程之间不会互相干扰；调度使得进程可以共享 CPU。但是我们现在还不知道进程之间是如何交换信息的。睡眠和唤醒实际上就提供了进程间通信的机制，它们可以让一个进程暂时休眠，等待某个特定事件的发生，然后当特定事件发生时，另一个进程会唤醒该进程。睡眠与唤醒通常被称为顺序合作（*sequence coordination*）或者有条件同步（*conditional synchronization*）机制，在操作系统的哲学中，还有很多类似的机制。

为了说明，假设有一个生产者/消费者队列。这个队列有些类似于 IDE 驱动用来同步处理器和设备驱动的队列（见第3章），不过下面所讲的更能概括 IDE 驱动中的代码。该队列允许一个进程将一个非零指针发送给另一个进程。假设只有一个发送者和一个接受者，并且它们运行在不同的 CPU 上，那么下面的实现显然是正确的：

```
struct q {
    void *ptr;
};

void*
send(struct q *q, void *p)
{
    while(q->ptr != 0)
        ;
    q->ptr = p;
}

void*
recv(struct q *q)
{
    void *p;
    while((p = q->ptr) == 0)
        ;
    q->ptr = 0;
    return p;
}
```

`send` 会不断循环，直到队列为空（`ptr == 0`），然后将指针 `p` 放到队列中。`recv` 会不断循环，直到队列非空然后取出指

针。当不同的进程运行时，`send` 和 `recv` 会同时修改 `q->ptr`，不过 `send` 只在队列空时写入指针，而 `recv` 只在队列非空时拿出指针，这样他们之间是不会互相干扰的。

上面这种实现方法固然正确，但是代价是巨大的。如果发送者很少发送，那么接受者就会消耗大量的时间在 `while` 循环中苦苦等待一个指针的出现。而实际上如果有一种方法使得 `send` 放入指针时，能够通知接受者。那么接受者所在的 CPU 就能在这段时间找到更有意义的事情做。

让我们来考虑一对调用 `sleep` 和 `wakeup`，其工作方式如下。`sleep(chan)` 让进程在任意的 `chan` 上休眠，称之为等待队列（*wait channel*）。`sleep` 让调用进程休眠，释放所占 CPU。`wakeup(chan)` 则唤醒在 `chan` 上休眠的所有进程，让他们的 `sleep` 调用返回。如果没有进程在 `chan` 上等待唤醒，`wakeup` 就什么也不做。让我们用 `sleep` 和 `wakeup` 来重新实现上面的代码：

```
void*
send(struct q *q, void *p)
{
    while(q->ptr != 0)
        ;
    q->ptr = p;
    wakeup(q); /*wake recv*/
}

void*
recv(struct q *q)
{
    void *p;
    while((p = q->ptr) == 0)
        sleep(q);
    q->ptr = 0;
    return p;
}
```

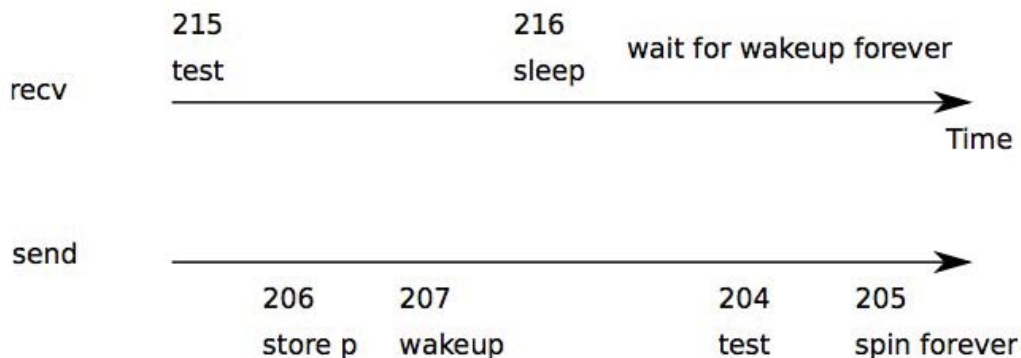


Figure 5-2. Example lost wakeup problem

令人动容的是，现在 `recv` 能够让出 CPU 而不是空等浪费资源了。但对于图表5-2中出现的“遗失的唤醒”问题，我们却很难通过已有的接口下，直观地设计出能够避免该问题的 `sleep` 和 `wakeup` 机制。假设在第215行 `recv` 发现 `q->ptr == 0`，然后决定调用 `sleep`，但是在 `recv` 调用 `sleep` 之前（譬如这时处理器突然收到一个中断然后开始执行中断处理，延迟了对 `sleep` 的调用），`send` 又在另一个 CPU 上运行了，它将 `q->ptr` 置为非零，然后调用 `wakeup`，发现没有进程在休眠，于是什么也没有做。接着，`recv` 从第216行继续执行了：它调用 `sleep` 进入休眠。这就出现问题了，休眠的 `recv` 实际上在等待一个已经到达的指针。而下一个 `send` 又在睡眠中等着 `recv` 取出队列中的指针。这种情况就被称为死锁（*deadlock*）。

这个问题的根源在于没有维持好一个固定状态，即由于 `send` 在错误的时机运行了，而使得 `recv` 只能在 `q->ptr == 0` 时睡眠这个行为被妨碍了。下面我们还将看到一段能保护该固定状态但仍有问题的代码：

```
struct q {
    struct spinlock lock;
```



```

    void *ptr;
};

void *
send(struct q *q, void *p)
{
    acquire(&q->lock);
    while(q->ptr != 0)
        ;
    q->ptr = p;
    wakeup(q);
    release(&q->lock);
}

void*
recv(struct q *q)
{
    void *p;
    acquire(&q->lock);
    while((p = q->ptr) == 0)
        sleep(q);
    q->ptr = 0;
    release(&q->lock);
    return p;
}

```

由于要调用 `sleep` 的进程是持有锁 `q->lock` 的，而 `send` 想要调用 `wakeup` 也必须获得锁，所以这种方案能够保护上面讲到的固定状态。但是这种方案也会出现死锁：当 `recv` 带着锁 `q->lock` 进入睡眠后，发送者就会在希望获得锁时一直阻塞。

所以想要解决问题，我们必须改变 `sleep` 的接口。`sleep` 必须将锁作为一个参数，然后在进入睡眠状态后释放之；这样就能避免上面提到的“遗失的唤醒”问题。一旦进程被唤醒了，`sleep` 在返回之前还需要重新获得锁。于是我们应该使用下面的代码：

```

struct q {
    struct spinlock lock;
    void *ptr;
};

void *
send(struct q *q, void *p)
{
    acquire(&q->lock);
    while(q->ptr != 0)
        ;
    q->ptr = p;
    wakeup(q);
    release(&q->lock);
}

void*
recv(struct q *q)
{
    void *p;
    acquire(&q->lock);
    while((p = q->ptr) == 0)
        sleep(q, &q->lock);
    q->ptr = 0;
    release(&q->lock);
    return p;
}

```

`recv` 持有 `q->lock` 就能防止 `send` 在 `recv` 检查 `q->ptr` 与调用 `sleep` 之间调用 `wakeup` 了。当然，为了避免死锁，接收进程最好别在睡眠时仍持有锁。所以我们希望 `sleep` 能用原子操作释放 `q->lock` 并让接收进程进入休眠状态。

完整的发送者/接收者的实现还应该让发送者在等待接收者拿出前一个 `send` 放入的值时处于休眠状态。

代码：睡眠与唤醒

接下来让我们看看 xv6 中 `sleep` 和 `wakeup` 的实现。总体思路是希望 `sleep` 将当前进程转化为 `SLEEPING` 状态并调用 `sched` 以释放 CPU，而 `wakeup` 则寻找一个睡眠状态的进程并把它标记为 `RUNNABLE`。

`sleep` 首先会进行几个必要的检查：必须存在当前进程（2555）并且 `sleep` 必须持有锁（2558-2559）。接着 `sleep` 要求持有 `ptable.lock`（2568）。于是该进程就会同时持有锁 `ptable.lock` 和 `lk` 了。调用者（例如 `recv`）是必须持有 `lk` 的，这样可以保证其他进程（例如一个正在运行的 `send`）无法调用 `wakeup(chan)`。而如今 `sleep` 已经持有了 `ptable.lock`，那么它现在就能安全地释放 `lk` 了：这样即使别的进程调用了 `wakeup(chan)`，`wakeup` 也不可能在没有持有 `ptable.lock` 的情况下运行，所以 `wakeup` 必须等待 `sleep` 让进程睡眠后才能运行。这样一来，`wakeup` 就不会错过 `sleep` 了。

这里有一个复杂一点的情况：即 `lk` 就是 `ptable.lock` 的时候，这样 `sleep` 在要求持有 `ptable.lock` 然后又把它作为 `lk` 释放的时候会出现死锁。这种情况下，`sleep` 就会直接跳过这两个步骤（2567）。例如，当 `wait`（2403）持有 `&ptable.lock` 时调用 `sleep`。

现在仅有该进程的 `sleep` 持有 `ptable.lock`，于是它通过记录睡眠队列，改变进程状态，调用 `sched`（2573-2575）让进程进入睡眠。

稍后，进程会调用 `wakeup(chan)`。`wakeup`（2603）要求获得 `ptable.lock` 并调用 `wakeup1`，其中，实际工作是由 `wakeup1` 完成的。对于 `wakeup` 来说持有 `ptable.lock` 也是很重要的，因为它也要修改进程的状态并且要保证 `sleep` 和 `wakeup` 不会错过彼此。而之所以要单独实现一个 `wakeup1`，是因为有时调度器会在持有 `ptable.lock` 的情况下唤醒进程，稍后我们会看到这样的例子。当 `wakeup` 找到了对应 `chan` 中处于 `SLEEPING` 的进程时，它将进程状态修改为 `RUNNABLE`。于是下一次调度器在运行时，就可以调度该进程了。

`wakeup` 必须在有一个监视唤醒条件的锁的时候才能被调用：在上面的例子中这个锁就是 `q->lock`。至于为什么睡眠中的进程不会错过唤醒，则是因为从 `sleep` 检查进程状态之前，到进程进入睡眠之后，`sleep` 都持有进程状态的锁或者 `ptable.lock` 或者是两者兼有。由于 `wakeup` 必须在持有这两个锁的时候运行，所以它必须在 `sleep` 检查状态之前和一个进程已经完全进入睡眠后才能执行。

有些情况下可能有多个进程在同一队列中睡眠；例如，有多个进程想要从管道中读取数据时。那么单独一个 `wakeup` 的调用就能将它们全部唤醒。他们的其中一个会首先运行并要求获得 `sleep` 被调用时所持的锁，然后读取管道中的任何数据。而对于其他进程来说，即使被唤醒了，它们也读不到任何数据，所以唤醒它们其实是徒劳的，它们还得接着睡。正是由于这个原因，我们在一个检查状态的循环中不断调用 `sleep`。

`sleep` 和 `wakeup` 的调用者可以使用任何方便使用的数字作为队列号码；而实际上，xv6 通常使用内核中和等待相关的数据结构的地址，譬如磁盘缓冲区。即使两组 `sleep / wakeup` 使用了相同的队列号码，也是无妨的：对于那些无用的唤醒，它们会通过不断检查状态忽略之。`sleep / wakeup` 的优点主要是其轻量级（不需另定义一个结构来作为睡眠队列），并且提供了一层抽象（调用者不需要了解与之交互的是哪一个进程）。

代码：管道

上面我们提到的队列只是一个简单的模型，实际上在 xv6 中有两个使用 `sleep / wakeup` 来同步读者写者的队列。一个在 IDE 驱动中：进程将未完成的磁盘请求加入队列，然后调用 `sleep`。中断处理程序会使用 `wakeup` 告诉进程其磁盘请求已经完成了。

更为复杂的一个例子是管道。我们在第0章已经介绍了管道的接口：我们从管道的一端写入数据字节，然后数据被拷贝到内核缓冲区中，接着就能从管道的另一端读取数据了。后面的章节中还会讲到文件系统是怎样帮助我们实现管道的。不过现在先让我们来看看 `pipewrite` 和 `piperead` 的实现。

每个管道由一个结构体 `struct pipe` 表示，其中有一个锁 `lock` 和内存缓冲区。其中的域 `nread` 和 `nwrite` 表示从缓冲区读出和写入的字节数。`pipe` 对缓冲区做了包装，使得虽然计数器继续增长，但实际上在 `buf[PIPE_SIZE - 1]` 之后写入的字节存放在 `buf[0]`。这样就让我们可以区分一个满的缓冲区（`nwrite == nread + PIPE_SIZE`）和一个空的缓冲区（`nwrite == nread`），但这也意味着我们必须使用 `buf[nread % PIPE_SIZE]` 而不是 `buf[nread]` 来读出/写入数据。现在假设 `piperead` 和 `pipewrite` 分别在两个 CPU 上连续执行。

`pipewrite`（6080）首先请求获得管道的锁，以保护计数器、数据以及相关不变量。`piperead`（6101）接着也请求获得锁，

结果当然是无法获得。于是它停在了 `acquire` (1474) 上等待锁的释放。与此同时, `pipewrite` 在循环中依次写入 `addr[0]`, `addr[1]`, ..., `addr[n-1]` 并添加到管道中 (6904)。在此循环中, 缓冲区可能被写满 (6086), 这时 `pipewrite` 会调用 `wakeup` 通知睡眠中的读者缓冲区中有数据可读, 然后使得在 `&p->nwrite` 队列中睡眠的读者从缓冲区中读出数据。注意, `sleep` 在让 `pipewrite` 的进程进入睡眠时还会释放 `p->lock`。

现在 `p->lock` 被释放了, `piperead` 尝试获得该锁然后开始执行: 此时它会检查到 `p->nread != p->nwrite` (6106) (正是在 `nwrite == nread + PIPESIZE` (6086) 的时候 `pipewrite` 进入了睡眠), 于是 `piperead` 跳出 `for` 循环, 将数据从管道中拷贝出来 (6113-6117), 然后将 `nread` 增加读取字节数。现在缓冲区又多出了很多可写的字节, 所以 `piperead` 调用 `wakeup` (6118) 唤醒睡眠中的写者, 然后返回到调用者中。 `wakeup` 会找到在 `&p->nwrite` 队列上睡眠的进程, 正是该进程之前在运行 `pipewrite` 时由于缓冲区满而停止了。然后 `wakeup` 将该进程标记为 `RUNNABLE`。

代码: `wait`, `exit`, `kill`

`sleep` 和 `wakeup` 可以在很多需要等待一个可检查状态的情况中使用。如我们在第0章中所见, 父进程可以调用 `wait` 来等待一个子进程退出。在 `xv6` 中, 当一个子进程要退出时它并不是直接死掉, 而是将状态转变为 `ZOMBIE`, 然后当父进程调用 `wait` 时才能发现子进程可以退出了。所以父进程要负责释放退出的子进程相关的内存空间, 并修改对应的 `struct proc` 以便重用。每个进程结构体都会在 `p->parent` 中保存指向其父进程的指针。如果父进程在子进程之前退出了, 初始进程 `init` 会接收其子进程并等待它们退出。我们必须这样做以保证可以为退出的进程做好子进程的清理工作。另外, 所有的进程结构都是被 `ptable.lock` 保护的。

`wait` 首先要求获得 `ptable.lock`, 然后查看进程表中是否有子进程, 如果找到了子进程, 并且没有子进程已经退出, 那么就调用 `sleep` 等待其中一个子进程退出 (2439), 然后不断循环。注意, 这里 `sleep` 中释放的锁是 `ptable.lock`, 也就是我们之前提到过的特殊情况。

`exit` 首先要求获得 `ptable.lock` 然后唤醒当前进程的父进程 (2376)。这一步似乎为时过早, 但由于 `exit` 这时还没有把当前进程标记为 `ZOMBIE`, 所以这样并不会出错: 即使父进程已经是 `RUNNABLE` 的了, 但在 `exit` 调用 `sched` 以释放 `ptable.lock` 之前, `wait` 是无法运行其中的循环的。所以说只有在子进程被标记为 `ZOMBIE` (2388)之后, `wait` 才可能找到要退出的子进程。在退出并重新调度之前, `exit` 会把所有子进程交给 `initproc` (2378-2385)。最后, `exit` 调用 `sched` 来让出 CPU。

退出进程的父进程本来是通过调用 `wait` (2439) 处于睡眠状态中, 不过现在它就可以被调度器调度了。对 `sleep` 的调用返回时仍持有 `ptable.lock`; `wait` 接着会重新查看进程表并找到 `state == ZOMBIE` (2382) 的已退出子进程。它会记录该子进程的 `pid` 然后清理其 `struct proc`, 释放相关的内存空间 (2418-2426)。

子进程在 `exit` 时已经做好了大部分的清理工作, 但父进程一定要为其释放 `p->kstack` 和 `p->pgdir`; 当子进程运行 `exit` 时, 它正在利用 `p->kstack` 分配到的栈以及 `p->pgdir` 对应的页表。所以这两者只能在子进程结束运行后, 通过调用 `sched` 中的 `swtch` 被清理。这就是为什么调度器要运行在自己单独的栈上, 而不能运行在调用 `sched` 的线程的栈上。

`exit` 让一个应用程序可以自我终结; `kill` (2625) 则让一个应用程序可以终结其他进程。在实现 `kill` 时有两个困难: 1) 被终结的进程可能正在另一个 CPU 上运行, 所以它必须在被终结之前把 CPU 让给调度器; 2) 被终结的进程可能正在 `sleep` 中, 并持有内核资源。 `kill` 很轻松地解决了这两个难题: 它在进程表中设置被终结的进程的 `p->killed`, 如果这个进程在睡眠中则唤醒之。如果被终结的进程正在另一个处理器上运行, 它总会通过系统调用或者中断 (例如时钟中断) 进入内核。当它离开内核时, `trap` 会检查它的 `p->killed`, 如果被设置了, 该进程就会调用 `exit`, 终结自己。

如果被终结的进程在睡眠中, 调用 `wakeup` 会让被终结的进程开始运行并从 `sleep` 中返回。此举有一些隐患, 因为进程可能是在它等待的状态尚为假的时候就从 `sleep` 中返回了。所以 `xv6` 谨慎地在调用 `sleep` 时使用了 `while` 循环, 检查 `p->killed` 是否被设置了, 若是, 则返回到调用者。调用者也必须再次检查 `p->killed` 是否被设置, 若是, 返回到再上一级调用者, 依此下去。最后进程的栈展开 (`unwind`) 到了 `trap`, `trap` 若检查到 `p->killed` 被设置了, 则调用 `exit` 终结自己。我们已经在管道的实现中 (6087) 看到了在 `sleep` 外的 `while` 循环中检查 `p->killed`。

有一处的 `while` 没有检查 `p->killed`。 `ide` 驱动 (3979) 直接重新调用了 `sleep`。之所以可以确保能被唤醒, 是因为它在等待一个磁盘中断。如果它不是在等待磁盘中断的话, `xv6` 就搞不清楚它在做什么了。如果有第二个进程在中断之前调用了 `iderw`, `ideintr` 会唤醒该进程 (第二个), 而非原来等待中断的那一个 (第一个) 进程。第二个进程会认为它收到了它正在等待的数据, 但实际上它收到的是第一个进程想要读的数据。

现实情况

xv6 所实现的调度算法非常朴素，仅仅是让每个进程轮流执行。这种算法被称作轮转法（*round robin*）。真正的操作系统使用了更为复杂的算法，例如，让每个进程都有优先级。主要思想是优先处理高优先级的可运行进程。但是由于要权衡多项指标，例如要保证公平性和高的吞吐量，调度算法往往很快变得复杂起来。另外，复杂的调度算法还会无意中导致像优先级倒转（*priority inversion*）和护航（*convoy*）这样的现象。优先级倒转是指当高优先级进程和低优先级进程共享一个锁时，如果锁已经被低优先级进程获得了，高优先级进程就无法运行了。护航则是指当很多高优先级进程等待一个持有锁的低优先级进程的情况，护航一旦发生，则可能持续很久。如果要避免这些问题，就必须在复杂的调度器中设计更多的机制。

`sleep` 和 `wakeup` 是非常普通但有效的同步方法，当然还有很多其他的同步方法。同步要解决的第一个问题是本章开始我们看到的“丢失的唤醒”问题。原始的 Unix 内核的 `sleep` 只是简单地屏蔽中断，由于 Unix 只在单处理器上运行，所以这样已经足够了。但是由于 xv6 要在多处理器上运行，所以它给 `sleep` 增加了一个现实的锁。FreeBSD 的 `msleep` 也使用了同样的方法。Plan 9 的 `sleep` 使用了一个回调函数，并在其返回到 `sleep` 中之前持有调度用的锁；这个函数对睡眠状态作了最后检查，以避免丢失的唤醒。Linux 内核的 `sleep` 用一个显式的进程队列代替 xv6 中的等待队列（wait channel）；而该队列本身内部还有锁。

在 `wakeup` 中遍历整个进程表来寻找对应 `chan` 的进程是非常低效的。更好的办法是用另一个结构体代替 `sleep` 和 `wakeup` 中的 `chan`，该结构体中维护了一个睡眠进程的链表。Plan 9 的 `sleep` 和 `wakeup` 将该结构体称为集合点（*rendezvous point*）或者 *Rendez*。许多线程库都把相同的一个结构体作为一个状态变量；如果是这样的话，`sleep` 和 `wakeup` 操作则被称为 `wait` 和 `signal`。所有此类机制都有同一个思想：使得睡眠状态可以被某种执行原子操作的锁保护。

在 `wakeup` 的实现中，它唤醒了特定队列中的所有进程，而有可能很多进程都在同一个队列中等待被唤醒。操作系统会调度这里的所有进程，它们会互相竞争以检查其睡眠状态。这样的一堆进程被称作惊群（*thundering herd*），而我们最好是避免这种情况的发生。大多数的状态变量都有两种不同的 `wakeup`，一种唤醒一个进程，即 `signal`；另一种唤醒所有进程，即 `broadcast`。

信号量是另一种合作机制。一个信号量是一个提供两种操作，即加和减的整数。我们可以不断对一个信号量进行加操作，但是我们只能在信号量为正时进行减操作，否则当我们对为零的信号量减时，进程会进入睡眠，直到另一个进程对其进行加操作，这对加减操作就会抵消掉。这个整数实际上反映了一个真正的计数值，例如管道缓冲区中可写的字节数，或者进程拥有的僵尸子进程数。在这个抽象中显式地计数可以避免“丢失的唤醒”问题：通过显式地记录唤醒次数。这样计数同时还能避免无谓的唤醒和惊群问题。

终结、清理进程是 xv6 中较为复杂的内容。而大多数操作系统则更为复杂。例如，被终结的进程可能在内核深处睡眠，我们需要谨慎地编写代码以展开其栈。许多操作系统使用显式的处理异常的机制来展开栈，例如使用 `longjmp`。另外，有些事件可能让睡眠的进程被唤醒，即使这些事件尚未发生。例如，当进程在睡眠时，另一个进程向它发送一个信号。进程就会从被打断的系统调用中返回-1并在 `EINTR` 中存放错误代码。应用程序可以查看这些值以决定下一步怎么做。xv6 并不支持信号，所以没有这么复杂。

练习

1. `sleep` 必须检查 `lk != &ptable.lock` 以避免死锁（2567-2570）。我们可以通过把代码

```
if (lk != &ptable.lock){
    acquire(&ptable.lock);
    release(lk);
}
```

替换为

```
release(lk);
acquire(&ptable.lock);
```

来避免对特殊情况的判断。这样做会打断 `sleep` 吗？如果会，是在什么情况下？

2.大部分的进程清理工作可以由 `exit` 或 `wait` 完成，但我们在上面看到 `exit` 不能释放 `p->stack`。不过必须由 `exit` 来关闭打开文件。这是为什么？答案会涉及到管道（pipe）。

3.在 xv6 中实现信号量。你可以使用互斥锁，但不能使用睡眠和唤醒。用信号量代替 xv6 中的睡眠和唤醒并评价你的结果。

第六章

文件系统

文件系统的目的是组织和存储数据，典型的文件系统支持用户和程序间的数据共享，并提供数据持久化的支持（即重启之后数据仍然可用）。

xv6 的文件系统中使用了类似 Unix 的文件，文件描述符，目录和路经名（请参阅第零章），并且把数据存储到一块 IDE 磁盘上（请参阅第三章）。这个文件系统解决了几大难题：

- 该文件系统需要磁盘上数据结构来表示目录树和文件，记录每个文件用于存储数据的块，以及磁盘上哪些区域是空闲的。
- 该文件系统必须支持崩溃恢复，也就是说，如果系统崩溃了（比如掉电了），文件系统必须保证在重启后仍能正常工作。问题在于一次系统崩溃可能打断一连串的更新操作，从而使得磁盘上的数据结构变得不一致（例如：有的块同时被标记为使用中和空闲）。
- 不同的进程可能同时操作文件系统，要保证这种并行不会破坏文件系统的正常工作。
- 访问磁盘比访问内存要慢几个数量级，所以文件系统必须要维护一个内存内的 cache 用于缓存常被访问的块。

这一章的剩余内容将阐述 xv6 是如何解决这些问题的。

概述

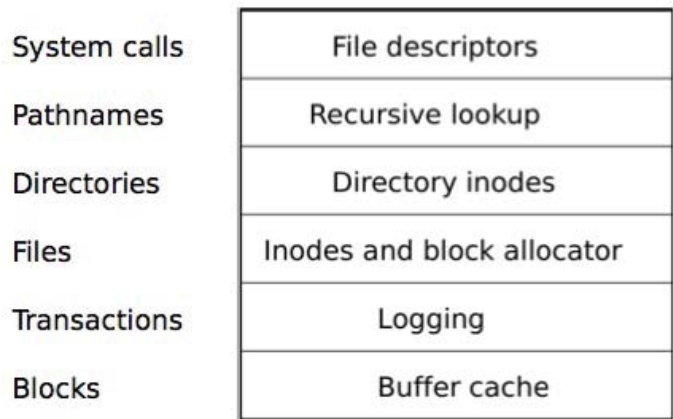


Figure 6-1. Layers of the xv6 file system.

xv6 的文件系统分6层实现，如图6-1所示。最下面一层通过块缓冲读写 IDE 硬盘，它同步了对磁盘的访问，保证同时只有一个内核进程可以修改磁盘块。第二层使得更高层的接口可以将对磁盘的更新按会话打包，通过会话的方式来保证这些操作是原子操作（要么都被应用，要么都不被应用）。第三层提供无名文件，每一个这样的文件由一个 i 节点和一连串的数据块组成。第四层将目录实现为一种特殊的 i 节点，它的内容是一连串的目录项，每一个目录项包含一个文件名和对应的 i 节点。第五层提供了层次路经名（如usr/rtn/xv6/fs.c这样的），这一层通过递归的方式来查询路径对应的文件。最后一层将许多 UNIX 的资源（如管道，设备，文件等）抽象为文件系统的接口，极大地简化了程序员的工作。

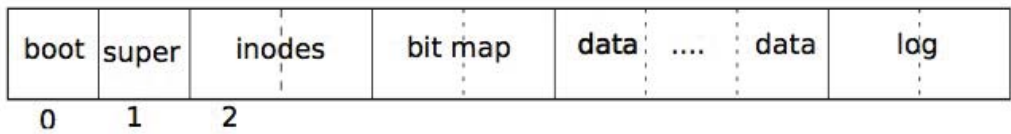


Figure 6-2. Structure of the xv6 file system. The header fs.h (3650) contains constants and data structures describing the exact layout of the file system.

文件系统必须设计好在磁盘上的什么地方放置 i 节点和数据块。xv6 把磁盘划分为几个区块，如图6-2所示。文件系统不使用第0块（第0块存有 bootloader）。第1块叫做超级块；它包含了文件系统的元信息（如文件系统的总块数，数据块块数，i 节点数，以及日志的块数）。从第2块开始存放 i 节点，每一块能够存放多个 i 节点。接下来的块存放空闲块位图。剩下的大部分块是数据块，它们保存了文件和目录的内容。在磁盘的最后是日志块，它们是会话层的一部分，将在后面详述。

块缓冲层

块缓冲有两个任务：（1）同步对磁盘的访问，使得对于每一个块，同一时间只有一份拷贝放在内存中并且只有一个内核线程使用这份拷贝；（2）缓存常用的块以提升性能。代码参见 `bio.c`。

块缓冲提供的主要接口是 `bread` 和 `bwrite`；前者从磁盘中取出一块放入缓冲区，后者把缓冲区中的一块写到磁盘上正确的地方。当内核处理完一个缓冲块之后，需要调用 `brelease` 释放它。

块缓冲仅允许最多一个内核线程引用它，以此来同步对磁盘的访问，如果一个内核线程引用了一个缓冲块，但还没有释放它，那么其他调用 `bread` 的进程就会阻塞。文件系统的更高几层正是依赖块缓冲层的同步机制来保证其正确性。

块缓冲有固定数量的缓冲区，这意味着如果文件系统请求一个不在缓冲中的块，必须换出一个已经使用的缓冲区。这里的置换策略是 LRU，因为我们假设最近未使用的块近期内最不可能再被使用。

代码：块缓冲

块缓冲是缓冲区的双向链表。`binit`（1231）会从一个静态数组 `buf` 中构建出一个有 `NBUF` 个元素的双向链表。所有对块缓冲的访问都通过链表而非静态数组。

一个缓冲区有三种状态：`B_VALID` 意味着这个缓冲区拥有磁盘块的有效内容。`B_DIRTY` 意味着缓冲区的内容已经被改变并且需要写回磁盘。`B_BUSY` 意味着有某个内核线程持有这个缓冲区且尚未释放。

`bread`（4102）调用 `bget` 获得指定扇区的缓冲区（4106）。如果缓冲区需要从磁盘中读出，`bread` 会在返回缓冲区前调用 `iderw`。

`bget`（4066）扫描缓冲区链表，通过给定的设备号和扇区号找到对应的缓冲区（4073-4084）。如果存在这样一个缓冲区，并且它还不是处于 `B_BUSY` 状态，`bget` 就会设置它的 `B_BUSY` 位并且返回。如果找到的缓冲区已经在使用中，`bget` 就会睡眠等待它被释放。当 `sleep` 返回的时候，`bget` 并不能假设这块缓冲区现在可用了，事实上，`sleep` 时释放了

`buf_table_lock`，醒来后重新获取了它，这就不能保证 `b` 仍然是可用的缓冲区：它有可能被用来缓冲另外一个扇区。`bget` 非常无奈，只能重新扫描一次（4082），希望这次能够找到可用的缓冲区。

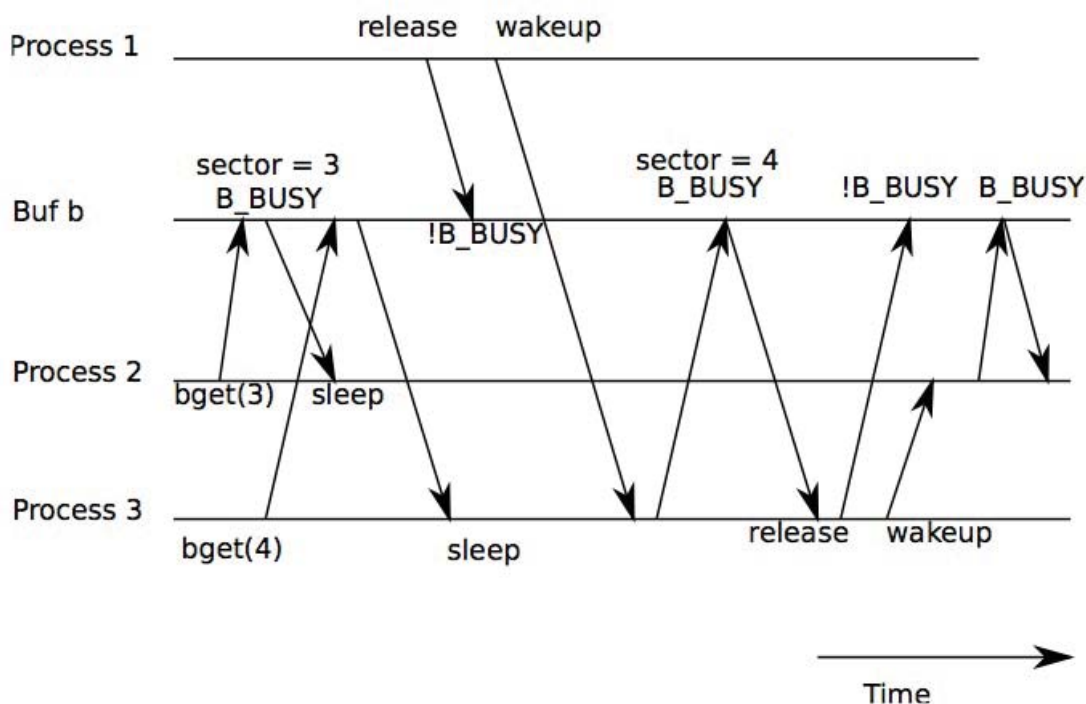


Figure 6-3. A race resulting in process 3 receiving a buffer containing block 4, even though it asked for block 3.

如果 `bget` 中没有那句 `goto` 语句的话，那么就可能会产生图6-3中的竞争。第1个进程有一块缓冲了扇区3的缓冲区。现在另外两个进程来了，第1个进程 `get` 缓冲区3并且为之睡眠（缓冲区3使用中）。第2个进程 `get` 缓冲区4，并且也可能在同一块缓冲区上睡眠，但这次是在等待新分配的循环中睡眠的，因为已经没有空闲的缓冲区，而持有扇区3的缓冲区处于链表头，因此被重用。第一个进程释放了这块缓冲区，`wakeup` 恰巧安排进程3先运行，而后它拿到这块缓冲区后把扇区4读了进来。进程3之后也释放了缓冲区并且唤醒了进程2。如果没有 `goto` 语句的话，进程2就会在把拿到的缓冲区标记为 `BUSY` 后从 `bget` 返回，但实际上这块缓冲区装的是扇区4而不是3。这样的错误可能导致各种各样的麻烦，因为扇区3和4的内容是不同的；实际上，xv6中它们存储着 `i` 节点。

如果所请求的扇区还未被缓冲，`bget` 必须分配一个缓冲区，可能是重用某一个缓冲区。它再次扫描缓冲区列表，寻找一块不是忙状态的块，任何这样的块都可以被拿去使用。`bget` 修改这个块的元数据来记录新的设备号和扇区号并且标记这个块为 `BUSY`，最后返回它（4091-4093）。需要注意的是，对标记位的赋值（4089-4091）不仅设置了 `B_BUSY` 位，也清除了 `B_VALID` 位和 `B_DIRTY` 位，用来保证 `bread` 会用磁盘的内容来填充缓冲区，而不是继续使用块之前的内容。

因为块缓冲是用于同步的，保证任何时候对于每一个扇区都只有一块缓冲区是非常重要的，`bget` 第一个循环确认没有缓冲区已经加载了所需扇区的内容，并且在此之后 `bget` 都没有释放 `buf_table_lock`，因此 `bget` 的操作是安全的。

如果所有的缓冲区都处于忙碌状态，那么就出问题了，`bget` 就会报错。不过一个更优雅的反应是进入睡眠状态，直到有一块缓冲区变为空闲状态。虽然这有可能导致死锁。

一旦 `bread` 给自己的调用者返回了一块缓冲区，调用者就独占了这块缓冲区。如果调用者写了数据，他必须调用 `bwrite`（4114）在释放缓冲区之前将修改了的数据写入磁盘，`bwrite` 设置 `B_DIRTY` 位并且调用的 `iderw` 将缓冲区的内容写到磁盘。

当调用者使用完了一块缓冲区，他必须调用 `brelease` 来释放它，（关于 `brelease` 这个名字，它是 `b-release` 的缩写，它源自 Unix 并且在 BSD, Linux 和 Solaris 中被广泛使用）。`brelease`（4125）将一块缓冲区移动到链表的头部（4132-4137），清除 `B_BUSY`，唤醒睡眠在这块缓冲区上的进程。移动缓冲区的作用在于使得链表按照最近被使用的情况排序，链表中的第一块是最近被用的，最后一块是最早被用的。`bget` 中的两个循环就利用这一点：寻找已经存在的缓冲区在最坏情况下必须遍历整个链表，但是由于数据局部性，从最近使用的块开始找（从 `bcache.head` 开始，然后用 `next` 指针遍历）会大大减少扫描的时间。反之，找一块可重用的缓冲区是从链表头向前找，相当于从尾部往头部通过 `prev` 指针遍历，从而找到的就是最近不被使用的块。

日志层

文件系统设计中最有趣的问题之一就是错误恢复，产生这样的问题是因为大多数的文件系统都涉及到对磁盘多次的写操作，如果在写操作的过程中系统崩溃了，就会使得磁盘上的文件系统处于不一致的状态中。举例来说，根据写的顺序的不同，上述错误可能会导致一个目录项指向一个空闲的 `i` 节点，或者产生一个已被分配但是未被引用的块。后一种情况相对来说好一些，但在前一种情况中，目录项指向了一个空闲的 `i` 节点，重启之后就会导致非常严重的问题。

xv6 通过简单的日志系统来解决文件操作过程中崩溃所导致的问题。一个系统调用并不直接导致对磁盘上文件系统的写操作，相反，他会把一个对磁盘写操作的描述包装成一个日志写在磁盘中。当系统调用把所有的写操作都写入了日志，它就会写一个特殊的提交记录到磁盘上，代表一次完整的操作。从那时起，系统调用就会把日志中的数据写入磁盘文件系统的数据结构中。在那些写操作都成功完成后，系统调用就会删除磁盘上的日志文件。

为什么日志可以解决文件系统操作中出现的崩溃呢？如果崩溃发生在操作提交之前，那么磁盘上的日志文件就不会被标记为已完成，恢复系统的代码就会忽视它，磁盘的状态正如这个操作从未执行过一样。如果是在操作提交之后崩溃的，恢复程序会重演所有的写操作，可能会重复之前已经进行了的对磁盘文件系统的写操作。在任何一种情况下，日志文件都使得磁盘操作对于系统崩溃来说是原子操作：在恢复之后，要么所有的写操作都完成了，要么一个写操作都没有完成。

日志设计

日志存在于磁盘末端已知的固定区域。它包含了一个起始块，紧接着一连串的数据块。起始块包含了一个扇区号的数组，每一个对应于日志中的数据块，起始块还包含了日志数据块的计数。xv6 在提交后修改日志的起始块，而不是之前，并且在将日志中的数据块都拷贝到文件系统之后将数据块计数清0。提交之后，清0之前的崩溃就会导致一个非0的计数值。

每一个系统调用都可能包含一个必须从头到尾原子完成的写操作序列，我们称这样的一个序列为一个会话，虽然他比数据库中的会话要简单得多。任何时候只能有一个进程在一个会话之中，其他进程必须等待当前会话中的进程结束。因此同一时刻日志最多只记录一次会话。

xv6 不允许并发会话，目的是为了下面几种问题。假设会话 X 把一个对 `i` 节点的修改写入了会话中。并发的会话 Y 从同一块中读出了另一个 `i` 节点，更新了它，把 `i` 节点块写入了日志并且提交。这就会导致可怕的后果：Y 的提交导致被 X 修改过的 `i` 节点块被写入磁盘，而 X 此时并没有提交它的修改。如果这时候发生崩溃会使得 X 的修改只应用了一部分而不是全部，从而打破会话是原子的这一性质。有一些复杂的办法可以解决这个问题，但 xv6 直接通过不允许并行的会话来回避这个问题。

xv6 允许只读的系统调用在一次会话中并发执行。`i` 节点锁会使得会话对只读系统调用看上去是原子性的。

xv6 使用固定量的磁盘空间来保存日志。系统调用写入日志的块的总大小不能大于日志的总大小。对于大多数系统调用来说这都不是个问题，但是其中两个可能会写大量的块：`write` 和 `unlink`。写一个大文件可能会写很多的数据块、位图块，以及 `i` 节点块。移除对一个大文件的链接可能会写很多的位图块以及一个 `i` 节点块。xv6 的写系统调用将大的写操作拆分成几个小的写操作，使得被修改的块能放入日志中。`unlink` 不会导致问题因为实际上 xv6 只使用一个位图块。

代码：日志

对日志的常见使用方法像下面这样

```
begin_trans();
...
bp = bread(...);
bp->data[...] = ...;
log_write(bp);
...
commit_trans();
```

`begin_trans` (4277) 会一直等到它独占了日志的使用权后返回。

`log_write` (4325) 像是 `bwrite` 的一个代理；它把块中新的内容记录到日志中，并且把块的扇区号记录在内存

中。 `log_write` 仍将修改后的块留在内存中的缓冲区中，因此相继的本会话中对这一块的读操作都会返回已修改的内容。 `log_write` 能够知道在一次会话中对同一块进行了多次读写，并且覆盖之前同一块的日志。

`commit_trans` (4301) 将日志的起始块写到磁盘上，这样在这个时间点之后的系统崩溃就能够恢复，只需将磁盘中的内容用日志中的内容改写。 `commit_trans` 调用 `install_trans` (4221) 来从日志中逐块的读并把他们写到文件系统中合适的地方。最后 `commit_trans` 会把日志起始块中的计数改为0，这样在下次会话之前的系统崩溃就会使得恢复代码忽略日志。

`recover_from_log` (4268) 在 `initlog` (4205) 中被调用，而 `initlog` 在第一个用户进程开始前的引导过程中被调用。它读取日志的起始块，如果起始块说日志中有一个提交了会话，它就会仿照 `commit_trans` 的行为执行，从而从错误中恢复。

`filewrite` (5352) 中有一个使用了日志的例子：

```
begin_trans();
ilock(f->ip);
r = writei(f->ip, ...);
iunlock(f->ip);
commit_trans();
```

我们在一个用于将一次大的写操作拆分成一些会话的循环中找到了这段代码，在每一次会话中这段只会写部分块，因为日志的大小是有限固定的。对 `writei` 的调用会在一次会话中写很多的块：文件的 `i` 节点，一个或多个位图块，以及一些数据块。在 `begin_trans` 之后再执行 `ilock` 是一种避免死锁的办法：因为每次会话都已经有一个锁保护了，因此在持有两个锁的时候，要保证一定的加锁顺序。

代码：块分配器

文件和目录的内容存在磁盘块中，磁盘块都从一个空闲块池中分配出来。xv6 的块分配器包含一个磁盘上的空闲块位图，每个块占一个位。引导区，超级块，`i` 节点块和位图块的位永远都被置为有效。

块分配器提供两个功能：`balloc` 分配一个新的磁盘块，`bfree` 释放一个块。`balloc` (4454) 最开始调用 `readsb` 从磁盘读出超级块（或者从块缓冲中）到 `sb` 中。`balloc` 会算出位图块的位置，计算的方法是计算多少块被引导区、超级块和 `i` 节点块占用（用 `BBLOCK`）。循环 (4462) 从第0块开始一直到 `sb.size`（文件系统块总数），寻找一个在位图中的位是0的块。为了提高效率，这个循环被分成两块。外层循环读位图的每一块。内层循环检查这一块中的所有 `BPB` 那么多个位。两个进程可能同时申请空闲块，这就有可能导致竞争，但事实上块缓冲只允许一个进程同时只使用一个块。

`bfree` (4481) 找到了空闲的块之后就会清空位图中对应的位。同样，`bread` 和 `breise` 的块的互斥使用使得无需再特意加锁。

i 节点

`i` 节点这个术语可以有两个意思。它可以指的是磁盘上的记录文件大小、数据块扇区号的数据结构。也可以指内存中的一个 `i` 节点，它包含了一个磁盘上 `i` 节点的拷贝，以及一些内核需要的附加信息。

所有的磁盘上的 `i` 节点都被打包在一个称为 `i` 节点块的连续区域中。每一个 `i` 节点的大小都是一样的，所以对于一个给定的数字 `n`，很容易找到磁盘上对应的 `i` 节点。事实上这个给定的数字就是操作系统中 `i` 节点的编号。

磁盘上的 `i` 节点由结构体 `dinode` (3676) 定义。`type` 域用来区分文件、目录和特殊文件的 `i` 节点。如果 `type` 是0的话就意味着这是一个空闲的 `i` 节点。`nlink` 域用来记录指向了这个 `i` 节点的目录项，这是用于判断一个 `i` 节点是否应该被释放的。`size` 域记录了文件的字节数。`addrs` 数组用于这个文件的数据块的块号。

内核在内存中维护活动的 `i` 节点。结构体 `inode` (3762) 是磁盘中的结构体 `dinode` 在内存中的拷贝。内核只会在有 `C` 指针指向一个 `i` 节点的时候才会把这个 `i` 节点保存在内存中。`ref` 域用于统计有多少个 `C` 指针指向它。如果 `ref` 变为0，内核就会丢掉这个 `i` 节点。`iget` 和 `iput` 两个函数申请和释放 `i` 节点指针，修改引用计数。`i` 节点指针可能从文件描述符产生，从当前工作目录产生，也有可能从一些内核代码如 `exec` 中产生。

持有 `iget` 返回的 `i` 节点的指针将保证这个 `i` 节点会留在缓存中，不会被删掉（特别地不会被用于缓存另一个文件）。因此

`iget` 返回的指针相当一种较弱的锁，虽然它并不要求持有者真的锁上这个 `i` 节点。文件系统的许多部分都依赖于这个特性，一方面是为了长期地持有对 `i` 节点的引用（比如打开的文件和当前目录），一方面是在操纵多个 `i` 节点的程序中避免竞争和死锁（比如路径名查找）。

`iget` 返回 `i` 节点可能没有任何有用的内容。为了保证它持有一个磁盘上 `i` 节点的有效拷贝，程序必须调用 `ilock`。它会锁住 `i` 节点（从而其他进程就无法使用它）并从磁盘中读出 `i` 节点的信息（如果它还没有被读出的话）。`iunlock` 释放 `i` 节点上的锁。把对 `i` 节点指针的获取和 `i` 节点的锁分开避免了某些情况下的死锁，比如在目录查询的例子中，数个进程都可以通过 `iget` 获得一个 `i` 节点的 `C` 指针，只有一个进程可以锁住一个 `i` 节点。

`i` 节点缓存只会缓存被 `C` 指针指向的 `i` 节点。它主要的工作是同步多个进程对 `i` 节点的访问而非缓存。如果一个 `i` 节点频繁被使用，块缓冲可能会把它保留在内存中，即使 `i` 节点缓存没有缓存它。

代码：`i` 节点

要申请一个新的 `i` 节点（比如创建文件的时候），`xv6` 会调用 `ialloc`（4603）。`ialloc` 同 `balloc` 类似：它逐块遍历磁盘上的 `i` 节点数据结构，寻找一个标记为空闲的 `i` 节点。当它找到一个时，就会把它的 `type` 修改掉（变为非0），最后调用 `iget`（4620）使得它从 `i` 节点缓存中返回。由于每个时刻只有一个进程能访问 `bp`，`ialloc` 可以保证其他进程不会同时认为这个 `i` 节点是可用的并且获取到它。

`iget`（4654）遍历 `i` 节点缓存寻找一个指定设备和 `i` 节点号的活动中的项（`ip->ref > 0`）。如果它找到一项，它就返回对这个 `i` 节点的引用（4663-4667）。在 `iget` 扫描的时候，它会记录扫描到的第一个空槽（4668-4669），之后如果需要可以用这个空槽来分配一个新的缓存项。

调用者在读写 `i` 节点的元数据或内容之前必须用 `ilock` 锁住 `i` 节点。`ilock`（4703）用一个类似的睡眠循环（这种循环在 `bget` 中见过）来等待 `ip->flag` 的 `I_BUSY` 位被清除，而后由自己再设置它（4712-4714）。一旦 `ilock` 拥有了对 `i` 节点的独占，他可以根据需要从磁盘中读取 `i` 节点的元数据。函数 `iunlock`（4735）清除 `I_BUSY` 位并且唤醒睡眠在 `ilock` 中的其他进程。

`iput`（4756）释放指向 `i` 节点的 `C` 指针，实际上就是将引用计数减1。如果减到了0，那么 `i` 节点缓存中的这个 `i` 节点槽就会变为空闲状态，并且可以被另一个 `i` 节点重用。

如果 `iput` 发现没有指针指向一个 `i` 节点并且也没有任何目录项指向它（不在目录中出现的一个文件），那么这个 `i` 节点和它关联的数据块都应该被释放。`iput` 重新锁上这个 `i` 节点，调用 `itrunc` 来把文件截断为0字节，释放掉数据块；把 `i` 节点的类型设置为0（未分配）；把变化写到磁盘中；最后解锁 `i` 节点（4759-4771）。

`iput` 中对锁的使用方法值得我们研究。第一个值得研究的地方是当锁上 `ip` 的时候，`put` 简单地认为它是没有被锁过的，而非使用一个睡眠循环。这种情况是正常的，因为调用者被要求在调用 `iput` 之前解锁 `ip`，而调用者拥有对它的唯一引用（`ip->ref == 1`）。第二个值得研究的部分是 `iput` 临时释放后又重新获取了缓存的锁（4764）。这是必要的因为 `itrunc` 和 `iupdate` 可能会在磁盘 `i/o` 中睡眠，但是我们必须考虑在没有锁的这段时间都发生了什么。例如，一旦 `iupdate` 结束了，磁盘上的数据结构就被标注为可用的，而并发的一个 `ialloc` 的调用就可能找到它并且重新分配它，而这一切都在 `iput` 结束之前发生。`ialloc` 会通过调用 `iget` 返回对这一块的引用，而 `iget` 此时找到了 `ip`，但看见它的 `I_BUSY` 位是设置了的，从而睡眠。现在内存中的 `i` 节点就和磁盘上的不同步了：`ialloc` 重新初始化了磁盘上的版本，但需要其他的调用者通过 `ilock` 来将它调入内存，可是 `iget` 因为现在的 `ip` 的 `I_BUSY` 位被设置而进入了睡眠。为了保证这件事发生，`iput` 必须在释放锁之前把 `I_BUSY` 和 `I_VALID` 位都清除，所以它将 `flags` 清零（4769）。

代码：`i` 节点内容

磁盘上的 `i` 节点结构，结构体 `dinode`，记录了 `i` 节点的大小和数据块的块号数组（见图6-4）。`i` 节点数据能够在 `dinode` 的 `addrs` 数组中被找到。最开始的 `NDIRECT` 个块存在于这个数组的前 `NDIRECT` 项；这些块被称作直接块。接下来的 `NINDIRECT` 个块的数据在 `i` 节点中列了出来但并没有直接存在 `i` 节点中，它们存在于一个叫做间接块的数据块中。`addrs` 数组的最后一项就是间接块的地址。因此一个文件的前 6KB（`NDIRECT * BSIZE`）个自己可以直接从 `i` 节点中取出，而后 64KB（`NINDIRECT*BSIZE`）只能在访问了间接块后取出。在磁盘上这样保存是一种比较好的表示方法，但对于用户来说显得复杂了一些。函数 `bmap` 负责这层表示使得高层的像 `readi` 和 `writeti` 这样的接口更易于编写，我们马上就会看到这一点。`bmap` 返回 `i` 节点 `ip` 中的第 `bn` 个数据块，如果 `ip` 还没有这样一个数据块，`bmap` 就会分配一个。

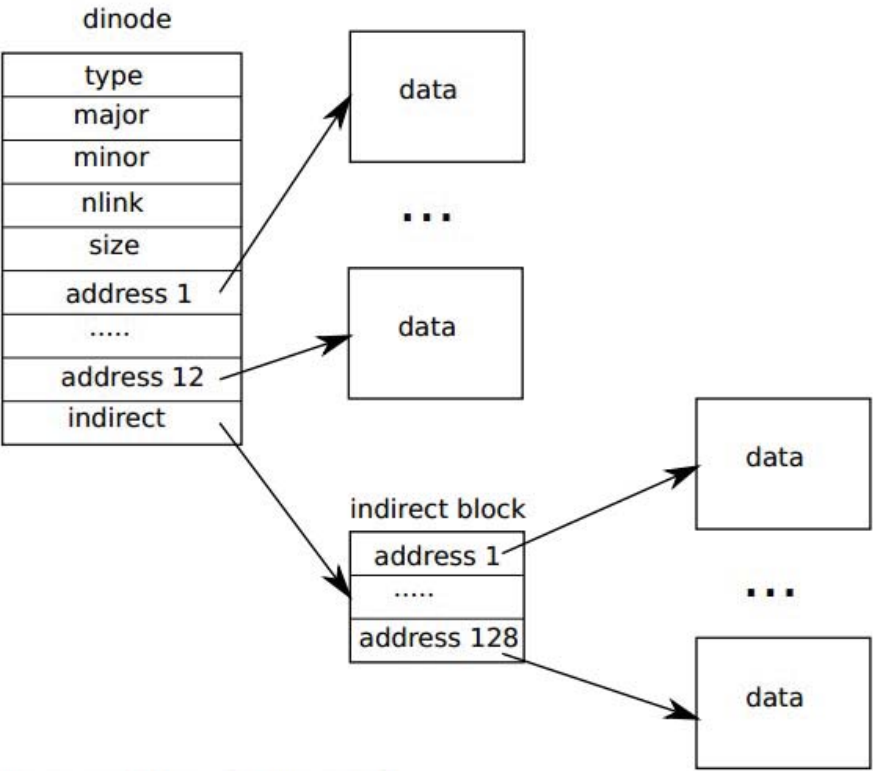


Figure 6-4. The representation of a file on disk.

函数 `bmap` (4810) 从最简单的情况开始：前面的 `NDIRECT` 个块的数据就在 `i` 节点本身中 (4815-4819)。后面的 `NINDIRECT` 个块在 `ip->addrs[NDIRECT]` 指向的间接块中 (4826)。 `bmap` 读出间接块然后再从正确的位置读出一个块号 (4827)。如果这个块号超出了 `NDIRECT+NINDIRECT`， `bmap` 就报错：调用者要负责不访问越界的块号。

当需要的时候， `bmap` 分配块。未申请的块用块号0表示。当 `bmap` 遇到0的时候，它就把它替换为新的块号 (4816-4817, 4824-4825)。

`bmap` 随着 `i` 节点的增长按需分配块，而 `itrunc` 释放它们。它把 `i` 节点的大小重新设置为0。 `itrunc` (4856) 从直接块开始释放 (4862-4867)，然后开始释放间接块中列出的块 (4872-4875)，最后释放间接块本身 (4877-4878)。

`bmap` 使得书写需要访问 `i` 节点数据流的函数变得非常轻松，比如 `readi` 和 `writei`。 `readi` (4902) 从 `i` 节点中读出数据。它最开始要保证给定的偏移和读出的量没有超出文件的末尾。从超出文件末尾的地方读会直接返回错误 (4913-4914)，如果是读的过程当中超出了文件末尾就会返回比请求的数据量少的数据 (4915-4916) (从读开始的地方到文件末尾的数据，这是所有的能返回的数据)。一个循环处理文件的每一块，从缓冲区中拷贝数据到 `dst` 中。函数 `writei` (4952) 和 `readi` 几乎是一样的，只有三个不同：1) 从文件超出文件末尾的地方开始的写或者写的过程中超出文件末尾的话会增长这个文件，直到达到最大的文件大小 (4965-4966)。2) 这个循环把数据写入缓冲区而非拷出缓冲区 (4971)。3) 如果写操作延伸了这个文件， `writei` 必须更新它的大小 (4976-4979)。

`readi` 和 `writei` 最初都会检查 `ip->type == T_DEV`。这是为了处理一些数据不存在文件系统上的特殊设备；我们会在文件描述符层重新回到这个问题。

函数 `stati` (4887) 把 `i` 节点的元数据拷贝到 `stat` 结构体中，这个结构体可通过系统调用 `stat` 暴露给用户程序。

代码：目录层

在xv6中，目录的实现和文件的实现过程很像。目录的 `i` 节点的类型 `T_DIR`，它的数据是一系列的目录条目。每个条目是一个 `struct dirent` (3700)结构体，包含一个名字和一个 `i` 节点编号。这个名字最多有 `DIRSIZ` (14)个字符；如果比较短，它将以 `NUL` (0) 作为结尾字符。`i` 节点编号是0的条目都是可用的。

函数 `dirlookup` (5011) 用于查找目录中指定名字的条目。如果找到，它将返回一个指向相关 `i` 节点的指针，未被锁的，并设置目录内条目的字节偏移 `*poff`，从而用于调用者编辑。如果 `dirlookup` 正好找到对应名字的一个条目，它会更新 `*poff`，释放块，并通过 `iget` 返回一个未被锁的 `i` 节点。`dirlookup` 是 `iget` 返回未被锁的 `i` 节点的原因。调用者已经对 `dp` 上锁，如果查找 `.`，当前目录的别名，并在返回之前尝试去锁上该 `i` 节点会导致二次锁上 `dp` 的可能并造成死锁。（`.` 不是唯一的问题，还有很多更复杂的死锁问题，比如多进程和 `..`，当前目录的父目录的别名。）调用者会先解锁 `dp` 然后锁上 `ip`，以保证每次只能维持一个锁。

函数 `dirlink` (5052) 会写入一个新的目录条目到目录 `dp` 下，用指定的名字和 `i` 节点编号。如果这个名字已经存在，`dirlink` 会返回 `error` (5058-5062)。循环的主体是读取目录条目以查找尚未分配的条目。当找到一个未分配的条目后，循环会提前退出 (5022-5023)，并设置可用条目的偏移 `off`。否则，循环以设置 `dp->size` 的 `off` 方式结束。无论哪种方式，`dirlink` 之后会通过写入偏移 `off` 来给目录添加一个新的条目。

代码：路径名

路径名查询会对每一个路径的每一个元素调用 `dirlookup`。`namei` (5189) 解析 `path` 并返回对应的 `i` 节点。函数 `nameiparent` 是一个变种；它在最后一个元素之前停止，返回上级目录的 `i` 节点并且把最后一个元素拷贝到 `name` 中。这两个函数都使用 `namex` 来实现。

`namex` (5154) 会计算路径解析从什么地方开始。如果路径以反斜杠开始，解析就会从根目录开始；其他情况下则会从当前目录开始 (5161)。然后它使用 `skipelem` 来依次考虑路径中的每一个部分 (5163)。每一次循环迭代都必须在当前的 `i` 节点 `ip` 中找 `name`。循环的开始会把 `ip` 锁住，然后检查它是否确实是一个目录。如果它不是目录的话，查询就宣告失败。（锁住 `ip` 是必须的，不是因为 `ip->type` 随时有可能变（事实上它不会变），而是因为如果不调用 `ilock` 的话，`ip->type` 可能还没有从磁盘加载出来）。如果是调用 `nameiparent` 而且这是最后一个路径元素，那么循环就直接结束了。因为最后一个路径元素已经拷贝到了 `name` 中，所以 `namex` 只需要返回解锁的 `ip` (5169-5173)。最后，循环用 `dirlookup` 寻找路径元素并且令 `ip=next`，准备下一次的循环 (5174-5179)。当循环处理了每一个路径元素后，它返回 `ip`。

文件描述符层

UNIX 接口很爽的一点就是大多数的资源都可以用文件来表示，包括终端这样的设备、管道，当然，还有真正的文件。文件描述符层就是实现这种统一性的一层。

xv6 给每个进程都有自己的打开文件表，正如我们在第零章中所见。每一个打开文件都由结构体 `file` (3750) 表示，它是一个对 `i` 节点或者管道和文件偏移的封装。每次调用 `open` 都会创建一个新的打开文件（一个新的 `file` 结构体）。如果多个进程相互独立地打开了同一个文件，不同的实例将拥有不同的 `i/o` 偏移。另一方面，同一个文件可以（同一个 `file` 结构体）可以在一个进程的文件表中多次出现，同时也可以多个进程的文件表中出现。当一个进程用 `open` 打开了一个文件而后使用 `dup`，或者把这个文件和子进程共享就会导致这一点发生。对每一个打开的文件都有一个引用计数，一个文件可以被打开用于读、写或者二者。`readable` 域和 `writable` 域记录这一点。

系统中所有的打开文件都存在于一个全局的文件表 `ftable` 中。这个文件表有一个分配文件的函数（`filealloc`），有一个重复引用文件的函数（`filedup`），释放对文件引用的函数（`fileclose`），读和写文件的函数（`fileread` 和 `filewrite`）。

前三个的形式我们已经很熟悉了。`filealloc` (5225) 扫描整个文件表来寻找一个没有被引用的文件（`file->ref == 0`）并且返回一个新的引用；`filedup` (5252) 增加引用计数；`fileclose` (5264) 减少引用计数。当一个文件的引用计数变为 0 的时候，`fileclose` 就会释放掉当前的管道或者 `i` 节点（根据文件类型的不同）。

函数 `filestat`，`fileread`，`filewrite` 实现了对文件的 `stat`，`read`，`write` 操作。`filestat` (5302) 只允许作用在 `i` 节点上，它通过调用 `stat` 实现。`fileread` 和 `filewrite` 检查这个操作被文件的打开属性所允许然后把执行让渡给 `i` 节点的实现或者管道的实现。如果这个文件代表的是一个 `i` 节点，`fileread` 和 `filewrite` 就会把 `i/o` 偏移作为该操作的偏移并且往前移 (5325-5326, 5365-5366)。管道没有偏移这个概念。回顾一下 `i` 节点的函数需要调用者来处理锁 (5305-5307, 5324-5327, 5364-5378)。`i` 节点锁有一个方便的副作用那就是读写偏移会自动更新，所以同时对一个文件写并不会覆盖各自的文件，但是写的顺序是不被保证的，因此写的结果可能是交织的（在一个写操作的过程中插入了另一个写操作）。

代码：系统调用

有了底层的这些函数，大多数的系统调用的实现都是很简单的（参见 `sysfile.c`）。还有少数几个调用值得一说。

函数 `sys_link` 和 `sys_unlink` 修改目录文件，可能会创建或者移除对 `i` 节点的引用。它们是使用会话的另一个佳例。`sys_link`（5513）最开始获取自己的参数 `old` 和 `new` 两个字符串。假设 `old` 是存在的并且不是一个目录文件（5520-5530），`sys_link` 增加它的 `ip->nlink` 计数。然后 `sys_link` 调用 `nameiparent(new)` 来寻找上级目录和最终的目录元素（5536），并且创建一个目录项指向 `old` 的 `i` 节点（5539）。`new` 的上级目录必须和已有 `old` 的 `i` 节点在同一个设备上；`i` 节点号只在同一个磁盘上有意义。如果这样的错误发生了，`sys_link` 必须回溯并且还原引用计数。

`sys_link` 为一个已有的 `i` 节点创建一个新的名字。而函数 `create`（5657）为一个新的 `i` 节点创建新名字。它是三个文件创建系统调用的综合：用 `O_CREATE` 方式 `open` 一个文件创建一个新的普通文件，`mkdir` 创建一个新的目录文件，`mkdev` 创建一个新的设备文件。就像 `sys_link` 一样，`create` 调用 `nameiparent` 获取上级目录的 `i` 节点。然后调用 `dirlookup` 来检查同名文件是否已经存在（5667）。如果的确存在，`create` 的行为就由它服务的系统调用所决定，`open` 和 `mkdir` 以及 `mkdev` 的语义是不同的。如果是 `open`（`type==T_FILE`）调用的 `create` 并且按指定文件名找到的文件是一个普通文件，那么就认为打开成功，因此 `create` 中也认为是成功。在其他情况下，这就是一个错误（5672-5673）。如果文件名并不存在，`create` 就会用 `ialloc`（5676）分配一个新的 `i` 节点。如果新的 `i` 节点是一个目录，`create` 就会初始化 `.` 和 `..` 两个目录项。最后所有的数据都初始化妥当了，`create` 就可以把它连接到它的上级目录（5689）。`create`，正如 `sys_link` 一样，同时拥有两个 `i` 节点锁：`ip` 和 `dp`。这不可能导致死锁，因为 `i` 节点 `ip` 是刚被分配的：系统中没有其他进程会持有 `ip` 的锁并且尝试锁 `dp`。

使用 `create`，就能轻易地实现 `sys_open` 和 `sys_mkdir`，以及 `sys_mknod`。`sys_open` 是最复杂的，创建一个新文件只是他能做的很少一部分事。如果 `open` 以 `O_CREATE` 调用，它就会调用 `create`（5712）。否则，它就会调用 `namei`（5717）。`create` 会返回一个带锁的 `i` 节点，但是 `namei` 并不会，所以 `sys_open` 必须要自己锁上这个 `i` 节点。这样提供了一个合适的地方来检查目录只被打开用于读，而不是写。总之我们获得了一个 `i` 节点（不管是用 `create` 还是用 `namei`），`sys_open` 分配了一个文件和文件描述符（5726），接着填充了这个文件（5734-5738）。我们要记住没有其他进程能够访问初始化尚未完成的文件，因为他只存在于当前进程的文件表中。

第五章研究了管道的实现，在那时我们甚至还没有一个文件系统。函数 `sys_pipe` 通过管道对的方式把管道的实现和文件系统连接起来。它的参数是一个指向可装入两个整数的数组指针，这个数组将用于记录两个新的文件描述符。然后它分配管道，将新的文件描述符存入这个数组中。

现实情况

现实情况中操作系统的块缓冲比xv6中要复杂的多，但依旧为两个主要目的服务：缓冲和同步到磁盘。xv6的块缓冲，与V6的类似，使用简单的近期最少使用算法（Least Recently Used, or LRU）的回收策略；当然现实中还有很多其他相对复杂的策略，每种策略都在某些方面有优势而在其他方面有劣势。一种更有效的LRU缓冲策略会减少链表的使用，它将用哈希表来实现查找、用堆来实现LRU的回收。现代块缓冲一般会结合虚拟内存系统从而可以支持映射到内存的文件。

xv6的日志系统相当不高效。它不支持并发的、可更新的系统调用，即使当系统调用命令在文件系统完全不同的部分间执行。它要记录整个块，即使一个块只有很少一部分字节改变。它可以实现日志记录的同步，但每次只有一个块，每次都可能需要整个磁盘的运转时间。现实的日志系统可以处理这些问题。

日志记录不是唯一的崩溃后的恢复机制，早期文件系统在重新启动的时候使用清扫机制（比如，UNIX系统中的 `fsck` 命令）来检查每个文件和目录以及各个块和 `i` 节点可用的链表，查找并解决出现的不一致问题。在大型文件系统中清扫机制可能耗费数小时，并且存在一些情况，无法猜测出正确解决不一致问题的方案。相对来说，通过日志系统来恢复更快捷准确。

xv6使用基于磁盘分层结构的 `i` 节点和目录，和早期的UNIX系统相同；这种机制直到最近这些年仍旧保留。BSD的UFS/FFS和Linux的ext2/ext3都使用了类似的数据结构。文件系统结构最不高效的部分是目录，每次查找都需要在磁盘块上线性查找。当目录只存在几个磁盘块上时，这样是合理的，但对于有很多文件的目录来说，这样不高效。微软Windows的NTFS，Mac OS X的HFS，以及Solaris的ZFS,这些文件系统，仅仅命名了一些目录，把目录实现为磁盘上块的平衡树结构。这样虽然复杂，但保证了log级别时间复杂度的查找。

xv6处理磁盘任务失败的方式是too simple, sometimes naive：如果磁盘操作失败，xv6报警。这种方式是否有效要取决于硬件：如果一个操作系统是在某些特殊的硬件上方，该硬件使用冗余来掩饰磁盘错误，那么或许该操作系统不会经常出错，报警机制是有效的。另一方面，如果操作系统是用普通平凡的硬件，则可能会遇到更多操作失败，需要经常处理，才能够在—

个文件中的块丢失情况下不会影响到文件系统其他的使用。

xv6需要文件系统固定在一个磁盘设备上，大小不改变。由于大型数据库和多媒体文件需要更大的储存容量，操作系统需要攻破“一个磁盘一个文件系统”的瓶颈。基础的方法是把许多磁盘联合到一个逻辑磁盘上。硬件解决方案如RAID依旧是最流行的，但目前的趋势是在软件上尽可能地实现这种逻辑。这种软件的实现可以允许很多丰富的功能，如通过快速添加和移除命令来实现逻辑设备的增长和缩减。当然，存储层的快速增加和减少需要类似机制的文件系统：UNIX文件系统使用的i节点块结构是固定空间大小的数组结构，它们不能很好的实现上述功能。把文件系统和磁盘管理系统分开也许是最干净的设计，但是两个系统的复杂接口产生了新的系统，如Sun的ZFS，可以把它们联系起来。

xv6文件系统缺少很多现今其他文件系统的特征；比如，缺乏对备份的快照和增加机制的支持。

xv6有两种不同的文件实现方式：管道和i节点。现代的UNIX系统有很多方式：管道，网络连接，和许多来自不同类型文件系统的i节点，包括网络文件系统。与在 `fileread` 和 `filewrite` 中 `if` 声明不同的是，这些系统一般给每个打开的文件一个由函数指针构成的表，每次操作之后，调用函数指针来实现相关i节点的调用。网络文件系统和用户层面的文件系统提供了这些函数，调用到网络的RPC并在返回前等待答复。

练习

1.为什么在 `balloc` 中会 `panic` ？我们能恢复吗？

2.为什么在 `ialloc` 中会 `panic` ？我们能恢复吗？

3.i节点产生编号

4.为什么当 `filealloc` 用完文件的时候不会 `panic` ？为什么这是很普通但是值得处理的？

5.假设 `ip` 对应的文件在 `sys_link` 调用到 `iunlock(ip)` 和 `dirlink` 时被另一个进程去除对应的联系，这个联系还会被正确的创建吗？并说明原因。

6. `create` 有四个函数调用（其中一个是调用到 `ialloc`，其他三个调用到 `dirlink`）需要成功创建。如果不能，`create` C产生 `panic`。为什么这是可取的？为什么这四个调用中的任意一个都不会失败？

7. `sys_chdir` 在调用 `iput(cp->cwd)` 之前调用 `iunlock(ip)`，这个过程可能会锁住 `cp->cwd`，然而推迟 `iunlock(ip)` 调用直到 `iput` 调用结束，将不会造成死锁，为什么不会？

附录 A

PC 硬件

本文介绍供 x86 运行的个人计算机(PC)硬件平台。

PC 是指遵守一定工业标准的计算机，它的目标是使得不同厂家生产的机器都能够运行一定范围内的软件。这些标准随时间迁移不断变化，因此90年代的 PC 与今日的 PC 看起来已是大不相同。

从外观来看，PC 是一个配置有键盘、屏幕和各种设备的“盒子”。盒子内部则是一块集成电路——主板，上面有 CPU 芯片，内存芯片，显卡芯片，I/O 控制器芯片，以及负责芯片间通信的总线。总线会遵守某种标准（如 PCI 或 USB），从而能够兼容不同厂家的设备。

我们可以把 PC 抽象为三部分：CPU、内存和 I/O 设备。CPU 负责计算，内存用于保存计算用的指令和数据，其他设备用于实现存储、通讯等其他功能。

你可以想象主存以一组导线与 CPU 相连接，其中一些是地址线，一些是数据线，还有一些则是控制线。CPU 要从主存读出一个值，需要向地址线上输出一系列表示0和1的电压，并在规定的时间内在“读”线上发出信号1，接下来再从数据线上的高低电压中获取数据。CPU 若要向内存中写入一个值，则向数据线和地址线上写入合适的值，并在规定时间内在“写”位上发出信号1。真实的内存接口比这复杂的多，但除非你是在追求高性能，否则你不必考虑这么多的细节。

处理器和内存

CPU（中央处理单元，或处理器）其实只是在执行一个非常简单的循环：从一个被称为『程序计数器』的寄存器中获取一个内存地址，从这个地址读出机器指令，增加程序计数器的值，执行机器指令，不断反复。某些机器指令如分支和函数调用会改变程序计数器，如果执行机器指令没有改变程序计数器，这个循环就会从程序计数器开始一条一条地执行指令。

如果不能保存和修改程序数据，那么执行指令就是毫无意义的。速度最快的数据存储单元是处理器的寄存器组。一个寄存器是处理器内的一个存储单元，能够保存一个字大小的值（按照机器不同，一个字通常会16, 32或者64位）。寄存器内的值能在一个 CPU 周期内被快速地读写。

PC 处理器实现了 x86 指令集，该指令集由 Intel 发布并成为了一种标准，一些厂商生产实现了该指令集的处理器。和其他的 PC 标准一样，这个标准也在不断更新，但是新的标准是向前兼容的。由于 PC 处理器启动时都是模拟1981年 IBM PC 上使用的芯片 Intel 8088，所以 boot loader 需要作出改变以应对标准的更新。但是，对于 xv6 的绝大部分内容，你只需要关心现代 x86 指令集。

现代 x86 提供了8个32位通用寄存器-- %eax, %ebx, %ecx, %edx, %edi, %esi, %ebp, %esp 和一个程序计数器 %eip (*instruction pointer*)。前缀e是指扩展的 (*extended*)，表示它们是16位寄存器 %ax, %bx, %cx, %dx, %di, %si, %bp, %sp 的32位扩展。这两套寄存器其实是相互的别名，例如 %ax 是 %eax 的低16位：我们在写 %ax 的时候也会改变 %eax，反之亦然。前四个寄存器的两个低8位还有自己的名字：%al, %ah 分别表示 %ax 的低8位和高8位，%bl, %bh, %cl, %ch, %dl, %dh 同理。另外，x86 还有8个80位的浮点寄存器，以及一系列特殊用途的寄存器如控制寄存器 %cr0, %cr2, %cr3, %cr4，调试寄存器 %dr0, %dr1, %dr2, %dr3；段寄存器 %cs, %ds, %es, %fs, %gs, %ss；还有全局和局部描述符表的伪寄存器 %gdtr, %ldtr。控制寄存器和段寄存器对于任何操作系统都是非常重要的。浮点寄存器和调试寄存器则没那么有意思，并且也没有在 xv6 中使用。

寄存器非常快但是也非常昂贵。大多数处理器都会提供至多数十个通用寄存器。下一个层次的存储器是随机存储器（RAM）。主存的速度大概比寄存器慢10到100倍，但要便宜得多，所以容量可以更大。主存较慢的一个原因是它不在处理器芯片上。一个 x86 处理器只有十多个寄存器，但今天的 PC 通常有 GB 级的主存。由于寄存器和主存在读写速度和大小上的巨大差异，大多数处理器，包括 x86，都在芯片上的缓存中保存了最近使用的主存数据。缓存是主存和寄存器在速度和大小上的折衷。现在的 x86 处理器通常有二级缓存，第一级较小，读写速率接近处理器的时钟周期，第二级较大，读写速率在第一级缓存和主存之间。下表显示了 Intel Core 2 Duo 系统的实际数据：

Intel Core 2 Duo E7200 at 2.53 GHz 备忘：换上真实数字！| 存储器 | 读写时间 | 大小 | |-----|-----|-----| 寄存器 | 0.6ns | 64

字节| |L1缓存|0.5ns|64K 字节| |L2缓存|10ns|4M 字节| |主存|100ns|4G 字节|

通常 x86 对操作系统隐藏了缓存，所以我们只需要考虑寄存器和主存两种存储器，不用担心主存的层次结构引发的差异。

I/O

处理器必须像和主存交互一样同设备交互。x86 处理提供了特殊的 `in, out` 指令来在设备地址（称为'I/O 端口'）上读写。这两个指令的硬件实现本质上和读写内存是相同的。早期的 x86 处理器有一条附加的地址线：0 表示从 I/O 端口读写，1 则表示从主存读写。每个硬件设备会处理它所在 I/O 端口所接收到的读写操作。设备的端口使得软件可以配置设备，检查状态，使用设备；例如，软件可以通过对 I/O 端口的读写，使磁盘接口硬件对磁盘扇区进行读写。

很多计算机体系结构都没有单独的设备访问指令，取而代之的是让设备拥有固定的内存地址，然后通过内存读写实现设备读写。实际上现代 x86 体系结构就在大多数高速设备上（如网络、磁盘、显卡控制器）使用了该技术，叫做内存映射 I/O。但由于向前兼容的原因，`in, out` 指令仍能使用，而比较老的设备如 xv6 中使用的 IDE 磁盘控制器仍使用两个指令。

附录 B

引导加载器（boot loader）

当 x86 PC 启动时，它执行的是一个叫 BIOS 的程序。BIOS 存放在非易失存储器中，BIOS 的作用是在启动时进行硬件的准备工作，接着把控制权交给操作系统。具体来说，BIOS 会把控制权交给从引导扇区（用于引导的磁盘的第一个 512 字节的数据区）加载的代码。引导扇区中包含引导加载器——负责内核加载到内存中。BIOS 会把引导扇区加载到内存 0x7c00 处，接着（通过设置寄存器 `%ip`）跳转至该地址。引导加载器开始执行后，处理器处于模拟 Intel 8088 处理器的模式下。而接下来的工作就是把处理器设置为现代的操作模式，并从磁盘把 xv6 内核载入到内存中，然后将控制权交给内核。xv6 引导加载器包括两个源文件，一个由 16 位和 32 位汇编混合编写而成（`bootasm.s`；(8400)），另一个由 C 写成（`bootmain.c`；(8500)）。

代码：汇编引导程序

引导加载器的第一条指令 `cli` (8412) 屏蔽处理器中断。硬件可以通过中断触发中断处理程序，从而调用操作系统的功能。BIOS 作为一个小型操作系统，为了初始化硬件设备，可能设置了自己的中断处理程序。但是现在 BIOS 已经没有了控制权，而是引导加载器正在运行，所以现在还允许中断不合理也不安全。当 xv6 准备好了后（详见第 3 章），它会重新允许中断。

现在处理器处在模拟 Intel 8088 的实模式下，有 8 个 16 位通用寄存器可用，但实际上处理器发送给内存的是 20 位的地址。这时，多出来的 4 位其实是由段寄存器 `%cs`，`%ds`，`%es`，`%ss` 提供的。当程序用到一个内存地址时，处理器会自动在该地址上加上某个 16 位段寄存器值的 16 倍。因此，内存引用中其实隐含地使用了段寄存器的值：取指令会用到 `%cs`，读写数据会用到 `%ds`，读写栈会用到 `%ss`。

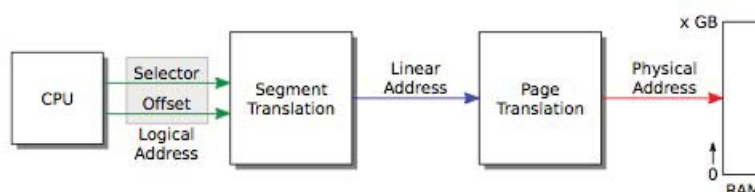


Figure B-1. The relationship between logical, linear, and physical addresses.

xv6 假设 x86 指令在做内存操作时使用的是虚拟地址，但实际上 x86 指令使用的是逻辑地址（见表 B-1）。逻辑地址由段选择器和偏移组成，有时又被写作 `segment:offset`。更多时候，段是隐含的，所以程序会直接使用偏移。分段硬件会完成上述处理，从而产生一个线性地址。如果允许分页硬件工作（见第 2 章），分页硬件则会把线性地址翻译为物理地址；否则处理器直接把线性地址看作物理地址。

引导加载器还没有允许分页硬件工作；它通过分段硬件把逻辑地址转化为线性地址，然后直接作为物理地址使用。xv6 会配置分段硬件，使之不对逻辑地址做任何改变，直接得到线性地址，所以线性地址和逻辑地址是相等的。由于历史原因我们用虚拟地址这个术语来指程序操作时用的地址。xv6 的虚拟地址等于 X86 的逻辑地址，同样也等于分段硬件映射的线性地址。等到开启了分页后，系统中值得关心的就只有从线性地址到物理地址的映射。

BIOS 完成工作后，`%ds`，`%es`，`%ss` 的值是未知的，所以在屏蔽中断后，引导加载器的第一个工作就是将 `%ax` 置零，然后把这个零值拷贝到三个段寄存器中（8415-8418）。

虚拟地址 `segment:offset` 可能产生 21 位物理地址，但 Intel 8088 只能向内存传递 20 位地址，所以它截断了地址的最高位：`0xffff0 + 0xffff = 0x10fff`，但在 8088 上虚拟地址 `0xffff:0xffff` 则是引用物理地址 `0x0fff`。早期的软件依赖硬件来忽略第 21 位地址位，所以当 Intel 研发出使用超过 20 位物理地址的处理器时，IBM 就想出了一个技巧来保证兼容性。那就是，如果键盘控制器输出端口的第 2 位是低位，则物理地址的第 21 位被清零；否则，第 21 位可以正常使用。引导加载器用 I/O 指令控制端口 `0x64` 和 `0x60` 上的键盘控制器，使其输出端口的第 2 位为高位，来使第 21 位地址正常工作（8436）。

对于使用内存超过65536字节的程序而言，实模式的16位寄存器和段寄存器就显得非常困窘了，显然更不可能使用超过 1M 字节的内存。x86系列处理器在80286之后就有了保护模式。保护模式下可以使用更多位的地址，并且（80386之后）有了“32位”模式使得寄存器，虚拟地址和大多数的整型运算都从16位变成了32位。xv6 引导程序依次允许了保护模式和32位模式。

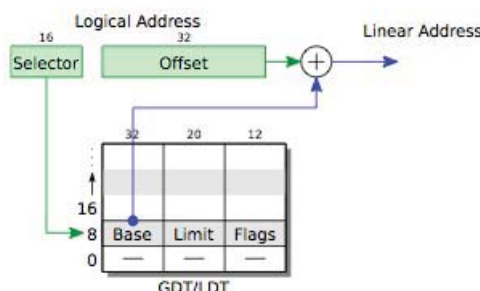


Figure B-2. Segments in protected mode.

在保护模式下，段寄存器保存着段描述符表的索引（见图表 B-2）。每一个表项都指定了一个基物理地址，最大虚拟地址（称为限制），以及该段的权限位。这些权限位在保护模式下起着保护作用，内核可以根据它们来保证一个程序只使用属于自己的内存。

xv6 几乎没有使用段；取而代之的是第2章讲述的分页。引导加载器将段描述符表 `gdt`（8482-8485）中的每个段的基址都置零，并让所有段都有相同的内存限制（4G字节）。该表中有一个空指针表项，一个可执行代码的表项，一个数据的表项。代码段描述符的标志位中指示了代码只能在32位模式下执行（0660）。正是由于这样的设置，引导加载器在进入保护模式时，逻辑地址才会直接映射为物理地址。

引导加载器执行 `lgdt`（8441）指令来把指向 `gdt` 的指针 `gdt_desc`（8487-8489）加载到全局描述符表（GDT）寄存器中。

加载完毕后，引导加载器将 `%cr0` 中的 `CR0_PE` 位置为1，从而开启保护模式。允许保护模式并不会马上改变处理器把逻辑地址翻译成物理地址的过程；只有当某个段寄存器加载了一个新的值，然后处理器通过这个值读取 GDT 的一项从而改变了内部的段设置。我们没法直接修改 `%cs`，所以使用了一个 `ljmp` 指令（8453）。跳转指令会接着在下一行（8456）执行，但这样做实际上将 `%cs` 指向了 `gdt` 中的一个代码描述符表项。该描述符描述了一个32位代码段，这样处理器就切换到了32位模式下。就这样，引导加载器让处理器从8088进化到80286，接着进化到了80386。

在32位模式下，引导加载器首先用 `SEG_KDATA`（8458-8461）初始化了数据段寄存器。逻辑地址现在是直接映射到物理地址的。运行 C 代码之前的最后一个步骤是在空闲内存中建立一个栈。内存 `0xa0000` 到 `0x100000` 属于设备区，而 xv6 内核则是放在 `0x100000` 处。引导加载器自己是在 `0x7c00` 到 `0x7d00`。本质上来讲，内存的其他任何部分都能用来存放栈。引导加载器选择了 `0x7c00`（在该文件中即 `$start`）作为栈顶；栈从此处向下增长，直到 `0x0000`，不断远离引导加载器代码。

最后加载器调用 C 函数 `bootmain`（8468）。`bootmain` 的工作就是加载并运行内核。只有在出错时该函数才会返回，这时它会向端口 `0x8a00`（8470-8476）输出几个字。在真实硬件中，并没有设备连接到该端口，所以这段代码相当于什么也没有做。如果引导加载器是在 PC 模拟器上运行，那么端口 `0x8a00` 则会连接到模拟器并把控制权交还给模拟器本身。无论是否使用模拟器，这段代码接下来都会执行一个死循环（8477-8478）。而一个真正的引导加载器则应该会尝试输出一些调试信息。

代码：C 引导程序

引导加载器的 C 语言部分 `bootmain.c`（8500）目的是在磁盘的第二个扇区开头找到内核程序。如我们在第2章所见，内核是 ELF 格式的二进制文件。为了读取 ELF 头，`bootmain` 载入 ELF 文件的前4096字节（8514），并将其拷贝到内存中 `0x10000` 处。

下一步要通过 ELF 头检查这是否的确是一个 ELF 文件。`bootmain` 从磁盘中 ELF 头之后 `off` 字节处读取扇区的内容，并写到内存中地址 `paddr` 处。`bootmain` 调用 `readseg` 将数据从磁盘中载入（8538），并调用 `stosb` 将段的剩余部分置零（8540）。`stosb`（0492）使用 x86 指令 `rep stsb` 来初始化内存块中的每个字节。

在内核编译和链接后，我们应该能在虚拟地址 0x80100000 处找到它。因此，函数调用指令使用的地址都是 0xf01xxxx 的形式；你可以在 `kernel.asm` 中找到类似的例子。这个地址是在 `kernel.ld` 中设置的。0x80100000 是一个比较高的地址，差不多处在32位地址空间的尾部；至于原因，我们在第2章中对此作出了详细解释。当然，实际的物理内存中可能并没有这么高的地址。一旦内核开始运行，它会开启分页硬件来将虚拟地址 0x80100000 映射到物理地址 0x00100000。引导程序运行到现在，分页机制尚未被开启。在 `kernel.ld` 中指明了内核的 `paddr` 是 0x00100000，也就是说，引导加载器将内核拷贝到的低地址正是分页硬件最终会映射的物理地址。

引导加载器的最后一项工作是调用内核的入口指令，即内核第一条指令的执行地址。在 xv6 中入口指令的地址是 0x10000c：

```
# objdump -f kernel

kernel:      file format elf32-i386
architecture: i386, flags 0x00000112:
EXEC_P, HAS_SYMS, D_PAGED
start address 0x0010000c
```

按照惯例，在 `entry.S`（1036）中定义的 `_start` 符号即 ELF 入口。由于 xv6 还没有建立虚拟内存，xv6 的入口即 `entry`（1040）的物理地址。

现实情况

该附录中谈到的引导加载器编译后大概有470字节的机器码，具体大小取决于编译优化。为了放入比较小的空间中，xv6 引导加载器做了一个简单的假设：内核放在引导磁盘中从扇区1开始的连续空间中。通常内核就放在普通的文件系统中，而且可能不是连续的。也有可能内核是通过网络加载的。这种复杂性就使得引导加载器必须要能够驱动各种磁盘和网络控制器，并能够解读不同的文件系统和网络原型。也就是说，引导加载器本身就已经成为了一个小操作系统。显然这样的引导加载器不可能只有512字节，大多数的 PC 操作系统的引导过程分为2步。首先，一个类似于该附录介绍的简单的引导加载器会从一个已知的磁盘位置上把完整的引导加载器加载进来，通常这一步会依靠空间权限更大的 BIOS 来操作磁盘。接下来，这个超过512字节的完整加载器就有足够的能力定位、加载并执行内核了。也许在更现代的设计中，会直接用 BIOS 从磁盘中读取完整的引导加载器（并在保护模式和32位模式下启动之）。

本文假设在开机后，引导加载器运行前，唯一发生的事即 BIOS 加载引导扇区。但实际上 BIOS 会做相当多的初始化工作来确保现代计算机中结构复杂的硬件能像传统标准中的 PC 一样工作。

练习

1. 基于扇区大小，文中提到的调用 `readseg` 的作用和 `readseg((uchar*)0x100000, 0xb500, 0x1000)` 的作用是相同的。实际上，这个草率的实现并不会导致错误。这是为什么呢？
2. 一些关于 BIOS 存在时长与安全性的问题。
3. 假设你希望 `bootmain()` 能把内核加载到 0x200000 而非 0x100000，于是你在 `bootmain()` 中把每个 ELF 段的 `va` 都加上了 0x100000。这样做是会导致错误发生的，请说明会发生什么错误。
4. 引导加载器把 ELF 头拷贝到了一个随意的地址 0x10000 上，这样做看起来似乎有些危险。那么为什么不调用 `malloc` 来分配它所需要的空间呢？

词汇表

统一特殊词汇的翻译标准，请按字典序记录。

a

| 英文 | 中文 |
|---------------|------|
| address space | 地址空间 |
| allocator | 分配器 |

b

| 英文 | 中文 |
|--------------|-------------|
| boot loader | 引导加载器 |
| buffer cache | 缓冲器高速缓存，块缓冲 |

c

| 英文 | 中文 |
|------------------------|-----------|
| (sleep/wakeup) channel | (睡眠/唤醒)队列 |
| console | 控制台 |
| context | 上下文 |
| convoy | 护航 |
| coroutine | 共行程序 |

d

| 英文 | 中文 |
|------|----|
| disk | 磁盘 |

e

f

| 英文 | 中文 |
|-----------|------|
| free list | 空闲链表 |

g

| 英文 | 中文 |
|------------|-----|
| guard page | 保护页 |

h

| 英文 | 中文 |
|--------------------|------|
| hardware privilege | 硬件特权 |

i

| 英文 | 中文 |
|-------|------|
| inode | i 节点 |

j

k

| 英文 | 中文 |
|--------|----|
| kernel | 内核 |

l

| 英文 | 中文 |
|-----------------|-------|
| logical address | 逻辑地址 |
| lost wakeup | 丢失的唤醒 |

m

| 英文 | 中文 |
|-------------|------|
| multiplex | 多路复用 |
| motherboard | 主板 |

n

o

| 英文 | 中文 |
|-----------|-----|
| on-disk | 磁盘上 |
| in-memory | 内存中 |

p

| 英文 | 中文 |
|-------------------------|------------|
| page table | 页表 |
| page table entry | 页表条目 |
| paging hardwre | 分页硬件 |
| per-cpu | per-cpu |
| pipe | 管道 |
| pre-arranged | 预定义 |
| previlege level | 特权级 |
| PC-relative direct jump | PC 相关的直接跳转 |

q

r

| 英文 | 中文 |
|------------------|------|
| real world | 现实情况 |
| rendezvous point | 集合点 |

s

| 英文 | 中文 |
|------------------|-------|
| segment hardware | 分段硬件 |
| semaphore | 信号量 |
| shell | shell |
| system call | 系统调用 |

t

| 英文 | 中文 |
|-----------------|----|
| thundering herd | 惊群 |
| time-sharing | 分时 |

u

| 英文 | 中文 |
|----------------|-------|
| unwind (stack) | 展开（栈） |

v

w

x

y

z

原文档错误

1. 50页：“To avoid this situation ...”段中，“(cli is the x86 ...”缺少后括号。
2. 65页：最后一段倒数第4行，“process may send a signal to it.”排版有误。
3. 51页：“For example ...”段中，倒数第4行的“than”应为“then”。

xv6 is a re-implementation of Dennis Ritchie's and Ken Thompson's Unix Version 6 (v6). xv6 loosely follows the structure and style of v6, but is implemented for a modern x86-based multiprocessor using ANSI C.

ACKNOWLEDGMENTS

xv6 is inspired by John Lions's Commentary on UNIX 6th Edition (Peer to Peer Communications; ISBN: 1-57398-013-7; 1st edition (June 14, 2000)). See also <http://pdos.csail.mit.edu/6.828/2012/v6.html>, which provides pointers to on-line resources for v6.

xv6 borrows code from the following sources:

- JOS (asm.h, elf.h, mmu.h, bootasm.S, ide.c, console.c, and others)
- Plan 9 (entryother.S, mp.h, mp.c, lapic.c)
- FreeBSD (ioapic.c)
- NetBSD (console.c)

The following people have made contributions:

- Russ Cox (context switching, locking)
- Cliff Frey (MP)
- Xiao Yu (MP)
- Nickolai Zeldovich
- Austin Clements

In addition, we are grateful for the patches contributed by Greg Price, Yandong Mao, and Hitoshi Mitake.

The code in the files that constitute xv6 is Copyright 2006-2012 Frans Kaashoek, Robert Morris, and Russ Cox.

ERROR REPORTS

If you spot errors or have suggestions for improvement, please send email to Frans Kaashoek and Robert Morris (kaashoek,rtm@csail.mit.edu).

BUILDING AND RUNNING XV6

To build xv6 on an x86 ELF machine (like Linux or FreeBSD), run "make". On non-x86 or non-ELF machines (like OS X, even on x86), you will need to install a cross-compiler gcc suite capable of producing x86 ELF binaries. See <http://pdos.csail.mit.edu/6.828/2012/tools.html>. Then run "make TOOLPREFIX=i386-jos-elf-".

To run xv6, you can use the Bochs or QEMU PC simulators. Bochs makes debugging easier, but QEMU is much faster. To run in Bochs, run "make bochs" and then type "c" at the bochs prompt. To run in QEMU, run "make qemu".

To create a typeset version of the code, run "make xv6.pdf". This requires the "mpage" utility. See <http://www.mesa.nl/pub/mpage/>.

The numbers to the left of the file names in the table are sheet numbers. The source code has been printed in a double column format with fifty lines per column, giving one hundred lines per sheet (or page). Thus there is a convenient relationship between line numbers and sheet numbers.

| # basic headers | # system calls | # string operations |
|-----------------|----------------|----------------------|
| 01 types.h | 29 traps.h | 61 string.c |
| 01 param.h | 29 vectors.pl | |
| 02 memlayout.h | 30 trapasm.S | # low-level hardware |
| 02 defs.h | 30 trap.c | 63 mp.h |
| 04 x86.h | 32 syscall.h | 64 mp.c |
| 06 asm.h | 32 syscall.c | 66 lapic.c |
| 07 mmu.h | 34 sysproc.c | 68 ioapic.c |
| 09 elf.h | | 69 picirq.c |
| | # file system | 70 kbd.h |
| # entering xv6 | 35 buf.h | 71 kbd.c |
| 10 entry.S | 35 fcntl.h | 72 console.c |
| 11 entryother.S | 36 stat.h | 75 timer.c |
| 12 main.c | 36 fs.h | 76 uart.c |
| | 37 file.h | |
| # locks | 38 ide.c | # user-level |
| 14 spinlock.h | 40 bio.c | 77 initcode.S |
| 14 spinlock.c | 41 log.c | 77 usys.S |
| | 44 fs.c | 78 init.c |
| # processes | 52 file.c | 78 sh.c |
| 16 vm.c | 54 sysfile.c | |
| 20 proc.h | 59 exec.c | # bootloader |
| 21 proc.c | | 84 bootasm.S |
| 27 switch.S | # pipes | 85 bootmain.c |
| 27 kalloc.c | 60 pipe.c | |

The source listing is preceded by a cross-reference that lists every defined constant, struct, global variable, and function in xv6. Each entry gives, on the same line as the name, the line number (or, in a few cases, numbers) where the name is defined. Successive lines in an entry list the line numbers where the name is used. For example, this entry:

```
switch 2658
      0374 2428 2466 2657 2658
```

indicates that switch is defined on line 2658 and is mentioned on five lines on sheets 03, 24, and 26.

```
acquire 1474
0377 1474 1478 2210 2373 2408
2467 2524 2568 2583 2616 2629
2826 2843 3116 3472 3492 3907
3965 4070 4130 4279 4310 4658
4691 4711 4740 4758 4768 5229
5254 5268 6063 6084 6105 7260
7416 7458 7506
allocproc 2205
2205 2257 2310
allocum 1853
0422 1853 1867 2287 5943 5953
alltraps 3004
2959 2967 2980 2985 3003 3004
ALT 7010
7010 7038 7040
argfd 5419
5419 5456 5471 5483 5494 5506
argint 3295
0395 3295 3308 3324 3432 3456
3470 5424 5471 5483 5708 5776
5777 5826
argptr 3304
0396 3304 5471 5483 5506 5857
argstr 3321
0397 3321 5518 5608 5708 5757
5775 5806 5826
BACK 7861
7861 7974 8120 8389
backcmd 7896 8114
7896 7909 7975 8114 8116 8242
8355 8390
BACKSPACE 7350
7350 7367 7394 7426 7432
ballToc 4454
4454 4476 4817 4825 4829
BBLock 3695
3695 4463 4488
begin_trans 4277
0333 4277 5283 5374 5523 5613
5711 5756 5774
bfree 4481
4481 4864 4874 4877
bget 4066
4066 4096 4106
binit 4038
0261 1231 4038
bmap 4810
4810 4836 4919 4969
bootmain 8517
8517 7358
```

```
clearpteu 1929
0431 1929 1935 5955
cli 0557
0557 0559 1126 1560 7310 7389
8412
cmd 7865
7865 7877 7886 7887 7892 7893
7898 7902 7906 7915 7918 7923
7931 7937 7941 7951 7975 7977
8052 8055 8057 8058 8059 8060
8063 8064 8066 8068 8069 8070
8071 8072 8073 8074 8075 8076
8079 8080 8082 8084 8085 8086
8087 8088 8089 8100 8101 8103
8105 8106 8107 8108 8109 8110
8113 8114 8116 8118 8119 8120
8121 8122 8212 8213 8214 8215
8217 8221 8224 8230 8231 8234
8237 8239 8242 8246 8248 8250
8253 8255 8258 8260 8263 8264
8275 8278 8281 8285 8300 8303
8308 8312 8313 8316 8321 8322
8328 8337 8338 8344 8345 8351
8352 8361 8364 8366 8372 8373
8378 8384 8390 8391 8394
COM1 7613
7613 7623 7626 7627 7628 7629
7630 7631 7634 7640 7641 7657
7659 7667 7669
commit_trans 4301
0334 4301 5285 5379 5528 5546
5555 5645 5652 5713 5758 5762
5779 5783
CONSOLE 3787
3787 7521 7522
consoleinit 7516
0267 1227 7516
consoleintr 7412
0269 7198 7412 7675
consoleread 7451
7451 7522
consolewrite 7501
7501 7521
conspc 7386
7216 7247 7268 7286 7289 7293
7294 7386 7426 7432 7439 7508
context 2093
0251 0374 2056 2093 2111 2238
2239 2240 2241 2478 2516 2678
copyout 2018
```

0287 5011 5017 5059 5174 5621 5264 5266 5302 5315 5352 5413
5667 5419 5422 5438 5453 5467 5479
DIRSIZ 3698 5492 5503 5705 5854 6006 6021
3698 3702 5005 5072 5128 5129 7210 7608 7878 7933 7934 8064
5191 5515 5605 5661 8072 8272
DPL_USER 0779 filealloc 5225
0779 1627 1628 2264 2265 3073 0276 5225 5726 6027
3168 3177 fileclose 5264
E0ESC 7016 0277 2365 5264 5270 5497 5728
7016 7170 7174 7175 7177 7180 5865 5866 6054 6056
elfhdr 0955 filedup 5252
0955 5915 8519 8524 0278 2329 5252 5256 5460
ELF_MAGIC 0952 fileinit 5218
0952 5928 8530 0279 1232 5218
ELF_PROG_LOAD 0986 fileread 5315
0986 5939 0280 5315 5330 5473
entry 1040 filestat 5302
0961 1036 1039 1040 2952 2953 0281 5302 5508
5987 6321 8521 8545 8546 filewrite 5352
E0I 6614 0282 5352 5384 5389 5485
6614 6684 6725 FL_IF 0710
ERROR 6635 0710 1562 1568 2268 2513 6708
6635 6677 fork 2304
6617 6680 6681 0360 2304 3411 7760 7823 7825
8043 8045
EXEC 7857 fork1 8039
7857 7922 8059 8365 7900 7942 7954 7961 7976 8024
8039
exec 5910 forkret 2533
0273 5842 5910 7768 7829 7830 2167 2241 2533
7926 7927 freerange 2801
execcmd 7869 8053 2761 2784 2790 2801
7869 7910 7923 8053 8055 8321 freemv 1910
8327 8328 8356 8366 0424 1910 1915 1977 2421 5990
exit 2354 5995
0359 2354 2390 3105 3109 3169 gatedesc 0901
3178 3417 7716 7719 7761 7826 0523 0526 0901 3061
7831 7916 7925 7935 7980 8028 getcallerpcs 1526
8035 0378 1488 1526 2678 7315
EXTWEM 0202 0202 0208 1729 getcmd 7984
0202 0208 1729 7984 8015
fdalloc 5438 5438 5458 5726 5862 gettoken 8156
5438 5458 5726 5862 8156 8241 8245 8257 8270 8271
fetchint 3267 8307 8311 8333
0398 3267 3297 5833 growproc 2281
0399 3279 3326 5839 0361 2281 3459
0399 3279 3326 5839 havedisk1 3828
file 3750 3828 3864 3962
0252 0276 0277 0278 0280 0281 holding 1544
0282 0351 2114 3750 4420 5208
5214 5224 5227 5230 5251 5252 0379 1477 1504 1544 2507

ialloc 4603 0453 3837 3863 6554 7164 7167
0288 4603 4624 5676 5677 7361 7363 7634 7640 7641 7657
IBLOCK 3689 7667 7669 8423 8431 8554
3689 4613 4634 4718 initlock 1462
ICRHI 6628 0380 1462 2175 2782 3075 3855
6628 6687 6756 6768 0402 4211 4570 5220 6035 7518
ICRLO 6618 7519
6618 6688 6689 6757 6759 6769 initlog 4205
ID 6611 0331 2544 4205 4208
inituvm 1803
ideinit 3851 0425 1803 1808 2261
inode 3762
ideintr 3902 0304 1234 3851
0305 3124 3902 0253 0286 0287 0288 0289 0291
ideLock 3825 0292 0293 0294 0295 0297 0298
3825 3855 3907 3909 3928 3965 0299 0300 0301 0426 1818 2115
3979 3982 3756 3762 3781 3782 4423 4564
iderw 3954 4573 4602 4629 4653 4656 4662
0306 3954 3959 3961 3963 4108 4688 4689 4703 4735 4756 4778
4119 4810 4856 4887 4902 4952 5010
idestart 3875 5195 5516 5561 5603 5656 5660
3829 3875 3878 3926 3975 5706 5754 5769 5804 5916 7451
7501
idewait 3833 INPUT_BUF 7400
3833 3858 3880 3916 7400 7403 7424 7436 7438 7440
IDE_BSY 3813 7468
3813 3837 insl 0462
IDE_CMD_READ 3818 0462 0464 3917 8573
3818 3891 install_trans 4221
IDE_CMD_WRITE 3819 4221 4271 4305
3819 3888 INT_DISABLED 6819
IDE_DF 3815 6819 6867
3815 3839 IOAPIC 6808
IDE_DRDY 3814 6808 6858
3814 3837 ioapic 6827
IDE_ERR 3816 6507 6529 6530 6824 6827 6836
3816 3839 6837 6843 6844 6858
idtinit 3079 ioapicenable 6873
0406 1265 3079 0309 3857 6873 7526 7643
idup 4689 ioapicid 6417
0289 2330 4689 5161 0310 6417 6530 6547 6861 6862
iget 4654 4573 4620 4654 4674 5029 5159 ioapicinit 6851
4573 4620 4654 4674 5029 5159 0311 1226 6851 6862
iinit 4568 0290 1233 4568
ioapicread 6834
0290 1233 4568 6834 6859 6860
iLock 4703 0291 4703 4709 4729 5164 5305 ioapicwrite 6841
5324 5375 5525 5538 5551 5615 6841 6867 6868 6881 6882
5623 5665 5669 5679 5719 5808 IO_PIC1 6907
5922 7463 7483 7510 6907 6920 6935 6944 6947 6952
inb 0453 6962 6976 6977

IO_PIC2 6908
6908 6921 6936 6965 6966 6967
6970 6979 6980
IO_RTC 6735
6735 6748 6749
IO_TIMER1 7559
7559 7568 7578 7579
IPB 3686
3686 3689 3695 4614 4635 4719
iput 4756
0292 2370 4756 4762 4781 5060
5182 5284 5544 5814
IRQ_COM1 2933
2933 3134 7642 7643
IRQ_ERROR 2935
2935 6677
IRQ_IDE 2934
2934 3123 3127 3856 3857
IRQ_KBD 2932
2932 3130 7525 7526
IRQ_SLAVE 6910
6910 6914 6952 6967
IRQ_SPIRIOUS 2936
2936 3139 6657
IRQ_TIMER 2931
2931 3114 3173 6664 7580
isdtrempty 5561
5561 5568 5627
ismp 6415
0337 1235 6415 6512 6520 6540
6543 6855 6875
itrunc 4856
4423 4765 4856
iunlock 4735
0293 4735 4738 4780 5171 5307
5327 5378 5534 5732 5813 7456
7505
iunlockput 4778
0294 4778 5166 5175 5178 5527
5540 5543 5554 5628 5639 5643
5651 5668 5672 5696 5721 5729
5761 5782 5810 5948 5997
iupdate 4629
0295 4629 4767 4882 4978 5533
5553 5637 5642 5683 5687
I_BUSY 3775
3775 4712 4714 4737 4741 4761
4763
I_VALID 3776
3776 4717 4727 4759

kvmalloc 1757
0418 1220 1757
lapiceoi 6722
0325 3121 3125 3132 3136 3142
6722
lapicinit 6651
0326 1222 1256 6651
lapicstartap 6740
0327 1299 6740
lapicw 6644
6644 6657 6663 6664 6665 6668
6669 6674 6677 6680 6681 6684
6687 6688 6693 6725 6756 6757
6759 6768 6769
lcr3 0590
0590 1768 1783
lgdt 0512
0512 0520 1133 1633 8441
lidt 0526
0526 0534 3081
LINT0 6633
6633 6668
LINT1 6634
6634 6669
LIST 7860
7860 7940 8107 8383
listcmd 7890 8101
7890 7911 7941 8101 8103 8246
8357 8384
loadgs 0551
0551 1634
loaduvm 1818
0426 1818 1824 1827 5945
log 4190 4200
4190 4200 4211 4213 4214 4215
4225 4226 4227 4239 4242 4243
4244 4256 4259 4260 4261 4272
4279 4280 4281 4283 4284 4303
4306 4310 4311 4312 4313 4329
4331 4334 4335 4338 4339 4343
4344
logheader 4185
4185 4196 4207 4208 4240 4257
LOGSIZE 0160
0160 4187 4329 5367
log_write 4325
0332 4325 4444 4468 4494 4618
4642 4830 4972
ltr 0538
0538 0540 1780

mappages 1679
1679 1748 1811 1872 1971
MAXARC 0159
0159 5822 5914 5960
MAXARCS 7863
7863 7871 7872 8340
MAXFILE 3673
3673 4965
memcmp 6165
0386 6165 6445 6488
memmove 6181
0387 1285 1812 1970 2032 4228
4340 4432 4641 4725 4921 4971
5129 5131 6181 6204 7373
memset 6154
0388 1666 1744 1810 1871 2240
2263 2823 4443 4616 5632 5829
6154 7375 7987 8058 8069 8085
8106 8119
microdelay 6731
0328 6731 6758 6760 6770 7658
min 4422
4422 4920 4970
mp 6302
6302 6408 6437 6444 6445 6446
6455 6460 6464 6465 6468 6469
6480 6483 6485 6487 6494 6504
6510 6550
mpbcpu 6420
0338 6420
MPBUS 6352
6352 6533
mpconf 6313
6313 6479 6482 6487 6505
mpconfig 6480
6480 6510
mpenter 1252
1252 1296
mpinit 6501
0339 1221 6501 6519 6539
MPIOAPIC 6353
6353 6528
mpioapic 6339
6339 6507 6529 6531
MPIOINTR 6354
6354 6534
MPLINTR 6355
6355 6535
mpmain 1262
1209 1241 1257 1262

MPPROC 6351
6351 6516
mpproc 6328
6328 6506 6517 6526
mpsearch 6456
6456 6485
mpsearch1 6438
6438 6464 6468 6471
multiboot_header 1025
1024 1025
namecmp 5003
0296 5003 5024 5618
namei 5189
0297 2273 5189 5520 5717 5806
5920
nameiparent 5196
0298 5154 5169 5181 5196 5536
5610 5663
namex 5154
5154 5192 5198
NBUF 0155
0155 4030 4053
NCPU 0152
0152 2068 6413
ncpu 6416
1224 1287 2069 3857 6416 6518
6519 6523 6524 6525 6545
NDEV 0157
0157 4908 4958 5211
NDIRECT 3671
3671 3673 3682 3773 4815 4820
4824 4825 4862 4869 4870 4877
4878
NELEM 0434
0434 1747 2672 3380 5831
nextpid 2166
2166 2219
NFILE 0154
0154 5214 5230
NINDIRECT 3672
3672 3673 4822 4872
NINODE 0156
0156 4564 4662
NO 7006
7006 7052 7055 7057 7058 7059
7060 7062 7074 7077 7079 7080
7081 7082 7084 7102 7103 7105
7106 7107 7108
NOFILE 0153
0153 2114 2327 2363 5426 5442

NPDENTRIES 0821
0821 1311 1917
NPROC 0150
0150 2161 2211 2379 2412 2468
2607 2630 2669
NPTENTRIES 0822
0822 1894
NSEGS 2051
1611 2051 2058
nulterminate 8352
8215 8230 8352 8373 8379 8380
8385 8386 8391
NUMLOCK 7013
7013 7046
outb 0471
0471 3861 3870 3881 3882 3883
3884 3885 3886 3888 3891 6553
6554 6748 6749 6920 6921 6935
6936 6944 6947 6952 6962 6965
6966 6967 6970 6976 6977 6979
6980 7360 7362 7378 7379 7380
7381 7577 7578 7579 7623 7626
7627 7628 7629 7630 7631 7659
8428 8436 8564 8565 8566 8567
8568 8569
outs1 0483
0483 0485 3889
outw 0477
0477 1181 1183 8474 8476
O_CREATE 3553
3553 5710 8278 8281
O_RDONLY 3550
3550 5720 8275
O_RDWR 3552
3552 5738 7814 7816 8007
O_WRONLY 3551
3551 5737 5738 8278 8281
P2V 0218
0218 1219 1238 6462 6750 7352
panic 7305 8032
0270 1478 1505 1569 1571 1690
1746 1782 1808 1824 1827 1898
1915 1935 1964 1966 2260 2360
2390 2508 2510 2512 2514 2556
2559 2820 3155 3878 3959 3961
3963 4096 4117 4128 4208 4330
4332 4476 4492 4624 4674 4709
4729 4738 4762 4836 5017 5021
5067 5075 5256 5270 5330 5384
5389 5568 5626 5634 5677 5690

5694 7263 7305 7312 7901 7920
7953 8032 8045 8228 8272 8306
8310 8336 8341
panicked 7218
7218 7318 7388
parseblock 8301
8301 8306 8325
parsecmd 8218
7902 8025 8218
parseexec 8317
8214 8255 8317
parseline 8235
8212 8224 8235 8246 8308
parsepipe 8251
8213 8239 8251 8258
parseredirs 8264
8264 8312 8331 8342
PCINT 6632
6632 6674
pde_t 0103
0103 0420 0421 0422 0423 0424
0425 0426 0427 0430 0431 1210
1270 1311 1610 1654 1656 1679
1736 1739 1742 1803 1818 1853
1882 1910 1929 1952 1953 1955
2002 2018 2105 5918
PDX 0812
0812 1659
PDXSHIFT 0827
0812 0818 0827 1315
peek 8201
8201 8225 8240 8244 8256 8269
8305 8309 8324 8332
PGROUNDOWN 0830
0830 1684 1685 2025
PGROUNDUP 0829
0829 1863 1890 2804 5952
PGSIZE 0823
0823 0829 0830 1310 1666 1694
1695 1744 1807 1810 1811 1823
1825 1829 1832 1864 1871 1872
1891 1894 1962 1970 1971 2029
2035 2262 2269 2805 2819 2823
5953 5955
PHYSTOP 0203
0203 1238 1731 1745 1746 2819
picenable 6925
0343 3856 6925 7525 7580 7642
picinit 6932
0344 1225 6932

picsetmask 6917
6917 6927 6983
pinit 2173
0363 1229 2173
PIPE 7859
7859 7950 8086 8377
pipe 6011
0254 0352 0353 0354 3755 5281
5322 5359 6011 6023 6029 6035
6039 6043 6061 6080 6101 7763
7952 7953
pipealloc 6021
0351 5859 6021
pipeclose 6061
0352 5281 6061
pipecmd 7884 8080
7884 7912 7951 8080 8082 8258
8358 8378
piperead 6101
0353 5322 6101
PIPESIZE 6009
6009 6013 6086 6094 6116
pipewrite 6080
0354 5359 6080
popch 1566
0383 1521 1566 1569 1571 1784
printint 7226
7226 7276 7280
proc 2103
0255 0358 0428 1205 1458 1606
1638 1773 1779 2065 2080 2103
2109 2156 2161 2164 2204 2207
2211 2254 2285 2287 2290 2293
2294 2307 2314 2320 2321 2322
2328 2329 2330 2334 2356 2359
2364 2365 2366 2370 2371 2376
2379 2380 2388 2405 2412 2413
2433 2439 2460 2468 2475 2478
2483 2511 2516 2525 2555 2573
2574 2578 2605 2607 2627 2630
2665 2669 3055 3104 3106 3108
3151 3159 3160 3162 3168 3173
3177 3255 3269 3283 3286 3297
3310 3379 3381 3384 3385 3406
3440 3458 3475 3807 4416 5161
5411 5426 5443 5444 5496 5814
5815 5864 5904 5981 5984 5985
5986 5987 5988 5989 6004 6087
6107 6411 6506 6517 6518 6519
6522 7213 7461 7610


```
procdump 2654          6810 6860
  0364 2654 7420      REG_TABLE 6812
proghdr 0974          6812 6867 6868 6881 6882
  0974 5917 8520 8534  REG_VER 6811
  0844 1661 1828 1896 1919 1967  release 1502
  2011          0381 1502 1505 2214 2220 2427
PTE_P 0833          2434 2485 2527 2537 2569 2582
  0833 1313 1315 1660 1670 1689  2618 2636 2640 2831 2848 3119
  1691 1895 1918 1965 2007  3476 3481 3494 3909 3928 3982
PTE_PS 0840          4078 4092 4142 4284 4313 4665
  0840 1313 1315  4681 4693 4715 4743 4764 4773
pte_t 0847          5233 5237 5258 5272 5278 6072
  0847 1653 1657 1661 1663 1682  6075 6088 6097 6108 6119 7301
  1821 1884 1931 1956 2004  7448 7462 7482 7509
PTE_U 0835          ROOTDEV 0158
  0835 1670 1811 1872 1936 1971  0158 4212 4215 5159
  2009          ROOTINO 3660
PTE_W 0834          3660 5159
  0834 1313 1315 1670 1729 1731  run 2764
  1732 1811 1872 1971  2661 2764 2765 2771 2817 2827
PTX 0815          2840
  0815 1672          runcmd 7906
PTXSHIFT 0826          7906 7920 7937 7943 7945 7959
  0815 0818 0826  7966 7977 8025
pushcli 1555          RUNNING 2100
  0382 1476 1555 1775  2100 2477 2511 2661 3173
rcr2 0582          safestrncpy 6232
  0582 3154 3161  0389 2272 2334 5981 6232
readeflags 0544          sched 2503
  0544 1559 1568 2513 6708  0366 2389 2503 2508 2510 2512
  2514 2526 2575
readi 4902          2514 2526 2575
  0299 1833 4902 5020 5066 5325  scheduler 2458
  5567 5568 5926 5937  0365 1267 2056 2458 2478 2516
readsb 4427          SCROLLLOCK 7014
  0285 4212 4427 4461 4487 4610  7014 7047
readsect 8560          SECTSIZE 8512
  8560 8595          8512 8573 8586 8589 8594
readseg 8579          SEG 0769
  8514 8527 8538 8579  0769 1625 1626 1627 1628 1631
read_head 4237          SEG16 0773
  4237 4270          0773 1776
recover_from_log 4268  segdesc 0752
  4202 4216 4268  0509 0512 0752 0769 0773 1611
  2058
REDIR 7858          seginit 1616
  7858 7930 8070 8371  0417 1223 1255 1616
redircmd 7875 8064          SEG_ASM 0660
  7875 7913 7931 8064 8066 8275  0660 1190 1191 8484 8485
  8278 8281 8359 8372  SEG_KCODE 0741
REG_ID 6810
```

```
0741 1150 1625 3072 3073 8453  sti 0563
SEG_KCPU 0743          0563 0565 1573 2464
0743 1631 1634 3016  stosb 0492
SEG_KDATA 0742          0492 0494 6160 8540
  0742 1154 1626 1778 3013 8458  stosl 0501
  0501 0503 6158
SEG_NULLASM 0654          0501 6251
  0654 1189 8483  strlen 6251
SEG_TSS 0746          0390 5962 5963 6251 8019 8223
  0746 1776 1777 1780  strncmp 6208
SEG_UCODE 0744          0391 5005 6208
  0744 1627 2264  strncpy 6218
SEG_UDATA 0745          0392 5072 6218
  0745 1628 2265  STS_IC32 0800
SETGATE 0921          STS_IC32 0800
  0921 3072 3073  STS_T32A 0927
  0921 3072 3073  STS_T32A 0927
setupkvm 1737          0797 1776
  0420 1737 1759 1960 2259 5931  STS_TG32 0801
SHIFT 7008          0801 0927
  7008 7036 7037 7185  sum 6426
skipelem 5115          6426 6428 6430 6432 6433 6445
  5115 5163  6492
sleep 2553          superblock 3664
  0367 2439 2553 2556 2559 2659  0258 0285 3664 4210 4427 4458
  3479 3979 4081 4281 4713 6092  4484 4608
  6111 7466 7779  SVR 6615
  6615 6657
spinlock 1401          6615 6657
  0256 0367 0377 0379 0380 0381  switchkvm 1766
  0409 1401 1459 1462 1474 1502  0429 1254 1760 1766 2479
  1544 2157 2160 2553 2759 2769  switchvum 1773
  3058 3063 3810 3825 4025 4029  0428 1773 1782 2294 2476 5989
  4153 4191 4417 4563 5209 5213  swtch 2708
  6007 6012 7208 7221 7402 7606  0374 2478 2516 2707 2708
  start 1125 7708 8411  SYSCALL 7753 7760 7761 7762 7763 77
  1124 1125 1167 1175 1177 4192  7760 7761 7762 7763 7764 7765
  4213 4226 4239 4256 4339 7707  7766 7767 7768 7769 7770 7771
  7708 8410 8411 8467  7772 7773 7774 7775 7776 7777
startothers 1274          7778 7779 7780
  1208 1237 1274  syscall 3375
stat 3604          0400 3107 3257 3375
  0257 0281 0300 3604 4414 4887  SYS_chdir 3209
  5302 5409 5504 7803  3209 3359
  4887          sys_chdir 5801
  0300 4887 5306  3329 3359 5801
STA_R 0669 0786          SYS_close 3221
  0669 0786 1190 1625 1627 8484  3221 3371
STA_W 0668 0785          sys_close 5489
  0668 0785 1191 1626 1628 1631  3330 3371 5489
  8485  SYS_dup 3210
STA_X 0665 0782          3210 3360
  0665 0782 1190 1625 1627 8484  sys_dup 5451
```

3331 3360 5451
SYS_exec 3207 3212 3362
sys_sbrk 3451 3344 3362 3451
SYS_sleep 3213 3213 3363
sys_sleep 3465 3345 3363 3465
SYS_unlink 3218 3218 3368
sys_unlink 5601 3346 3368 5601
SYS_uptime 3214 3214 3364
sys_uptime 3488 3349 3364 3488
SYS_wait 3203 3203 3353
sys_wait 3422 3347 3353 3422
SYS_write 3216 3216 3366
sys_write 5477 3348 3366 5477
taskstate 0851 0851 2057
TDCR 6639 6639 6663
ticks 3064 0407 3064 3117 3118 3473 3474
3479 3493
tickslock 3063 0409 3063 3075 3116 3119 3472
3476 3479 3481 3492 3494
TICR 6637 6637 6665
TIMER 6629 6629 6664
timerinit 7574 0403 1236 7574
TIMER_16BIT 7571 7571 7577
TIMER_DIV 7566 7566 7578 7579
TIMER_FREQ 7565 7565 7566
TIMER_MODE 7568 7568 7577
TIMER_RATEGEN 7570 7570 7577
TIMER_SELO 7569

7569 7577
TPR 6613 3036 1239 2252 2260
6613 6693 uva2ka 2002
trap 3101 0421 2002 2026
2952 2954 3022 3101 3153 3155 V2P 0217
3158 0217 1730 1731
trapframe 0602 V2P_WO 0220
0602 2110 2231 3101 0220 1036 1046
trapret 3027 VER 6612
2168 2236 3026 3027 6612 6673
tvinit 3067 wait 2403
0408 1230 3067 0369 2403 3424 7762 7833 7944
T_DEV 3602 7970 7971 8026
3602 4907 4957 5778 waitdisk 8551
T_DIR 3600 8551 8563 8572
3600 5016 5165 5526 5627 5635 wakeup 2614
5685 5720 5757 5809 0370 2614 3118 3922 4140 4312
T_FILE 3601 4742 4770 6066 6069 6091 6096
3601 5670 5712 6118 7442
T_IRQ0 2929 wakeup1 2603
2929 3114 3123 3127 3130 3134 2170 2376 2383 2603 2617
3138 3139 3173 6657 6664 6677 walkpgdir 1654
6867 6881 6947 6966 1654 1687 1826 1892 1933 1963
T_SYSCALL 2926 2006
2926 3073 3103 7713 7718 7757 writei 4952
uart 7615 0301 4952 5074 5376 5633 5634
7615 7636 7655 7665 write_head 4254
uartgetc 7663 4254 4273 4304 4307
7663 7675 xchg 0569
uartinit 7618 0569 1266 1483 1519
0412 1228 7618 yield 2522
uartintr 7673 0371 2522 3174
0413 3135 7673 __attribute__ 1310
uartputc 7651 0270 0365 1209 1310
0414 7395 7397 7647 7651

```
0100 typedef unsigned int  uint;
0101 typedef unsigned short ushort;
0102 typedef unsigned char  uchar;
0103 typedef uint pde_t;
0104
0105
0106
0107
0108
0109
0110
0111
0112
0113
0114
0115
0116
0117
0118
0119
0120
0121
0122
0123
0124
0125
0126
0127
0128
0129
0130
0131
0132
0133
0134
0135
0136
0137
0138
0139
0140
0141
0142
0143
0144
0145
0146
0147
0148
0149
```

```
0150 #define NPROC      64 // maximum number of processes
0151 #define KSTACKSIZE 4096 // size of per-process kernel stack
0152 #define NCPU        8 // maximum number of CPUs
0153 #define NFILE       16 // open files per process
0154 #define NBUF        100 // open files per system
0155 #define NINODE      50 // size of disk block cache
0156 #define NDEV        10 // maximum number of active i-nodes
0157 #define ROOTDEV     1 // device number of file system root disk
0158 #define MAXARG      32 // max exec arguments
0159 #define LOGSIZE     10 // max data sectors in on-disk log
0160
0161
0162
0163
0164
0165
0166
0167
0168
0169
0170
0171
0172
0173
0174
0175
0176
0177
0178
0179
0180
0181
0182
0183
0184
0185
0186
0187
0188
0189
0190
0191
0192
0193
0194
0195
0196
0197
0198
0199
```

```

0200 // Memory layout
0201
0202 #define EXTMEM 0x100000 // Start of extended memory
0203 #define PHYSTOP 0xE000000 // Top physical memory
0204 #define DEVSPACE 0xFE000000 // Other devices are at high addresses
0205
0206 // Key addresses for address space layout (see knap in vm.c for layout)
0207 #define KERNBASE 0x80000000 // First kernel virtual address
0208 #define KERNLINK (KERNBASE+EXTMEM) // Address where kernel is linked
0209
0210 #ifndef __ASSEMBLER__
0211
0212 static inline uint v2p(void *a) { return ((uint) (a)) - KERNBASE; }
0213 static inline void *p2v(uint a) { return (void *) ((a) + KERNBASE); }
0214
0215 #endif
0216
0217 #define V2P(a) (((uint) (a)) - KERNBASE)
0218 #define P2V(a) (((void *) (a)) + KERNBASE)
0219
0220 #define V2P_WO(x) ((x) - KERNBASE) // same as V2P, but without casts
0221 #define P2V_WO(x) ((x) + KERNBASE) // same as P2V, but without casts
0222
0223
0224
0225
0226
0227
0228
0229
0230
0231
0232
0233
0234
0235
0236
0237
0238
0239
0240
0241
0242
0243
0244
0245
0246
0247
0248
0249

```

```

0250 struct buf;
0251 struct context;
0252 struct file;
0253 struct inode;
0254 struct pipe;
0255 struct proc;
0256 struct spinlock;
0257 struct stat;
0258 struct superblock;
0259
0260 // bio.c
0261 void binit(void);
0262 struct buf* bread(uint, uint);
0263 void brelse(struct buf*);
0264 void bwrite(struct buf*);
0265
0266 // console.c
0267 void consoleinit(void);
0268 void cprintf(char*, ...);
0269 void consoleintr(int (*)(void));
0270 void panic(char*) __attribute__((noreturn));
0271
0272 // exec.c
0273 int exec(char*, char**);
0274
0275 // file.c
0276 struct file* filealloc(void);
0277 void fileclose(struct file*);
0278 struct file* filedup(struct file*);
0279 void fileinit(void);
0280 int fileread(struct file*, char*, int n);
0281 int filestat(struct file*, struct stat*);
0282 int filewrite(struct file*, char*, int n);
0283
0284 // fs.c
0285 void readsb(int dev, struct superblock *sb);
0286 int dirlink(struct inode*, char*, uint);
0287 struct inode* dirlookup(struct inode*, char*, uint*);
0288 struct inode* ialloc(uint, short);
0289 struct inode* idup(struct inode*);
0290 void iinit(void);
0291 void ilock(struct inode*);
0292 void iput(struct inode*);
0293 void iunlock(struct inode*);
0294 void iunlockput(struct inode*);
0295 void iupdate(struct inode*);
0296 int namecmp(const char*, const char*);
0297 struct inode* namei(char*);
0298 struct inode* nameiparent(char*, char*);
0299 int readi(struct inode*, char*, uint, uint);

```

```

0300 void      stati(struct inode*, struct stat*);
0301 int        writei(struct inode*, char*, uint, uint);
0302
0303 // ide.c
0304 void      ideinit(void);
0305 void      ideintr(void);
0306 void      iderw(struct buf*);
0307
0308 // ioapic.c
0309 void      ioapicenable(int irq, int cpu);
0310 extern uchar ioapicid;
0311 void      ioapicinit(void);
0312
0313 // kalloc.c
0314 char*      kalloc(void);
0315 void      kfree(char*);
0316 void      kinit1(void*, void*);
0317 void      kinit2(void*, void*);
0318
0319 // kbd.c
0320 void      kbdtintr(void);
0321
0322 // lapic.c
0323 int        cpunum(void);
0324 extern volatile uint* lapic;
0325 void      lapiceoi(void);
0326 void      lapicinit(void);
0327 void      lapicstartap(uchar, uint);
0328 void      microdelay(int);
0329
0330 // log.c
0331 void      initlog(void);
0332 void      log_write(struct buf*);
0333 void      begin_trans();
0334 void      commit_trans();
0335
0336 // mp.c
0337 extern int ismp;
0338 int        mpbcpu(void);
0339 void      mpinit(void);
0340 void      mpstartthem(void);
0341
0342 // picirq.c
0343 void      plicenable(int);
0344 void      picinit(void);
0345
0346
0347
0348
0349

```

```

0350 // pipe.c
0351 int        pipeallloc(struct file**, struct file**);
0352 void      pipeclose(struct pipe*, int);
0353 int        piperead(struct pipe*, char*, int);
0354 int        pipewrite(struct pipe*, char*, int);
0355
0356
0357 // proc.c
0358 struct proc* copyproc(struct proc*);
0359 void      exit(void);
0360 int        fork(void);
0361 int        growproc(int);
0362 int        kill(int);
0363 void      pinit(void);
0364 void      procdump(void);
0365 void      scheduler(void) __attribute__((noreturn));
0366 void      sched(void);
0367 void      sleep(void*, struct spinlock*);
0368 void      userinit(void);
0369 int        wait(void);
0370 void      wakeup(void*);
0371 void      yield(void);
0372
0373 // swtch.S
0374 void      swtch(struct context**, struct context*);
0375
0376 // spinlock.c
0377 void      acquire(struct spinlock*);
0378 void      getcallerpcs(void*, uint*);
0379 int        holding(struct spinlock*);
0380 void      initlock(struct spinlock*, char*);
0381 void      release(struct spinlock*);
0382 void      pushcli(void);
0383 void      popcli(void);
0384
0385 // string.c
0386 int        memcmp(const void*, const void*, uint);
0387 void*      memmove(void*, const void*, uint);
0388 void*      memset(void*, int, uint);
0389 char*      safestrcpy(char*, const char*, int);
0390 int        strlen(const char*);
0391 int        strncmp(const char*, const char*, uint);
0392 char*      strncpy(char*, const char*, int);
0393
0394 // syscall.c
0395 int        argint(int, int*);
0396 int        argptr(int, char**, int);
0397 int        argstr(int, char**);
0398 int        fetchint(uint, int*);
0399 int        fetchstr(uint, char**);

```

```

0400 void      syscall(void);
0401
0402 // timer.c
0403 void      timerinit(void);
0404
0405 // trap.c
0406 void      idtinit(void);
0407 extern uint ticks;
0408 void      tvinit(void);
0409 extern struct spinlock tickslock;
0410
0411 // uart.c
0412 void      uartinit(void);
0413 void      uartintr(void);
0414 void      uartputc(int);
0415
0416 // vm.c
0417 void      seginit(void);
0418 void      kvmalloc(void);
0419 void      vmenable(void);
0420 pde_t*    setupkvm(void);
0421 char*     uva2ka(pde_t*, char*);
0422 int      allocvm(pde_t*, uint, uint);
0423 int      deallocvm(pde_t*, uint, uint);
0424 void      freevm(pde_t*);
0425 void      initvm(pde_t*, char*, uint);
0426 int      loadvm(pde_t*, char*, struct inode*, uint, uint);
0427 pde_t*    copyvm(pde_t*, uint);
0428 void      switchvm(struct proc*);
0429 void      switchvm(void);
0430 int      copyout(pde_t*, uint, void*, uint);
0431 void      clearpteu(pde_t *pgdir, char *uva);
0432
0433 // number of elements in fixed-size array
0434 #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
0435
0436
0437
0438
0439
0440
0441
0442
0443
0444
0445
0446
0447
0448
0449

```

```

0450 // Routines to let C code use special x86 instructions.
0451
0452 static inline uchar
0453 inb(ushort port)
0454 {
0455     uchar data;
0456
0457     asm volatile("in %1,%0" : "=a" (data) : "d" (port));
0458     return data;
0459 }
0460
0461 static inline void
0462 insl(int port, void *addr, int cnt)
0463 {
0464     asm volatile("cld; rep insl" :
0465         "=D" (addr), "=c" (cnt) :
0466         "d" (port), "0" (addr), "1" (cnt) :
0467         "memory", "cc");
0468 }
0469
0470 static inline void
0471 outb(ushort port, uchar data)
0472 {
0473     asm volatile("out %0,%1" : : "a" (data), "d" (port));
0474 }
0475
0476 static inline void
0477 outw(ushort port, ushort data)
0478 {
0479     asm volatile("out %0,%1" : : "a" (data), "d" (port));
0480 }
0481
0482 static inline void
0483 outsl(int port, const void *addr, int cnt)
0484 {
0485     asm volatile("cld; rep outsl" :
0486         "=S" (addr), "=c" (cnt) :
0487         "d" (port), "0" (addr), "1" (cnt) :
0488         "cc");
0489 }
0490
0491 static inline void
0492 stosb(void *addr, int data, int cnt)
0493 {
0494     asm volatile("cld; rep stosb" :
0495         "=D" (addr), "=c" (cnt) :
0496         "0" (addr), "1" (cnt), "a" (data) :
0497         "memory", "cc");
0498 }
0499

```



```

0500 static inline void
0501 stosl(void *addr, int data, int cnt)
0502 {
0503     asm volatile("cld; rep stosl" :
0504                 "=D" (addr), "=c" (cnt) :
0505                 "0" (addr), "1" (cnt), "a" (data) :
0506                 "memory", "cc");
0507 }
0508
0509 struct segdesc;
0510
0511 static inline void
0512 lgdt(struct segdesc *p, int size)
0513 {
0514     volatile ushort pd[3];
0515
0516     pd[0] = size-1;
0517     pd[1] = (uint)p;
0518     pd[2] = (uint)p >> 16;
0519
0520     asm volatile("lgdt (%0)" : : "r" (pd));
0521 }
0522
0523 struct gatedesc;
0524
0525 static inline void
0526 lidt(struct gatedesc *p, int size)
0527 {
0528     volatile ushort pd[3];
0529
0530     pd[0] = size-1;
0531     pd[1] = (uint)p;
0532     pd[2] = (uint)p >> 16;
0533
0534     asm volatile("lidt (%0)" : : "r" (pd));
0535 }
0536
0537 static inline void
0538 ltr(ushort sel)
0539 {
0540     asm volatile("ltr %0" : : "r" (sel));
0541 }
0542
0543 static inline uint
0544 readeflags(void)
0545 {
0546     uint eflags;
0547     asm volatile("pushfl; popl %0" : "=r" (eflags));
0548     return eflags;
0549 }

```

```

0550 static inline void
0551 loadgs(ushort v)
0552 {
0553     asm volatile("movw %0, %%gs" : "=r" (v));
0554 }
0555
0556 static inline void
0557 cli(void)
0558 {
0559     asm volatile("cli");
0560 }
0561
0562 static inline void
0563 sti(void)
0564 {
0565     asm volatile("sti");
0566 }
0567
0568 static inline uint
0569 xchg(volatile uint *addr, uint newval)
0570 {
0571     uint result;
0572
0573     // The + in "+m" denotes a read-modify-write operand.
0574     asm volatile("lock; xchgl %0, %1" :
0575                 "+m" (*addr), "=a" (result) :
0576                 "1" (newval) :
0577                 "cc");
0578     return result;
0579 }
0580
0581 static inline uint
0582 rcr2(void)
0583 {
0584     uint val;
0585     asm volatile("movl %%cr2,%0" : "=r" (val));
0586     return val;
0587 }
0588
0589 static inline void
0590 lcr3(uint val)
0591 {
0592     asm volatile("movl %0,%%cr3" : "=r" (val));
0593 }
0594
0595
0596
0597
0598
0599

```

```

0600 // Layout of the trap frame built on the stack by the
0601 // hardware and by trapasm.S, and passed to trap().
0602 struct trapframe {

```

```

0603     // registers as pushed by pusha

```

```

0604     uint edi;

```

```

0605     uint esi;

```

```

0606     uint ebp;

```

```

0607     uint oesp;           // useless & ignored

```

```

0608     uint ebx;

```

```

0609     uint edx;

```

```

0610     uint ecx;

```

```

0611     uint eax;

```

```

0612

```

```

0613     // rest of trap frame

```

```

0614     ushort gs;

```

```

0615     ushort padding1;

```

```

0616     ushort fs;

```

```

0617     ushort padding2;

```

```

0618     ushort es;

```

```

0619     ushort padding3;

```

```

0620     ushort ds;

```

```

0621     ushort padding4;

```

```

0622     uint trapno;

```

```

0623

```

```

0624     // below here defined by x86 hardware

```

```

0625     uint err;

```

```

0626     uint eip;

```

```

0627     ushort cs;

```

```

0628     ushort padding5;

```

```

0629     uint eflags;

```

```

0630

```

```

0631     // below here only when crossing rings, such as from user to kernel

```

```

0632     uint esp;

```

```

0633     ushort ss;

```

```

0634     ushort padding6;

```

```

0635 };

```

```

0636

```

```

0637

```

```

0638

```

```

0639

```

```

0640

```

```

0641

```

```

0642

```

```

0643

```

```

0644

```

```

0645

```

```

0646

```

```

0647

```

```

0648

```

```

0649

```

```

0650 //
0651 // assembler macros to create x86 segments
0652 //

```

```

0653

```

```

0654 #define SEG_NULLASM

```

```

0655     .word 0, 0;

```

```

0656     .byte 0, 0, 0, 0

```

```

0657

```

```

0658 // The 0xC0 means the limit is in 4096-byte units

```

```

0659 // and (for executable segments) 32-bit mode.

```

```

0660 #define SEG_ASM(type,base,lim)

```

```

0661     .word (((lim) >> 12) & 0xffff), ((base) & 0xffff);

```

```

0662     .byte (((base) >> 16) & 0xff), (0x90 | (type)),

```

```

0663           (0xC0 | (((lim) >> 28) & 0xf)), (((base) >> 24) & 0xff)

```

```

0664

```

```

0665 #define STA_X      0x8      // Executable segment

```

```

0666 #define STA_E      0x4      // Expand down (non-executable segments)

```

```

0667 #define STA_C      0x4      // Conforming code segment (executable only)

```

```

0668 #define STA_W      0x2      // Writeable (non-executable segments)

```

```

0669 #define STA_R      0x2      // Readable (executable segments)

```

```

0670 #define STA_A      0x1      // Accessed

```

```

0671

```

```

0672

```

```

0673

```

```

0674

```

```

0675

```

```

0676

```

```

0677

```

```

0678

```

```

0679

```

```

0680

```

```

0681

```

```

0682

```

```

0683

```

```

0684

```

```

0685

```

```

0686

```

```

0687

```

```

0688

```

```

0689

```

```

0690

```

```

0691

```

```

0692

```

```

0693

```

```

0694

```

```

0695

```

```

0696

```

```

0697

```

```

0698

```

```

0699

```

```

0700 // This file contains definitions for the
0701 // x86 memory management unit (MMU).
0702
0703 // EFlags register
0704 #define FL_CF 0x00000001 // Carry Flag
0705 #define FL_PF 0x00000004 // Parity Flag
0706 #define FL_AF 0x00000010 // Auxiliary carry Flag
0707 #define FL_ZF 0x00000040 // Zero Flag
0708 #define FL_SF 0x00000080 // Sign Flag
0709 #define FL_TF 0x00000100 // Trap Flag
0710 #define FL_IF 0x00000200 // Interrupt Enable
0711 #define FL_DF 0x00000400 // Direction Flag
0712 #define FL_OF 0x00000800 // Overflow Flag
0713 #define FL_IOPL_MASK 0x00003000 // I/O Privilege Level bitmask
0714 #define FL_IOPL_0 0x00000000 // IOPL == 0
0715 #define FL_IOPL_1 0x00001000 // IOPL == 1
0716 #define FL_IOPL_2 0x00002000 // IOPL == 2
0717 #define FL_IOPL_3 0x00003000 // IOPL == 3
0718 #define FL_NT 0x00004000 // Nested Task
0719 #define FL_RF 0x00010000 // Resume Flag
0720 #define FL_VM 0x00020000 // Virtual 8086 mode
0721 #define FL_AC 0x00040000 // Alignment Check
0722 #define FL_VIF 0x00080000 // Virtual Interrupt Flag
0723 #define FL_VIP 0x00100000 // Virtual Interrupt Pending
0724 #define FL_ID 0x00200000 // ID flag
0725
0726 // Control Register flags
0727 #define CR0_PE 0x00000001 // Protection Enable
0728 #define CR0_MP 0x00000002 // Monitor coProcessor
0729 #define CR0_EM 0x00000004 // Emulation
0730 #define CR0_TS 0x00000008 // Task Switched
0731 #define CR0_ET 0x00000010 // Extension Type
0732 #define CR0_NE 0x00000020 // Numeric Error
0733 #define CR0_WP 0x00010000 // Write Protect
0734 #define CR0_AM 0x00040000 // Alignment Mask
0735 #define CR0_NW 0x20000000 // Not Writethrough
0736 #define CR0_CD 0x40000000 // Cache Disable
0737 #define CR0_PG 0x80000000 // Paging
0738
0739 #define CR4_PSE 0x00000010 // Page size extension
0740
0741 #define SEG_KCODE 1 // kernel code
0742 #define SEG_KDATA 2 // kernel data+stack
0743 #define SEG_KCPU 3 // kernel per-cpu data
0744 #define SEG_UCODE 4 // user code
0745 #define SEG_UDATA 5 // user data+stack
0746 #define SEG_TSS 6 // this process's task state
0747
0748
0749

```

```

0750 #ifndef __ASSEMBLER__
0751 // Segment Descriptor
0752 struct segdesc {
0753     uint lim_15_0 : 16; // Low bits of segment limit
0754     uint base_15_0 : 16; // Low bits of segment base address
0755     uint base_23_16 : 8; // Middle bits of segment base address
0756     uint type : 4; // Segment type (see STS_ constants)
0757     uint s : 1; // 0 = system, 1 = application
0758     uint dpl : 2; // Descriptor Privilege Level
0759     uint p : 1; // Present
0760     uint lim_19_16 : 4; // High bits of segment limit
0761     uint avl : 1; // Unused (available for software use)
0762     uint rsv1 : 1; // Reserved
0763     uint db : 1; // 0 = 16-bit segment, 1 = 32-bit segment
0764     uint g : 1; // Granularity: limit scaled by 4K when set
0765     uint base_31_24 : 8; // High bits of segment base address
0766 };
0767
0768 // Normal segment
0769 #define SEG(type, base, lim, dpl) (struct segdesc) { \
0770     ((lim) >> 12) & 0xffff, (uint)(base) & 0xffff, \
0771     ((uint)(base) >> 16) & 0xfff, type, 1, dpl, 1, \
0772     (uint)(lim) >> 28, 0, 0, 1, 1, (uint)(base) >> 24 }
0773 #define SEG16(type, base, lim, dpl) (struct segdesc) { \
0774     ((lim) & 0xffff, (uint)(base) & 0xffff, \
0775     ((uint)(base) >> 16) & 0xfff, type, 1, dpl, 1, \
0776     (uint)(lim) >> 16, 0, 0, 1, 0, (uint)(base) >> 24 }
0777 #endif
0778
0779 #define DPL_USER 0x3 // User DPL
0780
0781 // Application segment type bits
0782 #define STA_X 0x8 // Executable segment
0783 #define STA_E 0x4 // Expand down (non-executable segments)
0784 #define STA_C 0x4 // Conforming code segment (executable only)
0785 #define STA_W 0x2 // Writable (non-executable segments)
0786 #define STA_R 0x2 // Readable (executable segments)
0787 #define STA_A 0x1 // Accessed
0788
0789 // System segment type bits
0790 #define STS_T16A 0x1 // Available 16-bit TSS
0791 #define STS_LDT 0x2 // Local Descriptor Table
0792 #define STS_T16B 0x3 // Busy 16-bit TSS
0793 #define STS_CG16 0x4 // 16-bit Call Gate
0794 #define STS_TG 0x5 // Task Gate / Coum Transmissions
0795 #define STS_IG16 0x6 // 16-bit Interrupt Gate
0796 #define STS_TG16 0x7 // 16-bit Trap Gate
0797 #define STS_T32A 0x9 // Available 32-bit TSS
0798 #define STS_T32B 0x8 // Busy 32-bit TSS
0799 #define STS_CG32 0xC // 32-bit Call Gate

```

```

0800 #define STS_IG32 0xE // 32-bit Interrupt Gate
0801 #define STS_TG32 0xF // 32-bit Trap Gate
0802
0803 // A virtual address 'la' has a three-part structure as follows:
0804 //
0805 // +-----10-----+-----10-----+-----12-----+
0806 // | Page Directory | Page Table | Offset within Page |
0807 // | Index | Index | |
0808 // +-----+-----+-----+
0809 // \--- PDX(va) ---/ \--- PTX(va) ---/
0810
0811 // page directory index
0812 #define PDX(va) (((uint)(va) >> PDXSHIFT) & 0x3FF)
0813
0814 // page table index
0815 #define PTX(va) (((uint)(va) >> PTXSHIFT) & 0x3FF)
0816
0817 // construct virtual address from indexes and offset
0818 #define PGADDR(d, t, o) ((uint)((d) << PDXSHIFT | (t) << PTXSHIFT | (o)))
0819
0820 // Page directory and page table constants.
0821 #define NPENTRIES 1024 // # directory entries per page directory
0822 #define NPTENTRIES 1024 // # PTEs per page table
0823 #define PGSIZE 4096 // bytes mapped by a page
0824
0825 #define PGSHIFT 12 // log2(PGSIZE)
0826 #define PTXSHIFT 12 // offset of PTX in a linear address
0827 #define PDXSHIFT 22 // offset of PDX in a linear address
0828
0829 #define PGROUNDUP(sz) (((sz)+PGSIZE-1) & ~(PGSIZE-1))
0830 #define PGROUNDOWN(a) (((a) & ~(PGSIZE-1))
0831
0832 // Page table/directory entry flags.
0833 #define PTE_P 0x001 // Present
0834 #define PTE_W 0x002 // Writable
0835 #define PTE_U 0x004 // User
0836 #define PTE_PWT 0x008 // Write-Through
0837 #define PTE_PCD 0x010 // Cache-Disable
0838 #define PTE_A 0x020 // Accessed
0839 #define PTE_D 0x040 // Dirty
0840 #define PTE_PS 0x080 // Page Size
0841 #define PTE_MBZ 0x180 // Bits must be zero
0842
0843 // Address in page table or page directory entry
0844 #define PTE_ADDR(pte) ((uint)(pte) & ~0xFFF)
0845
0846 #ifndef __ASSEMBLER__
0847 typedef uint pte_t;
0848
0849

```

```

0850 // Task state segment format
0851 struct taskstate {
0852   uint link; // Old ts selector
0853   uint esp0; // Stack pointers and segment selectors
0854   ushort ss0; // after an increase in privilege level
0855   ushort padding1;
0856   uint *esp1;
0857   ushort ss1;
0858   ushort padding2;
0859   uint *esp2;
0860   ushort ss2;
0861   ushort padding3;
0862   void *cr3; // Page directory base
0863   uint *eip; // Saved state from last task switch
0864   uint eflags;
0865   uint eax; // More saved state (registers)
0866   uint ecx;
0867   uint edx;
0868   uint ebx;
0869   uint *esp;
0870   uint *ebp;
0871   uint esi;
0872   uint edi; // Even more saved state (segment selectors)
0873   ushort es;
0874   ushort padding4;
0875   ushort cs;
0876   ushort padding5;
0877   ushort ss;
0878   ushort padding6;
0879   ushort ds;
0880   ushort padding7;
0881   ushort fs;
0882   ushort padding8;
0883   ushort gs;
0884   ushort padding9;
0885   ushort ldt;
0886   ushort padding10;
0887   ushort t; // Trap on task switch
0888   ushort iomb; // I/O map base address
0889 };
0890
0891
0892
0893
0894
0895
0896
0897
0898
0899

```

```

0900 // Gate descriptors for interrupts and traps
0901 struct gatedesc {
0902     uint off_15_0 : 16; // low 16 bits of offset in segment
0903     uint cs : 16; // code segment selector
0904     uint args : 5; // # args, 0 for interrupt/trap gates
0905     uint rsv1 : 3; // reserved (should be zero I guess)
0906     uint type : 4; // type(STS_{TG,IG32,TG32})
0907     uint s : 1; // must be 0 (system)
0908     uint dpl : 2; // descriptor (meaning new) privilege level
0909     uint p : 1; // Present
0910     uint off_31_16 : 16; // high bits of offset in segment
0911 };
0912
0913 // Set up a normal interrupt/trap gate descriptor.
0914 // -- istrap: 1 for a trap (= exception) gate, 0 for an interrupt gate.
0915 // -- interrupt gate clears FL_IF, trap gate leaves FL_IF alone
0916 // -- sel: Code segment selector for interrupt/trap handler
0917 // -- off: Offset in code segment for interrupt/trap handler
0918 // -- dpl: Descriptor Privilege Level -
0919 //       the privilege level required for software to invoke
0920 //       this interrupt/trap gate explicitly using an int instruction.
0921 #define SETGATE(gate, istrap, sel, off, d) \
0922 { \
0923     (gate).off_15_0 = (uint)(off) & 0xffff; \
0924     (gate).cs = (sel); \
0925     (gate).args = 0; \
0926     (gate).rsv1 = 0; \
0927     (gate).type = (istrap) ? STS_TG32 : STS_IG32; \
0928     (gate).s = 0; \
0929     (gate).dpl = (d); \
0930     (gate).p = 1; \
0931     (gate).off_31_16 = (uint)(off) >> 16; \
0932 }
0933
0934 #endif
0935
0936
0937
0938
0939
0940
0941
0942
0943
0944
0945
0946
0947
0948
0949

```

```

0950 // Format of an ELF executable file
0951
0952 #define ELF_MAGIC 0x464C457FU // "\x7FELF" in little endian
0953
0954 // File header
0955 struct elfhdr {
0956     uint magic; // must equal ELF_MAGIC
0957     uchar elf[12];
0958     ushort type;
0959     ushort machine;
0960     uint version;
0961     uint entry;
0962     uint phoff;
0963     uint shoff;
0964     uint flags;
0965     ushort ehsize;
0966     ushort phentsize;
0967     ushort phnum;
0968     ushort shentsize;
0969     ushort shnum;
0970     ushort shstrndx;
0971 };
0972
0973 // Program section header
0974 struct proghdr {
0975     uint type;
0976     uint off;
0977     uint vaddr;
0978     uint paddr;
0979     uint filesz;
0980     uint memsz;
0981     uint flags;
0982     uint align;
0983 };
0984
0985 // Values for Proghdr type
0986 #define ELF_PROG_LOAD 1
0987
0988 // Flag bits for Proghdr flags
0989 #define ELF_PROG_FLAG_EXEC 1
0990 #define ELF_PROG_FLAG_WRITE 2
0991 #define ELF_PROG_FLAG_READ 4
0992
0993
0994
0995
0996
0997
0998
0999

```

```

1000 # Multiboot header, for multiboot boot loaders like GNU Grub.
1001 # http://www.gnu.org/software/grub/manual/multiboot/multiboot.html
1002 #
1003 # Using GRUB 2, you can boot xv6 from a file stored in a
1004 # Linux file system by copying kernel or kernelmemfs to /boot
1005 # and then adding this menu entry:
1006 #
1007 # menuentry "xv6" {
1008 #   insmod ext2
1009 #   set root='(hd0,msdos1)'
1010 #   set kernel="/boot/kernel"
1011 #   echo "Loading ${kernel}..."
1012 #   multiboot ${kernel} ${kernel}
1013 #   boot
1014 # }
1015
1016 #include "asm.h"
1017 #include "memlayout.h"
1018 #include "mmu.h"
1019 #include "param.h"
1020
1021 # Multiboot header. Data to direct multiboot loader.
1022 .p2align 2
1023 .text
1024 .globl multiboot_header
1025 multiboot_header:
1026   #define magic 0x1badb002
1027   #define flags 0
1028   .long magic
1029   .long flags
1030   .long (-magic-flags)
1031
1032 # By convention, the _start symbol specifies the ELF entry point.
1033 # Since we haven't set up virtual memory yet, our entry point is
1034 # the physical address of 'entry'.
1035 .globl _start
1036 _start = V2P_W0(entry)
1037
1038 # Entering xv6 on boot processor, with paging off.
1039 .globl entry
1040 entry:
1041   # Turn on page size extension for 4Mbyte pages
1042   movl %cr4, %eax
1043   orl $(CR4_PSE), %eax
1044   movl %eax, %cr4
1045   # Set page directory
1046   movl $(V2P_W0(entrypgdir)), %eax
1047   movl %eax, %cr3
1048   # Turn on paging.
1049   movl %cr0, %eax

```

```

1050   orl $(CR0_PG|CR0_WP), %eax
1051   movl %eax, %cr0
1052
1053   # Set up the stack pointer.
1054   movl $(stack + KSTACKSIZE), %esp
1055
1056   # Jump to main(), and switch to executing at
1057   # high addresses. The indirect call is needed because
1058   # the assembler produces a PC-relative instruction
1059   # for a direct jump.
1060   mov $main, %eax
1061   jmp *%eax
1062
1063   .comm stack, KSTACKSIZE
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099

```



```

1100 #include "asm.h"
1101 #include "memlayout.h"
1102 #include "mmu.h"
1103
1104 # Each non-boot CPU ("AP") is started up in response to a STARTUP
1105 # IPI from the boot CPU. Section B.4.2 of the Multi-Processor
1106 # Specification says that the AP will start in real mode with CS:IP
1107 # set to XY00:0000, where XY is an 8-bit value sent with the
1108 # STARTUP. Thus this code must start at a 4096-byte boundary.
1109 #
1110 # Because this code sets DS to zero, it must sit
1111 # at an address in the low 2^16 bytes.
1112 #
1113 # Startothers (in main.c) sends the STARTUPs one at a time.
1114 # It copies this code (start) at 0x7000. It puts the address of
1115 # a newly allocated per-core stack in start-4, the address of the
1116 # place to jump to (mpenter) in start-8, and the physical address
1117 # of entrypgdir in start-12.
1118 #
1119 # This code is identical to bootasm.S except:
1120 # - it does not need to enable A20
1121 # - it uses the address at start-4, start-8, and start-12
1122
1123 .code16
1124 .globl start
1125 start:
1126 cli
1127
1128 xorw %ax,%ax
1129 movw %ax,%ds
1130 movw %ax,%es
1131 movw %ax,%ss
1132
1133 lgdt gdtdesc
1134 movl %cr0, %eax
1135 orl $CR0_PE, %eax
1136 movl %eax, %cr0
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149

```

```

1150 jmp1 $(SEG_KCODE<<3), $(start32)
1151
1152 .code32
1153 start32:
1154 movw $(SEG_KDATA<<3), %ax
1155 movw %ax,%ds
1156 movw %ax,%es
1157 movw %ax,%ss
1158 movw $0, %ax
1159 movw %ax,%fs
1160 movw %ax,%gs
1161
1162 # Turn on page size extension for 4Mbyte pages
1163 movl %cr4, %eax
1164 orl $(CR4_PSE), %eax
1165 movl %eax, %cr4
1166 # Use enterpgdir as our initial page table
1167 movl (start-12), %eax
1168 movl %eax, %cr3
1169 # Turn on paging.
1170 movl %cr0, %eax
1171 orl $(CR0_PE|CR0_PG|CR0_WP), %eax
1172 movl %eax, %cr0
1173
1174 # Switch to the stack allocated by startothers()
1175 movl (start-4), %esp
1176 # Call mpenter()
1177 call *(start-8)
1178
1179 movw $0x8a00, %ax
1180 movw %ax, %dx
1181 outw %ax, %dx
1182 movw $0x8ae0, %ax
1183 outw %ax, %dx
1184 spin:
1185 jmp spin
1186
1187 .p2align 2
1188 gdt:
1189 SEG_NULLASM
1190 SEG_ASM(STA_X|STA_R, 0, 0xffffffff)
1191 SEG_ASM(STA_W, 0, 0xffffffff)
1192
1193
1194 gdtdesc:
1195 .word (gdtdesc - gdt - 1)
1196 .long gdt
1197
1198
1199

```

```

1200 #include "types.h"
1201 #include "defs.h"
1202 #include "param.h"
1203 #include "memlayout.h"
1204 #include "mmu.h"
1205 #include "proc.h"
1206 #include "x86.h"
1207
1208 static void startothers(void);
1209 static void mpmain(void) __attribute__((noreturn));
1210 extern pde_t *kpgdir;
1211 extern char end[]; // first address after kernel loaded from ELF file
1212
1213 // Bootstrap processor starts running C code here.
1214 // Allocate a real stack and switch to it, first
1215 // doing some setup required for memory allocator to work.
1216 int
1217 main(void)
1218 {
1219     kinit1(end, P2V(4*1024*1024)); // phys page allocator
1220     kvmalloc(); // kernel page table
1221     mpinit(); // collect info about this machine
1222     lapicinit(); // set up segments
1223     seginit(); // interrupt controller
1224     printf("ncpu%d: starting xv6\n", ncpu->id);
1225     picinit(); // another interrupt controller
1226     ioapicinit(); // I/O devices & their interrupts
1227     consoleinit(); // serial port
1228     uartinit(); // process table
1229     pinit(); // trap vectors
1230     tvinit(); // buffer cache
1231     binit(); // file table
1232     fileinit(); // inode cache
1233     iinit(); // disk
1234     ideinit(); // disk
1235     if(!ismp)
1236         timerinit(); // uniprocessor timer
1237     startothers(); // start other processors
1238     kinit2(P2V(4*1024*1024), P2V(PHYSTOP)); // must come after startothers()
1239     userinit(); // first user process
1240     // Finish setting up this processor in mpmain.
1241     mpmain();
1242 }
1243
1244
1245
1246
1247
1248
1249

```

```

1250 // Other CPUs jump here from entryother.S.
1251 static void
1252 mpenter(void)
1253 {
1254     switchvm();
1255     seginit();
1256     lapicinit();
1257     mpmain();
1258 }
1259
1260 // Common CPU setup code.
1261 static void
1262 mpmain(void)
1263 {
1264     printf("cpu%d: starting\n", cpu->id);
1265     idtinit(); // load idt register
1266     xchg(&cpu->started, 1); // tell startothers() we're up
1267     scheduler(); // start running processes
1268 }
1269
1270 pde_t entrypgdir[]; // For entry.S
1271
1272 // Start the non-boot (AP) processors.
1273 static void
1274 startothers(void)
1275 {
1276     extern uchar _binary_entryother_start[], _binary_entryother_size[];
1277     uchar *code;
1278     struct cpu *c;
1279     char *stack;
1280
1281     // Write entry code to unused memory at 0x7000.
1282     // The linker has placed the image of entryother.S in
1283     // _binary_entryother_start.
1284     code = p2v(0x7000);
1285     memmove(code, _binary_entryother_start, (uint)_binary_entryother_size);
1286
1287     for(c = cpus; c < cpus+ncpu; c++){
1288         if(c == cpus+cpunum()) // We've started already.
1289             continue;
1290         // Tell entryother.S what stack to use, where to enter, and what
1291         // pgdir to use. We cannot use kpgdir yet, because the AP processor
1292         // is running in low memory, so we use entrypgdir for the APs too.
1293         stack = kalloc();
1294         *(void**) (code-4) = stack + KSTACKSIZE;
1295         *(void**) (code-8) = mpenter;
1296         *(int**) (code-12) = (void *) v2p(entrypgdir);
1297
1298
1299         lapicstartap(c->id, v2p(code));

```

```
1300 // wait for cpu to finish mpmain()
1301 while(c->started == 0)
1302 ;
1303 }
1304 }
1305
1306 // Boot page table used in entry.S and entryother.S.
1307 // Page directories (and page tables), must start on a page boundary,
1308 // hence the "__aligned__" attribute.
1309 // Use PTE_PS in page directory entry to enable 4Mbyte pages.
1310 __attribute__((__aligned__(PGSIZE)))
1311 pde_t entrypgdir[NPDENTRIES] = {
1312 // Map VA's [0, 4MB) to PA's [0, 4MB)
1313 [0] = (0) | PTE_P | PTE_W | PTE_PS,
1314 // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
1315 [KERNBASE>>PDXSHIFT] = (0) | PTE_P | PTE_W | PTE_PS,
1316 };
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
```

```
1350 // Blank page.
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
```

```

1400 // Mutual exclusion lock.
1401 struct spinlock {
1402     uint locked;           // Is the lock held?

1403     // For debugging:
1404     char *name;            // Name of lock.
1405     struct cpu *cpu;        // The cpu holding the lock.
1406     uint pcs[10];          // The call stack (an array of program counters)
1407                             // that locked the lock.
1408 };
1409 };
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449

```

```

1450 // Mutual exclusion spin locks.
1451
1452 #include "types.h"
1453 #include "defs.h"
1454 #include "param.h"
1455 #include "x86.h"
1456 #include "memlayout.h"
1457 #include "mmu.h"
1458 #include "proc.h"
1459 #include "spinlock.h"
1460
1461 void
1462 initlock(struct spinlock *lk, char *name)
1463 {
1464     lk->name = name;
1465     lk->locked = 0;
1466     lk->cpu = 0;
1467 }
1468
1469 // Acquire the lock.
1470 // Loops (spins) until the lock is acquired.
1471 // Holding a lock for a long time may cause
1472 // other CPUs to waste time spinning to acquire it.
1473 void
1474 acquire(struct spinlock *lk)
1475 {
1476     pushcli(); // disable interrupts to avoid deadlock.
1477     if(holding(lk))
1478         panic("acquire");
1479
1480     // The xchg is atomic.
1481     // It also serializes, so that reads after acquire are not
1482     // reordered before it.
1483     while(xchg(&lk->locked, 1) != 0)
1484         ;
1485
1486     // Record info about lock acquisition for debugging.
1487     lk->cpu = cpu;
1488     getcallerpcs(&lk, lk->pcs);
1489 }
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499

```

```

1500 // Release the lock.
1501 void
1502 release(struct spinlock *lk)
1503 {
1504     if(holding(lk))
1505         panic("release");
1506
1507     lk->pcs[0] = 0;
1508     lk->cpu = 0;
1509
1510     // The xchg serializes, so that reads before release are
1511     // not reordered after it. The 1996 PentiumPro manual (Volume 3,
1512     // 7.2) says reads can be carried out speculatively and in
1513     // any order, which implies we need to serialize here.
1514     // But the 2007 Intel 64 Architecture Memory Ordering White
1515     // Paper says that Intel 64 and IA-32 will not move a load
1516     // after a store. So lock->locked = 0 would work here.
1517     // The xchg being asm volatile ensures gcc emits it after
1518     // the above assignments (and after the critical section).
1519     xchg(&lk->locked, 0);
1520
1521     popcli();
1522 }
1523
1524 // Record the current call stack in pcs[] by following the %ebp chain.
1525 void
1526 getcallerpcs(void *v, uint pcs[])
1527 {
1528     uint *ebp;
1529     int i;
1530
1531     ebp = (uint*)v - 2;
1532     for(i = 0; i < 10; i++){
1533         if(ebp == 0 || ebp < (uint*)KERNBASE || ebp == (uint*)0xffffffff)
1534             break;
1535         pcs[i] = ebp[i]; // saved %ebp
1536         ebp = (uint*)ebp[0]; // saved %ebp
1537     }
1538     for(; i < 10; i++)
1539         pcs[i] = 0;
1540 }
1541
1542 // Check whether this cpu is holding the lock.
1543 int
1544 holding(struct spinlock *lock)
1545 {
1546     return lock->locked && lock->cpu == cpu;
1547 }
1548
1549

```

```

1550 // Pushcli/popcli are like cli/sti except that they are matched:
1551 // it takes two popcli to undo two pushcli. Also, if interrupts
1552 // are off, then pushcli, popcli leaves them off.
1553
1554 void
1555 pushcli(void)
1556 {
1557     int eflags;
1558
1559     eflags = readeflags();
1560     cli();
1561     if(cpu->ncli++ == 0)
1562         cpu->intena = eflags & FL_IF;
1563 }
1564
1565 void
1566 popcli(void)
1567 {
1568     if(readeflags() & FL_IF)
1569         panic("popcli - interruptible");
1570     if(--cpu->ncli < 0)
1571         panic("popcli");
1572     if(cpu->ncli == 0 && cpu->intena)
1573         sti();
1574 }
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599

```

```

1600 #include "param.h"
1601 #include "types.h"
1602 #include "defs.h"
1603 #include "x86.h"
1604 #include "memlayout.h"
1605 #include "mmu.h"
1606 #include "proc.h"
1607 #include "elf.h"
1608
1609 extern char data[]; // defined by kernel.ld
1610 pde_t *pgdir; // for use in scheduler()
1611 struct segdesc gdt[NSEGS];
1612
1613 // Set up CPU's kernel segment descriptors.
1614 // Run once on entry on each CPU.
1615 void
1616 seginit(void)
1617 {
1618     struct cpu *c;
1619
1620     // Map "logical" addresses to virtual addresses using identity map.
1621     // Cannot share a CODE descriptor for both kernel and user
1622     // because it would have to have DPL_USR, but the CPU forbids
1623     // an interrupt from CPL=0 to DPL=3.
1624     c = &cpu[cpunum()];
1625     c->gdt[SEG_KCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, 0);
1626     c->gdt[SEG_XDATA] = SEG(STA_W, 0, 0xffffffff, 0);
1627     c->gdt[SEG_UCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, DPL_USER);
1628     c->gdt[SEG_UDATA] = SEG(STA_W, 0, 0xffffffff, DPL_USER);
1629
1630     // Map cpu, and curproc
1631     c->gdt[SEG_KCPU] = SEG(STA_W, &c->cpu, 8, 0);
1632
1633     lgdt(c->gdt, sizeof(c->gdt));
1634     loadgs(SEG_KCPU << 3);
1635
1636     // Initialize cpu-local storage.
1637     cpu = c;
1638     proc = 0;
1639 }
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649

```

```

1650 // Return the address of the PTE in page table pgdir
1651 // that corresponds to virtual address va. If alloc!=0,
1652 // create any required page table pages.
1653 static pte_t *
1654 walkpgdir(pde_t *pgdir, const void *va, int alloc)
1655 {
1656     pde_t *pde;
1657     pte_t *pgtab;
1658
1659     pde = &pgdir[PDX(va)];
1660     if(*pde & PTE_P){
1661         pgtab = (pte_t*)p2v(PTE_ADDR(*pde));
1662     } else {
1663         if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
1664             return 0;
1665         // Make sure all those PTE_P bits are zero.
1666         memset(pgtab, 0, PGSIZE);
1667         // The permissions here are overly generous, but they can
1668         // be further restricted by the permissions in the page table
1669         // entries, if necessary.
1670         *pde = v2p(pgtab) | PTE_P | PTE_W | PTE_U;
1671     }
1672     return &pgtab[PTX(va)];
1673 }
1674
1675 // Create PTEs for virtual addresses starting at va that refer to
1676 // physical addresses starting at pa. va and size might not
1677 // be page-aligned.
1678 static int
1679 mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
1680 {
1681     char *a, *last;
1682     pte_t *pte;
1683
1684     a = (char*)PGROUNDDOWN((uint)va);
1685     last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
1686     for(;;){
1687         if((pte = walkpgdir(pgdir, a, 1)) == 0)
1688             return -1;
1689         if(*pte & PTE_P)
1690             panic("remap");
1691         *pte = pa | perm | PTE_P;
1692         if(a == last)
1693             break;
1694         a += PGSIZE;
1695         pa += PGSIZE;
1696     }
1697     return 0;
1698 }
1699

```



```

1700 // There is one page table per process, plus one that's used when
1701 // a CPU is not running any process (kpgdir). The kernel uses the
1702 // current process's page table during system calls and interrupts;
1703 // page protection bits prevent user code from using the kernel's
1704 // mappings.
1705 //
1706 // setupkvm() and exec() set up every page table like this:
1707 //
1708 // 0..KERNBASE: user memory (text+data+stack+heap), mapped to
1709 // phys memory allocated by the kernel
1710 // KERNBASE..KERNBASE+EXTMEM: mapped to 0..EXTMEM (for I/O space)
1711 // KERNBASE+EXTMEM..data: mapped to EXTMEM..V2P(data)
1712 // for the kernel's instructions and r/o data
1713 // data..KERNBASE+PHYSTOP: mapped to V2P(data)..PHYSTOP,
1714 // rw data + free physical memory
1715 // 0xfe000000..0: mapped direct (devices such as ioapic)
1716 //
1717 // The kernel allocates physical memory for its heap and for user memory
1718 // between V2P(end) and the end of physical memory (PHYSTOP)
1719 // (directly addressable from end..P2V(PHYSTOP)).
1720 //
1721 // This table defines the kernel's mappings, which are present in
1722 // every process's page table.
1723 static struct kmap {
1724   void *virt;
1725   uint phys_start;
1726   uint phys_end;
1727   int perm;
1728 } kmap[] = {
1729   { (void*)KERNBASE, 0, EXTMEM, PTE_W}, // I/O space
1730   { (void*)KERNLINK, V2P(KERNLINK), V2P(data), 0}, // kern text+rodata
1731   { (void*)data, V2P(data), PHYSTOP, PTE_W}, // kern data+memory
1732   { (void*)DEVSPACE, DEVSPACE, 0, PTE_W}, // more devices
1733 };
1734
1735 // Set up kernel part of a page table.
1736 pde_t*
1737 setupkvm(void)
1738 {
1739   pde_t *pgdir;
1740   struct kmap *k;
1741
1742   if((pgdir = (pde_t*)kalloc()) == 0)
1743     return 0;
1744   memset(pgdir, 0, PGSIZE);
1745   if (p2v(PHYSTOP) > (void*)DEVSPACE)
1746     panic("PHYSTOP too high");
1747   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1748     if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
1749               (uint)k->phys_start, k->perm) < 0)

```

```

1750     return 0;
1751     return pgdir;
1752 }
1753
1754 // Allocate one page table for the machine for the kernel address
1755 // space for scheduler processes.
1756 void
1757 kvmalloc(void)
1758 {
1759   kpgdir = setupkvm();
1760   switchkvm();
1761 }
1762
1763 // Switch h/w page table register to the kernel-only page table,
1764 // for when no process is running.
1765 void
1766 switchkvm(void)
1767 {
1768   lcr3(v2p(kpgdir)); // switch to the kernel page table
1769 }
1770
1771 // Switch TSS and h/w page table to correspond to process p.
1772 void
1773 switchvm(struct proc *p)
1774 {
1775   pushcli();
1776   cpu->gdt[SEG_TSS].s = 0;
1777   cpu->ts.s0 = SEG_KDATA << 3;
1778   cpu->ts.esp0 = (uint)proc->stack + KSTACKSIZE;
1779   ltr(SEG_TSS << 3);
1780   if(p->pgdir == 0)
1781     panic("switchvm: no pgdir");
1782   lcr3(v2p(p->pgdir)); // switch to new address space
1783   popcli();
1784 }
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799

```

```

1800 // Load the initcode into address 0 of pgdir.
1801 // sz must be less than a page.
1802 void
1803 initvm(pde_t *pgdir, char *init, uint sz)
1804 {
1805     char *mem;
1806
1807     if(sz >= PGSIZE)
1808         panic("initvm: more than a page");
1809     mem = kalloc();
1810     memset(mem, 0, PGSIZE);
1811     mappages(pgdir, 0, PGSIZE, v2p(mem), PTE_W|PTE_U);
1812     memmove(mem, init, sz);
1813 }
1814
1815 // Load a program segment into pgdir. addr must be page-aligned
1816 // and the pages from addr to addr+sz must already be mapped.
1817 int
1818 loadvm(pde_t *pgdir, char *addr, struct inode *ip, uint offset, uint sz)
1819 {
1820     uint i, pa, n;
1821     pte_t *pte;
1822
1823     if((uint) addr % PGSIZE != 0)
1824         panic("loadvm: addr must be page aligned");
1825     for(i = 0; i < sz; i += PGSIZE){
1826         if(pte = walkpgdir(pgdir, addr+i, 0)) == 0)
1827             panic("loadvm: address should exist");
1828         pa = PTE_ADDR(*pte);
1829         if(sz - i < PGSIZE)
1830             n = sz - i;
1831         else
1832             n = PGSIZE;
1833         if(readi(ip, p2v(pa), offset+i, n) != n)
1834             return -1;
1835     }
1836     return 0;
1837 }
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849

```

```

1850 // Allocate page tables and physical memory to grow process from oldsz to
1851 // newsz, which need not be page aligned. Returns new size or 0 on error.
1852 int
1853 allocvm(pde_t *pgdir, uint oldsz, uint newsz)
1854 {
1855     char *mem;
1856     uint a;
1857
1858     if(newsz >= KERNBASE)
1859         return 0;
1860     if(newsz < oldsz)
1861         return oldsz;
1862
1863     a = PGROUNDUP(oldsz);
1864     for(; a < newsz; a += PGSIZE){
1865         mem = kalloc();
1866         if(mem == 0){
1867             cprintf("allocvm out of memory\n");
1868             deallocvm(pgdir, newsz, oldsz);
1869             return 0;
1870         }
1871         memset(mem, 0, PGSIZE);
1872         mappages(pgdir, (char*)a, PGSIZE, v2p(mem), PTE_W|PTE_U);
1873     }
1874     return newsz;
1875 }
1876
1877 // Deallocate user pages to bring the process size from oldsz to
1878 // newsz. oldsz and newsz need not be page-aligned, nor does newsz
1879 // need to be less than oldsz. oldsz can be larger than the actual
1880 // process size. Returns the new process size.
1881 int
1882 deallocvm(pde_t *pgdir, uint oldsz, uint newsz)
1883 {
1884     pte_t *pte;
1885     uint a, pa;
1886
1887     if(newsz >= oldsz)
1888         return oldsz;
1889
1890     a = PGROUNDUP(newsz);
1891     for(; a < oldsz; a += PGSIZE){
1892         pte = walkpgdir(pgdir, (char*)a, 0);
1893         if(!pte)
1894             a += (NPENTRIES - 1) * PGSIZE;
1895         else if((*pte & PTE_P) != 0){
1896             pa = PTE_ADDR(*pte);
1897             if(pa == 0)
1898                 panic("kfree");
1899             char *v = p2v(pa);

```

```

1900     kfree(v);
1901     *pte = 0;
1902 }
1903 }
1904 return newsz;
1905 }
1906
1907 // Free a page table and all the physical memory pages
1908 // in the user part.
1909 void
1910 freevm(pde_t *pgdir)
1911 {
1912     uint i;
1913
1914     if(pgdir == 0)
1915         panic("freevm: no pgdir");
1916     deallocvm(pgdir, KERNBASE, 0);
1917     for(i = 0; i < NPENTRIES; i++){
1918         if(pgdir[i] & PTE_P){
1919             char *v = p2v(PTE_ADDR(pgdir[i]));
1920             kfree(v);
1921         }
1922     }
1923     kfree((char*)pgdir);
1924 }
1925
1926 // Clear PTE_U on a page. Used to create an inaccessible
1927 // page beneath the user stack.
1928 void
1929 clearpteu(pde_t *pgdir, char *uva)
1930 {
1931     pte_t *pte;
1932
1933     pte = walkpgdir(pgdir, uva, 0);
1934     if(pte == 0)
1935         panic("clearpteu");
1936     *pte &= ~PTE_U;
1937 }
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949

```

```

1950 // Given a parent process's page table, create a copy
1951 // of it for a child.
1952 pde_t*
1953 copyvm(pde_t *pgdir, uint sz)
1954 {
1955     pde_t *d;
1956     pte_t *pte;
1957     uint pa, i;
1958     char *mem;
1959
1960     if((d = setupvm()) == 0)
1961         return 0;
1962     for(i = 0; i < sz; i += PGSIZE){
1963         if(pte = walkpgdir(pgdir, (void *) i, 0)) == 0)
1964             panic("copyvm: pte should exist");
1965         if(!(*pte & PTE_P))
1966             panic("copyvm: page not present");
1967         pa = PTE_ADDR(*pte);
1968         if((mem = kalloc()) == 0)
1969             goto bad;
1970         memmove(mem, (char*)p2v(pa), PGSIZE);
1971         if(mappages(d, (void*)i, PGSIZE, v2p(mem), PTE_W|PTE_U) < 0)
1972             goto bad;
1973     }
1974     return d;
1975
1976 bad:
1977     freevm(d);
1978     return 0;
1979 }
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999

```

```

2000 // Map user virtual address to kernel address.
2001 char*
2002 uva2ka(pde_t *pgdir, char *uva)
2003 {
2004     pte_t *pte;
2005
2006     pte = walkpgdir(pgdir, uva, 0);
2007     if ((*pte & PTE_P) == 0)
2008         return 0;
2009     if ((*pte & PTE_U) == 0)
2010         return 0;
2011     return (char*)p2v(PTE_ADDR(*pte));
2012 }
2013
2014 // Copy len bytes from p to user address va in page table pgdir.
2015 // Most useful when pgdir is not the current page table.
2016 // uva2ka ensures this only works for PTE_U pages.
2017 int
2018 copyout(pde_t *pgdir, uint va, void *p, uint len)
2019 {
2020     char *buf, *pa0;
2021     uint n, va0;
2022
2023     buf = (char*)p;
2024     while(len > 0){
2025         va0 = (uint)PGROUNDDOWN(va);
2026         pa0 = uva2ka(pgdir, (char*)va0);
2027         if(pa0 == 0)
2028             return -1;
2029         n = PGSIZE - (va - va0);
2030         if(n > len)
2031             n = len;
2032         memmove(pa0 + (va - va0), buf, n);
2033         len -= n;
2034         buf += n;
2035         va = va0 + PGSIZE;
2036     }
2037     return 0;
2038 }
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049

```

```

2050 // Segments in proc->gdt.
2051 #define NSEGS 7
2052
2053 // Per-CPU state
2054 struct cpu {
2055     uchar id; // Local APIC ID; index into cpus[] below
2056     struct context *scheduler; // switch() here to enter scheduler
2057     struct taskstate ts; // Used by x86 to find stack for interrupt
2058     struct segdesc gdt[NSEGS]; // x86 global descriptor table
2059     volatile uint started; // Has the CPU started?
2060     int ncli; // Depth of pushcli nesting.
2061     int intena; // Were interrupts enabled before pushcli?
2062 }
2063
2064 // Cpu-local storage variables; see below
2065 struct cpu *cpu;
2066 struct proc *proc; // The currently-running process.
2067 };
2068
2069 extern struct cpu cpus[NCPU];
2070 extern int ncpu;
2071
2072 // Per-CPU variables, holding pointers to the
2073 // current cpu and to the current process.
2074 // The asm suffix tells gcc to use "%gs:0" to refer to cpu
2075 // and "%gs:4" to refer to proc. seginit sets up the
2076 // %gs segment register so that %gs refers to the memory
2077 // holding those two variables in the local cpu's struct cpu.
2078 // This is similar to how thread-local variables are implemented
2079 // in thread libraries such as Linux pthreads.
2080 extern struct cpu *cpu asm("%gs:0"); // &cpus[cpunum()]
2081 extern struct proc *proc asm("%gs:4"); // cpus[cpunum()].proc
2082
2083 // Saved registers for kernel context switches.
2084 // Don't need to save all the segment registers (%cs, etc),
2085 // because they are constant across kernel contexts.
2086 // Don't need to save %eax, %ecx, %edx, because the
2087 // x86 convention is that the caller has saved them.
2088 // Contexts are stored at the bottom of the stack they
2089 // describe; the stack pointer is the address of the context.
2090 // The layout of the context matches the layout of the stack in switch.S
2091 // at the "Switch stacks" comment. Switch doesn't save eip explicitly,
2092 // but it is on the stack and allocproc() manipulates it.
2093 struct context {
2094     uint edi;
2095     uint esi;
2096     uint ebx;
2097     uint ebp;
2098     uint eip;
2099 };

```

```
2100 enum procstate { UNUSED, EMBRYO, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
2101
2102 // Per-process state
2103 struct proc {
2104   uint sz; // Size of process memory (bytes)
2105   pde_t* pgdir; // Page table
2106   char *kstack; // Bottom of kernel stack for this process
2107   enum procstate state; // Process state
2108   volatile int pid; // Process ID
2109   struct proc *parent; // Parent process
2110   struct trapframe *tf; // Trap frame for current syscall
2111   struct context *context; // swtch() here to run process
2112   void *chan; // If non-zero, sleeping on chan
2113   int killed; // If non-zero, have been killed
2114   struct file *ofile[Nofile]; // Open files
2115   struct inode *cwd; // Current directory
2116   char name[16]; // Process name (debugging)
2117 };
2118
```

```
2119 // Process memory is laid out contiguously, low addresses first:
```

```
2120 // text
2121 // original data and bss
2122 // fixed-size stack
2123 // expandable heap
2124
```

```
2150 #include "types.h"
2151 #include "defs.h"
2152 #include "param.h"
2153 #include "memlayout.h"
2154 #include "mmu.h"
2155 #include "x86.h"
2156 #include "proc.h"
2157 #include "spinlock.h"
2158
2159 struct {
2160   struct spinlock lock;
2161   struct proc proc[NPROC];
2162 } ptable;
2163
2164 static struct proc *initproc;
2165
2166 int nextpid = 1;
2167 extern void forkret(void);
2168 extern void trapret(void);
2169
2170 static void wakeup1(void *chan);
2171
2172 void
2173 pinit(void)
2174 {
2175   initlock(&ptable.lock, "ptable");
2176 }
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
```

```

2200 // Look in the process table for an UNUSED proc.
2201 // If found, change state to EMBRYO and initialize
2202 // state required to run in the kernel.
2203 // Otherwise return 0.
2204 static struct proc*
2205 allocproc(void)
2206 {
2207     struct proc *p;
2208     char *sp;
2209
2210     acquire(&ptable.lock);
2211     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2212         if(p->state == UNUSED)
2213             goto found;
2214         release(&ptable.lock);
2215         return 0;
2216     found:
2217     p->state = EMBRYO;
2218     p->pid = nextpid++;
2219     release(&ptable.lock);
2220
2221     // Allocate kernel stack.
2222     if((p->kstack = kalloc()) == 0){
2223         p->state = UNUSED;
2224         return 0;
2225     }
2226     sp = p->kstack + KSTACKSIZE;
2227
2228     // Leave room for trap frame.
2229     sp -= sizeof *p->tf;
2230     p->tf = (struct trapframe*)sp;
2231
2232     // Set up new context to start executing at forkret,
2233     // which returns to trapret.
2234     sp -= 4;
2235     *(uint*)sp = (uint)trapret;
2236
2237     sp -= sizeof *p->context;
2238     p->context = (struct context*)sp;
2239     memset(p->context, 0, sizeof *p->context);
2240     p->context->eip = (uint)forkret;
2241
2242     return p;
2243 }
2244
2245
2246
2247
2248
2249

```

```

2250 // Set up first user process.
2251 void
2252 userinit(void)
2253 {
2254     struct proc *p;
2255     extern char _binary_initcode_start[], _binary_initcode_size[];
2256
2257     p = allocproc();
2258     initproc = p;
2259     if((p->pgdir = setupkvm()) == 0)
2260         panic("userinit: out of memory?");
2261     inituvm(p->pgdir, _binary_initcode_start, (int)_binary_initcode_size);
2262     p->sz = PGSIZE;
2263     memset(p->tf, 0, sizeof(*p->tf));
2264     p->tf->cs = (SEG_UCODE << 3) | DPL_USER;
2265     p->tf->ds = (SEG_UDATA << 3) | DPL_USER;
2266     p->tf->es = p->tf->ds;
2267     p->tf->ss = p->tf->ds;
2268     p->tf->eflags = FL_IF;
2269     p->tf->esp = PGSIZE;
2270     p->tf->eip = 0; // beginning of initcode.S
2271
2272     safestrcpy(p->name, "initcode", sizeof(p->name));
2273     p->cwd = namei("/");
2274
2275     p->state = RUNNABLE;
2276 }
2277
2278 // Grow current process's memory by n bytes.
2279 // Return 0 on success, -1 on failure.
2280 int
2281 growproc(int n)
2282 {
2283     uint sz;
2284
2285     sz = proc->sz;
2286     if(n > 0){
2287         if((sz = allocuvm(proc->pgdir, sz, sz + n)) == 0)
2288             return -1;
2289     } else if(n < 0){
2290         if((sz = deallocuvm(proc->pgdir, sz, sz + n)) == 0)
2291             return -1;
2292     }
2293     proc->sz = sz;
2294     switchuvm(proc);
2295     return 0;
2296 }
2297
2298
2299

```



```

2300 // Create a new process copying p as the parent.
2301 // Sets up stack to return as if from system call.
2302 // Caller must set state of returned proc to RUNNABLE.
2303 int
2304 fork(void)
2305 {
2306     int i, pid;
2307     struct proc *np;
2308
2309     // Allocate process.
2310     if((np = allocproc()) == 0)
2311         return -1;
2312
2313     // Copy process state from p.
2314     if((np->pgdir = copyvm(proc->pgdir, proc->sz)) == 0){
2315         kfree(np->kstack);
2316         np->kstack = 0;
2317         np->state = UNUSED;
2318         return -1;
2319     }
2320     np->sz = proc->sz;
2321     np->parent = proc;
2322     *np->tf = *proc->tf;
2323
2324     // Clear %eax so that fork returns 0 in the child.
2325     np->tf->eax = 0;
2326
2327     for(i = 0; i < NOFILE; i++)
2328         if(proc->ofile[i])
2329             np->ofile[i] = filedup(proc->ofile[i]);
2330     np->cwd = idup(proc->cwd);
2331
2332     pid = np->pid;
2333     np->state = RUNNABLE;
2334     safestrcpy(np->name, proc->name, sizeof(proc->name));
2335     return pid;
2336 }
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349

```

```

2350 // Exit the current process. Does not return.
2351 // An exited process remains in the zombie state
2352 // until its parent calls wait() to find out it exited.
2353 void
2354 exit(void)
2355 {
2356     struct proc *p;
2357     int fd;
2358
2359     if(proc == initproc)
2360         panic("init exiting");
2361
2362     // Close all open files.
2363     for(fd = 0; fd < NOFILE; fd++){
2364         if(proc->ofile[fd]){
2365             fileclose(proc->ofile[fd]);
2366             proc->ofile[fd] = 0;
2367         }
2368     }
2369
2370     iput(proc->cwd);
2371     proc->cwd = 0;
2372
2373     acquire(&table.lock);
2374
2375     // Parent might be sleeping in wait().
2376     wakeup1(proc->parent);
2377
2378     // Pass abandoned children to init.
2379     for(p = table.proc; p < &table.proc[NPROC]; p++){
2380         if(p->parent == proc){
2381             p->parent = initproc;
2382             if(p->state == ZOMBIE)
2383                 wakeup1(initproc);
2384         }
2385     }
2386
2387     // Jump into the scheduler, never to return.
2388     proc->state = ZOMBIE;
2389     sched();
2390     panic("zombie exit");
2391 }
2392
2393
2394
2395
2396
2397
2398
2399

```

```

2400 // Wait for a child process to exit and return its pid.
2401 // Return -1 if this process has no children.
2402 int
2403 wait(void)
2404 {
2405     struct proc *p;
2406     int havekids, pid;
2407
2408     acquire(&ptable.lock);
2409     for(;;){
2410         // Scan through table looking for zombie children.
2411         havekids = 0;
2412         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2413             if(p->parent != proc)
2414                 continue;
2415             havekids = 1;
2416             if(p->state == ZOMBIE){
2417                 // Found one.
2418                 pid = p->pid;
2419                 kfree(p->kstack);
2420                 p->kstack = 0;
2421                 freevm(p->pgdir);
2422                 p->state = UNUSED;
2423                 p->pid = 0;
2424                 p->parent = 0;
2425                 p->name[0] = 0;
2426                 p->killed = 0;
2427                 release(&ptable.lock);
2428                 return pid;
2429             }
2430         }
2431
2432         // No point waiting if we don't have any children.
2433         if(!havekids || p->killed){
2434             release(&ptable.lock);
2435             return -1;
2436         }
2437
2438         // Wait for children to exit. (See wakeup1 call in proc_exit.)
2439         sleep(proc, &ptable.lock);
2440     }
2441 }
2442
2443
2444
2445
2446
2447
2448
2449

```

```

2450 // Per-CPU process scheduler.
2451 // Each CPU calls scheduler() after setting itself up.
2452 // Scheduler never returns. It loops, doing:
2453 //  - choose a process to run
2454 //  - switch to start running that process
2455 //  - eventually that process transfers control
2456 //    via switch back to the scheduler.
2457 void
2458 scheduler(void)
2459 {
2460     struct proc *p;
2461
2462     for(;;){
2463         // Enable interrupts on this processor.
2464         sti();
2465
2466         // Loop over process table looking for process to run.
2467         acquire(&ptable.lock);
2468         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2469             if(p->state != RUNNABLE)
2470                 continue;
2471
2472             // Switch to chosen process. It is the process's job
2473             // to release ptable.lock and then reacquire it
2474             // before jumping back to us.
2475             proc = p;
2476             switchvm(p);
2477             p->state = RUNNING;
2478             switch(&cpu->scheduler, proc->context);
2479             switchvm();
2480
2481             // Process is done running for now.
2482             // It should have changed its p->state before coming back.
2483             proc = 0;
2484         }
2485         release(&ptable.lock);
2486     }
2487 }
2488 }
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499

```

```

2500 // Enter scheduler. Must hold only ptable.lock
2501 // and have changed proc->state.
2502 void
2503 sched(void)
2504 {
2505     int intena;
2506
2507     if(!holding(&ptable.lock))
2508         panic("sched ptable.lock");
2509     if(cpu->ncli != 1)
2510         panic("sched locks");
2511     if(proc->state == RUNNING)
2512         panic("sched running");
2513     if(readeflags() & FL_IF)
2514         panic("sched interruptible");
2515     intena = cpu->intena;
2516     switch(&proc->context, cpu->scheduler);
2517     cpu->intena = intena;
2518 }
2519
2520 // Give up the CPU for one scheduling round.
2521 void
2522 yield(void)
2523 {
2524     acquire(&ptable.lock);
2525     proc->state = RUNNABLE;
2526     sched();
2527     release(&ptable.lock);
2528 }
2529
2530 // A fork child's very first scheduling by scheduler()
2531 // will switch here. "Return" to user space.
2532 void
2533 forkret(void)
2534 {
2535     static int first = 1;
2536     // Still holding ptable.lock from scheduler.
2537     release(&ptable.lock);
2538
2539     if (first) {
2540         // Some initialization functions must be run in the context
2541         // of a regular process (e.g., they call sleep), and thus cannot
2542         // be run from main().
2543         first = 0;
2544         initlog();
2545     }
2546
2547     // Return to "caller", actually trapret (see allocproc).
2548 }
2549

```

```

2550 // Atomically release lock and sleep on chan.
2551 // Reacquires lock when awakened.
2552 void
2553 sleep(void *chan, struct spinlock *lk)
2554 {
2555     if(proc == 0)
2556         panic("sleep");
2557     if(lk == 0)
2558         panic("sleep without lk");
2559
2560     // Must acquire ptable.lock in order to
2561     // change p->state and then call sched.
2562     // Once we hold ptable.lock, we can be
2563     // guaranteed that we won't miss any wakeup
2564     // (wakeup runs with ptable.lock locked),
2565     // so it's okay to release lk.
2566     if(lk != &ptable.lock){
2567         acquire(&ptable.lock);
2568         release(lk);
2569     }
2570
2571     // Go to sleep.
2572     proc->chan = chan;
2573     proc->state = SLEEPING;
2574     sched();
2575
2576     // Tidy up.
2577     proc->chan = 0;
2578
2579     // Reacquire original lock.
2580     if(lk != &ptable.lock){
2581         release(&ptable.lock);
2582         acquire(lk);
2583     }
2584
2585 }
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599

```

```

2600 // Wake up all processes sleeping on chan.
2601 // The ptable lock must be held.
2602 static void
2603 wakeup1(void *chan)
2604 {
2605     struct proc *p;
2606
2607     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
2608         if(p->state == SLEEPING && p->chan == chan)
2609             p->state = RUNNABLE;
2610 }
2611
2612 // Wake up all processes sleeping on chan.
2613 void
2614 wakeup(void *chan)
2615 {
2616     acquire(&ptable.lock);
2617     wakeup1(chan);
2618     release(&ptable.lock);
2619 }
2620
2621 // Kill the process with the given pid.
2622 // Process won't exit until it returns
2623 // to user space (see trap in trap.c).
2624 int
2625 kill(int pid)
2626 {
2627     struct proc *p;
2628
2629     acquire(&ptable.lock);
2630     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++) {
2631         if(p->pid == pid) {
2632             p->kill = 1;
2633             // Wake process from sleep if necessary.
2634             if(p->state == SLEEPING)
2635                 p->state = RUNNABLE;
2636             release(&ptable.lock);
2637             return 0;
2638         }
2639     }
2640     release(&ptable.lock);
2641     return -1;
2642 }
2643
2644
2645
2646
2647
2648
2649

```

```

2650 // Print a process listing to console. For debugging.
2651 // Runs when user types ^P on console.
2652 // No lock to avoid wedging a stuck machine further.
2653 void
2654 procdump(void)
2655 {
2656     static char *states[] = {
2657         [UNUSED]    "unused",
2658         [EMBRYO]    "embryo",
2659         [SLEEPING]  "sleep ",
2660         [RUNNABLE]  "runble",
2661         [RUNNING]   "run   ",
2662         [ZOMBIE]    "zombie",
2663     };
2664     int i;
2665     struct proc *p;
2666     char *state;
2667     uint pc[10];
2668
2669     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++) {
2670         if(p->state == UNUSED)
2671             continue;
2672         if(p->state >= 0 && p->state < NELEM(states) && states[p->state])
2673             state = states[p->state];
2674         else
2675             state = "???";
2676         printf("%d %s", p->pid, state, p->name);
2677         if(p->state == SLEEPING) {
2678             getcallerpcs((uint*)p->context->ebp+2, pc);
2679             for(i=0; i<10 && pc[i] != 0; i++)
2680                 printf(" %p", pc[i]);
2681         }
2682         printf("\n");
2683     }
2684 }
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

```

```

2700 # Context switch
2701 #
2702 # void switch(struct context **old, struct context *new);
2703 #
2704 # Save current register context in old
2705 # and then load register context from new.
2706
2707 .globl switch
2708 switch:
2709     movl 4(%esp), %eax
2710     movl 8(%esp), %edx
2711
2712 # Save old callee-save registers
2713     pushl %ebp
2714     pushl %ebx
2715     pushl %esi
2716     pushl %edi
2717
2718 # Switch stacks
2719     movl %esp, (%eax)
2720     movl %edx, %esp
2721
2722 # Load new callee-save registers
2723     popl %edi
2724     popl %esi
2725     popl %ebx
2726     popl %ebp
2727     ret
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749

```

```

2750 // Physical memory allocator, intended to allocate
2751 // memory for user processes, kernel stacks, page table pages,
2752 // and pipe buffers. Allocates 4096-byte pages.
2753
2754 #include "types.h"
2755 #include "defs.h"
2756 #include "param.h"
2757 #include "memlayout.h"
2758 #include "mmu.h"
2759 #include "spinlock.h"
2760
2761 void freerange(void *vstart, void *vend);
2762 extern char end[]; // first address after kernel loaded from ELF file
2763
2764 struct run {
2765     struct run *next;
2766 };
2767
2768 struct {
2769     struct spinlock lock;
2770     int use_lock;
2771     struct run *freelist;
2772 } kmem;
2773
2774 // Initialization happens in two phases.
2775 // 1. main() calls kinit1() while still using entrypgdir to place just
2776 // the pages mapped by entrypgdir on free list.
2777 // 2. main() calls kinit2() with the rest of the physical pages
2778 // after installing a full page table that maps them on all cores.
2779 void
2780 kinit1(void *vstart, void *vend)
2781 {
2782     initlock(&kmem.lock, "kmem");
2783     kmem.use_lock = 0;
2784     freerange(vstart, vend);
2785 }
2786
2787 void
2788 kinit2(void *vstart, void *vend)
2789 {
2790     freerange(vstart, vend);
2791     kmem.use_lock = 1;
2792 }
2793
2794
2795
2796
2797
2798
2799

```

```
2800 void
2801 freerange(void *vstart, void *vend)
2802 {
2803     char *p;
2804     p = (char*)PGROUNDUP((uint)vstart);
2805     for(; p + PGSIZE <= (char*)vend; p += PGSIZE)
2806         kfree(p);
2807 }
2808
2809
2810 // Free the page of physical memory pointed at by v,
2811 // which normally should have been returned by a
2812 // call to kalloc(). (The exception is when
2813 // initializing the allocator; see kinit above.)
2814 void
2815 kfree(char *v)
2816 {
2817     struct run *r;
2818
2819     if((uint)v % PGSIZE || v < end || v2p(v) >= PHYSTOP)
2820         panic("kfree");
2821
2822     // Fill with junk to catch dangling refs.
2823     memset(v, 1, PGSIZE);
2824
2825     if(kmem.use_lock)
2826         acquire(&kmem.lock);
2827     r = (struct run*)v;
2828     r->next = kmem.freelist;
2829     kmem.freelist = r;
2830     if(kmem.use_lock)
2831         release(&kmem.lock);
2832 }
2833
2834 // Allocate one 4096-byte page of physical memory.
2835 // Returns a pointer that the kernel can use.
2836 // Returns 0 if the memory cannot be allocated.
2837 char*
2838 kalloc(void)
2839 {
2840     struct run *r;
2841
2842     if(kmem.use_lock)
2843         acquire(&kmem.lock);
2844     r = kmem.freelist;
2845     if(r)
2846         kmem.freelist = r->next;
2847     if(kmem.use_lock)
2848         release(&kmem.lock);
2849     return (char*)r;
```

```
2850 }
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
```



```

2900 // x86 trap and interrupt constants.
2901
2902 // Processor-defined:
2903 #define T_DIVIDE 0 // divide error
2904 #define T_DEBUG 1 // debug exception
2905 #define T_NMI 2 // non-maskable interrupt
2906 #define T_BRKPT 3 // breakpoint
2907 #define T_OFLOW 4 // overflow
2908 #define T_BOUND 5 // bounds check
2909 #define T_ILLOP 6 // illegal opcode
2910 #define T_DEVICE 7 // device not available
2911 #define T_DBLFLT 8 // double fault
2912 // #define T_COPROC 9 // reserved (not used since 486)
2913 #define T_TSS 10 // invalid task switch segment
2914 #define T_SEGNP 11 // segment not present
2915 #define T_STACK 12 // stack exception
2916 #define T_GPFLT 13 // general protection fault
2917 #define T_PGFLT 14 // page fault
2918 // #define T_RES 15 // reserved
2919 #define T_FPERR 16 // floating point error
2920 #define T_ALIGN 17 // alignment check
2921 #define T_MCHK 18 // machine check
2922 #define T_SIMDERR 19 // SIMD floating point error
2923
2924 // These are arbitrarily chosen, but with care not to overlap
2925 // processor defined exceptions or interrupt vectors.
2926 #define T_SYSCALL 64 // system call
2927 #define T_DEFAULT 500 // catchall
2928
2929 #define T_IRQ0 32 // IRQ 0 corresponds to int T_IRQ
2930
2931 #define IRQ_TIMER 0
2932 #define IRQ_KBD 1
2933 #define IRQ_COM1 4
2934 #define IRQ_IDE 14
2935 #define IRQ_ERROR 19
2936 #define IRQ_SPURIOUS 31
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949

```

```

2950 #!/usr/bin/perl -w
2951
2952 # Generate vectors.S, the trap/interrupt entry points.
2953 # There has to be one entry point per interrupt number
2954 # since otherwise there's no way for trap() to discover
2955 # the interrupt number.
2956
2957 print "# generated by vectors.pl - do not edit\n";
2958 print "# handlers\n";
2959 print ".globl alltraps\n";
2960 for(my $i = 0; $i < 256; $i++){
2961     print ".globl vector${i}\n";
2962     print "vector${i}:\n";
2963     if(!($i == 8 || ($i >= 10 && $i <= 14) || $i == 17)){
2964         print "    pushl \${0}\n";
2965     }
2966     print "    pushl \${$i}\n";
2967     print "    jmp alltraps\n";
2968 }
2969
2970 print "\n# vector table\n";
2971 print ".data\n";
2972 print ".globl vectors\n";
2973 print "vectors:\n";
2974 for(my $i = 0; $i < 256; $i++){
2975     print "    .long vector${i}\n";
2976 }
2977
2978 # sample output:
2979 # # handlers
2980 # .globl alltraps
2981 # .globl vector0
2982 # vector0:
2983 #     pushl $0
2984 #     pushl $0
2985 #     jmp alltraps
2986 # ...
2987 #
2988 # # vector table
2989 # .data
2990 # .globl vectors
2991 # vectors:
2992 #     .long vector0
2993 #     .long vector1
2994 #     .long vector2
2995 #     ...
2996
2997
2998
2999

```

```

3000 #include "mmu.h"
3001
3002 # vectors.S sends all traps here.
3003 .globl alltraps
3004 alltraps:
3005 # Build trap frame.
3006 pushl %ds
3007 pushl %es
3008 pushl %fs
3009 pushl %gs
3010 pushal
3011
3012 # Set up data and per-cpu segments.
3013 movw $(SEG_KDATA<<3), %ax
3014 movw %ax, %ds
3015 movw %ax, %es
3016 movw $(SEG_KCPU<<3), %ax
3017 movw %ax, %fs
3018 movw %ax, %gs
3019
3020 # Call trap(tf), where tf=%esp
3021 pushl %esp
3022 call trap
3023 addl $4, %esp
3024
3025 # Return falls through to trapret...
3026 .globl trapret
3027 trapret:
3028 popal
3029 popl %gs
3030 popl %fs
3031 popl %es
3032 popl %ds
3033 addl $0x8, %esp # trapno and errcode
3034 iret
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049

```

```

3050 #include "types.h"
3051 #include "defs.h"
3052 #include "param.h"
3053 #include "memlayout.h"
3054 #include "mmu.h"
3055 #include "proc.h"
3056 #include "x86.h"
3057 #include "traps.h"
3058 #include "spinlock.h"
3059
3060 // Interrupt descriptor table (shared by all CPUs).
3061 struct gatedesc idt[256];
3062 extern uint vectors[]; // in vectors.S: array of 256 entry pointers
3063 struct spinlock tickslock;
3064 uint ticks;
3065
3066 void
3067 tvinit(void)
3068 {
3069     int i;
3070
3071     for(i = 0; i < 256; i++)
3072         SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);
3073     SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL], DPL_USER);
3074
3075     initlock(&tickslock, "time");
3076 }
3077
3078 void
3079 idtinit(void)
3080 {
3081     lidt(idt, sizeof(idt));
3082 }
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099

```

```

3100 void
3101 trap(struct trapframe *tf)
3102 {
3103     if(tf->trapno == T_SYSCALL){
3104         if(proc->killed)
3105             exit();
3106         proc->tf = tf;
3107         syscall();
3108         if(proc->killed)
3109             exit();
3110         return;
3111     }
3112     switch(tf->trapno){
3113     case T_IRQ0 + IRQ_TIMER:
3114         if(cpu->id == 0){
3115             acquire(&tickslock);
3116             ticks++;
3117             wakeup(&ticks);
3118             release(&tickslock);
3119         }
3120         lapiceoi();
3121         break;
3122     case T_IRQ0 + IRQ_IDE:
3123         ideintr();
3124         lapiceoi();
3125         break;
3126     case T_IRQ0 + IRQ_IDE+1:
3127         // Bochs generates spurious IDE1 interrupts.
3128         break;
3129     case T_IRQ0 + IRQ_KBD:
3130         kbdrintr();
3131         lapiceoi();
3132         break;
3133     case T_IRQ0 + IRQ_COM1:
3134         uartintr();
3135         lapiceoi();
3136         break;
3137     case T_IRQ0 + 7:
3138     case T_IRQ0 + IRQ_SPURIOUS:
3139         cprintf("cpu%d: spurious interrupt at %x:%x\n",
3140             cpu->id, tf->cs, tf->eip);
3141         lapiceoi();
3142         break;
3143     }
3144 }
3145
3146
3147
3148
3149

```

```

3150 default:
3151     if(proc == 0 || (tf->cs&3) == 0){
3152         // In kernel, it must be our mistake.
3153         cprintf("unexpected trap %d from cpu %d eip %x (cr2=0x%x)\n",
3154             tf->trapno, cpu->id, tf->eip, rcr2());
3155         panic("trap");
3156     }
3157     // In user space, assume process misbehaved.
3158     cprintf("pid %d %s: trap %d err %d on cpu %d "
3159         "eip 0x%x addr 0x%x--kill proc\n",
3160         proc->pid, proc->name, tf->trapno, tf->err, cpu->id, tf->eip,
3161         rcr2());
3162     proc->killed = 1;
3163 }
3164
3165 // Force process exit if it has been killed and is in user space.
3166 // (If it is still executing in the kernel, let it keep running
3167 // until it gets to the regular system call return.)
3168 if(proc && proc->killed && (tf->cs&3) == DPL_USER)
3169     exit();
3170
3171 // Force process to give up CPU on clock tick.
3172 // If interrupts were on while locks held, would need to check nlock.
3173 if(proc && proc->state == RUNNING && tf->trapno == T_IRQ0+IRQ_TIMER)
3174     yield();
3175
3176 // Check if the process has been killed since we yielded
3177 if(proc && proc->killed && (tf->cs&3) == DPL_USER)
3178     exit();
3179 }
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199

```

```

3200 // System call numbers
3201 #define SYS_fork 1
3202 #define SYS_exit 2
3203 #define SYS_wait 3
3204 #define SYS_pipe 4
3205 #define SYS_read 5
3206 #define SYS_kill 6
3207 #define SYS_exec 7
3208 #define SYS_fstat 8
3209 #define SYS_chdir 9
3210 #define SYS_dup 10
3211 #define SYS_getpid 11
3212 #define SYS_sbrk 12
3213 #define SYS_sleep 13
3214 #define SYS_uptime 14
3215 #define SYS_open 15
3216 #define SYS_write 16
3217 #define SYS_mknod 17
3218 #define SYS_unlink 18
3219 #define SYS_link 19
3220 #define SYS_mkdir 20
3221 #define SYS_close 21
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249

```

```

3250 #include "types.h"
3251 #include "defs.h"
3252 #include "param.h"
3253 #include "memlayout.h"
3254 #include "mmu.h"
3255 #include "proc.h"
3256 #include "x86.h"
3257 #include "syscall.h"
3258
3259 // User code makes a system call with INT T_SYSCALL.
3260 // System call number in %eax.
3261 // Arguments on the stack, from the user call to the C
3262 // library system call function. The saved user %esp points
3263 // to a saved program counter, and then the first argument.
3264
3265 // Fetch the int at addr from the current process.
3266 int
3267 fetchint(uint addr, int *ip)
3268 {
3269     if(addr >= proc->sz || addr+4 > proc->sz)
3270         return -1;
3271     *ip = *(int*)(addr);
3272     return 0;
3273 }
3274
3275 // Fetch the nul-terminated string at addr from the current process.
3276 // Doesn't actually copy the string - just sets *pp to point at it.
3277 // Returns length of string, not including nul.
3278 int
3279 fetchstr(uint addr, char **pp)
3280 {
3281     char *s, *ep;
3282
3283     if(addr >= proc->sz)
3284         return -1;
3285     *pp = (char*)addr;
3286     ep = (char*)proc->sz;
3287     for(s = *pp; s < ep; s++)
3288         if(*s == 0)
3289             return s - *pp;
3290     return -1;
3291 }
3292
3293 // Fetch the nth 32-bit system call argument.
3294 int
3295 argint(int n, int *ip)
3296 {
3297     return fetchint(proc->tf->esp + 4 + 4*n, ip);
3298 }
3299

```

```

3300 // Fetch the nth word-sized system call argument as a pointer
3301 // to a block of memory of size n bytes. Check that the pointer
3302 // lies within the process address space.
3303 int
3304 argptr(int n, char **pp, int size)
3305 {
3306     int i;
3307
3308     if(argint(n, &i) < 0)
3309         return -1;
3310     if((uint)i >= proc->sz || (uint)i+size > proc->sz)
3311         return -1;
3312     *pp = (char*)i;
3313     return 0;
3314 }
3315
3316 // Fetch the nth word-sized system call argument as a string pointer.
3317 // Check that the pointer is valid and the string is nul-terminated.
3318 // (There is no shared writable memory, so the string can't change
3319 // between this check and being used by the kernel.)
3320 int
3321 argstr(int n, char **pp)
3322 {
3323     int addr;
3324     if(argint(n, &addr) < 0)
3325         return -1;
3326     return fetchstr(addr, pp);
3327 }
3328
3329 extern int sys_chdir(void);
3330 extern int sys_close(void);
3331 extern int sys_dup(void);
3332 extern int sys_exec(void);
3333 extern int sys_exit(void);
3334 extern int sys_fork(void);
3335 extern int sys_getpid(void);
3336 extern int sys_kill(void);
3337 extern int sys_link(void);
3338 extern int sys_link(void);
3339 extern int sys_mkdir(void);
3340 extern int sys_mknod(void);
3341 extern int sys_open(void);
3342 extern int sys_pipe(void);
3343 extern int sys_read(void);
3344 extern int sys_sbrk(void);
3345 extern int sys_sleep(void);
3346 extern int sys_unlink(void);
3347 extern int sys_wait(void);
3348 extern int sys_write(void);
3349 extern int sys_uptime(void);

```

```

3350 static int (*syscalls[])(void) = {
3351     [SYS_fork]   sys_fork,
3352     [SYS_exit]   sys_exit,
3353     [SYS_wait]   sys_wait,
3354     [SYS_pipe]   sys_pipe,
3355     [SYS_read]   sys_read,
3356     [SYS_kill]   sys_kill,
3357     [SYS_exec]   sys_exec,
3358     [SYS_fstat]  sys_fstat,
3359     [SYS_chdir]  sys_chdir,
3360     [SYS_dup]    sys_dup,
3361     [SYS_getpid] sys_getpid,
3362     [SYS_sbrk]   sys_sbrk,
3363     [SYS_sleep]  sys_sleep,
3364     [SYS_uptime] sys_uptime,
3365     [SYS_open]   sys_open,
3366     [SYS_write]  sys_write,
3367     [SYS_mknod]  sys_mknod,
3368     [SYS_unlink] sys_unlink,
3369     [SYS_link]   sys_link,
3370     [SYS_mkdir]  sys_mkdir,
3371     [SYS_close]  sys_close,
3372 };
3373
3374 void
3375 syscall(void)
3376 {
3377     int num;
3378
3379     num = proc->tf->eax;
3380     if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
3381         proc->tf->eax = syscalls[num]();
3382     } else {
3383         cprintf("%d %s: unknown sys call %d\n",
3384                 proc->pid, proc->name, num);
3385         proc->tf->eax = -1;
3386     }
3387 }
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399

```

```

3400 #include "types.h"
3401 #include "x86.h"
3402 #include "defs.h"
3403 #include "param.h"
3404 #include "memlayout.h"
3405 #include "mmu.h"
3406 #include "proc.h"
3407
3408 int
3409 sys_fork(void)
3410 {
3411   return fork();
3412 }
3413
3414 int
3415 sys_exit(void)
3416 {
3417   exit();
3418   return 0; // not reached
3419 }
3420
3421 int
3422 sys_wait(void)
3423 {
3424   return wait();
3425 }
3426
3427 int
3428 sys_kill(void)
3429 {
3430   int pid;
3431
3432   if(argint(0, &pid) < 0)
3433     return -1;
3434   return kill(pid);
3435 }
3436
3437 int
3438 sys_getpid(void)
3439 {
3440   return proc->pid;
3441 }
3442
3443
3444
3445
3446
3447
3448
3449

```

```

3450 int
3451 sys_sbrk(void)
3452 {
3453   int addr;
3454   int n;
3455
3456   if(argint(0, &n) < 0)
3457     return -1;
3458   addr = proc->sz;
3459   if(growproc(n) < 0)
3460     return -1;
3461   return addr;
3462 }
3463
3464 int
3465 sys_sleep(void)
3466 {
3467   int n;
3468   uint ticks0;
3469
3470   if(argint(0, &n) < 0)
3471     return -1;
3472   acquire(&tickslock);
3473   ticks0 = ticks;
3474   while(ticks - ticks0 < n){
3475     if(proc->killed){
3476       release(&tickslock);
3477       return -1;
3478     }
3479     sleep(&ticks, &tickslock);
3480   }
3481   release(&tickslock);
3482   return 0;
3483 }
3484
3485 // return how many clock tick interrupts have occurred
3486 // since start.
3487 int
3488 sys_uptime(void)
3489 {
3490   uint xticks;
3491
3492   acquire(&tickslock);
3493   xticks = ticks;
3494   release(&tickslock);
3495   return xticks;
3496 }
3497
3498
3499

```



```
3500 struct buf {
3501     int flags;
3502     uint dev;
3503     uint sector;
3504     struct buf *prev; // LRU cache list
3505     struct buf *next;
3506     struct buf *qnext; // disk queue
3507     uchar data[512];
3508 };
3509 #define B_BUSY 0x1 // buffer is locked by some process
3510 #define B_VALID 0x2 // buffer has been read from disk
3511 #define B_DIRTY 0x4 // buffer needs to be written to disk
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
```

```
3550 #define O_RDONLY 0x000
3551 #define O_WRONLY 0x001
3552 #define O_RDWR 0x002
3553 #define O_CREATE 0x200
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563
3564
3565
3566
3567
3568
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599
```

```

3600 #define T_DIR 1 // Directory
3601 #define T_FILE 2 // File
3602 #define T_DEV 3 // Device
3603
3604 struct stat {
3605     short type; // Type of file
3606     int dev; // File system's disk device
3607     uint ino; // Inode number
3608     short nlink; // Number of links to file
3609     uint size; // Size of file in bytes
3610 };
3611
3612
3613
3614
3615
3616
3617
3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649

```

```

3650 // On-disk file system format.
3651 // Both the kernel and user programs use this header file.
3652
3653 // Block 0 is unused.
3654 // Block 1 is super block.
3655 // Blocks 2 through sb.ninodes/IPB hold inodes.
3656 // Then free bitmap blocks holding sb.size bits.
3657 // Then sb.nblocks data blocks.
3658 // Then sb.nlog log blocks.
3659
3660 #define ROOTINO 1 // root i-number
3661 #define BSIZE 512 // block size
3662
3663 // File system super block
3664 struct superblock {
3665     uint size; // Size of file system image (blocks)
3666     uint nblocks; // Number of data blocks
3667     uint ninodes; // Number of inodes.
3668     uint nlog; // Number of log blocks
3669 };
3670
3671 #define NDIRECT 12
3672 #define NINDIRECT (BSIZE / sizeof(uint))
3673 #define MAXFILE (NDIRECT + NINDIRECT)
3674
3675 // On-disk inode structure
3676 struct dinode {
3677     short type; // File type
3678     short major; // Major device number (T_DEV only)
3679     short minor; // Minor device number (T_DEV only)
3680     short nlink; // Number of links to inode in file system
3681     uint size; // Size of file (bytes)
3682     uint addrs[NDIRECT+1]; // Data block addresses
3683 };
3684
3685 // Inodes per block.
3686 #define IPB (BSIZE / sizeof(struct dinode))
3687
3688 // Block containing inode i
3689 #define IBLOCK(i) ((i) / IPB + 2)
3690
3691 // Bitmap bits per block
3692 #define BPB (BSIZE*8)
3693
3694 // Block containing bit for block b
3695 #define BBLOCK(b, ninodes) (b/BPB + (ninodes)/IPB + 3)
3696
3697 // Directory is a file containing a sequence of dirent structures.
3698 #define DIRSIZ 14
3699

```

```
3700 struct dirent {
3701     ushort inum;
3702     char name[DIRSIZ];
3703 };
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725
3726
3727
3728
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
```

```
3750 struct file {
3751     enum { FD_NONE, FD_PIPE, FD_INODE } type;
3752     int ref; // reference count
3753     char readable;
3754     char writable;
3755     struct pipe *pipe;
3756     struct inode *ip;
3757     uint off;
3758 };
3759
3760
3761 // in-memory copy of an inode
3762 struct inode {
3763     uint dev;           // Device number
3764     uint inum;          // Inode number
3765     int ref;            // Reference count
3766     int flags;          // I_BUSY, I_VALID
3767
3768     short type;         // copy of disk inode
3769     short major;
3770     short minor;
3771     short nlink;
3772     uint size;
3773     uint addrs[NDIRECT+1];
3774 };
3775 #define I_BUSY 0x1
3776 #define I_VALID 0x2
3777
3778 // table mapping major device number to
3779 // device functions
3780 struct devsw {
3781     int (*read)(struct inode*, char*, int);
3782     int (*write)(struct inode*, char*, int);
3783 };
3784
3785 extern struct devsw devsw[];
3786
3787 #define CONSOLE 1
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799
```

```

3800 // Simple PIO-based (non-DMA) IDE driver code.
3801
3802 #include "types.h"
3803 #include "defs.h"
3804 #include "param.h"
3805 #include "memlayout.h"
3806 #include "mmu.h"
3807 #include "proc.h"
3808 #include "x86.h"
3809 #include "traps.h"
3810 #include "spinlock.h"
3811 #include "buf.h"
3812
3813 #define IDE_BSY 0x80
3814 #define IDE_DRDY 0x40
3815 #define IDE_DF 0x20
3816 #define IDE_ERR 0x01
3817
3818 #define IDE_CMD_READ 0x20
3819 #define IDE_CMD_WRITE 0x30
3820
3821 // idequeue points to the buf now being read/written to the disk.
3822 // idequeue->qnext points to the next buf to be processed.
3823 // You must hold ideLock while manipulating queue.
3824
3825 static struct spinlock ideLock;
3826 static struct buf *idequeue;
3827
3828 static int havedisk1;
3829 static void idestart(struct buf*);
3830
3831 // Wait for IDE disk to become ready.
3832 static int
3833 idewait(int checkerr)
3834 {
3835     int r;
3836
3837     while(((r = inb(0x1f7)) & (IDE_BSY|IDE_DRDY)) != IDE_DRDY)
3838         ;
3839     if(checkerr && (r & (IDE_DF|IDE_ERR)) != 0)
3840         return -1;
3841     return 0;
3842 }
3843
3844
3845
3846
3847
3848
3849

```

```

3850 void
3851 ideinit(void)
3852 {
3853     int i;
3854
3855     initlock(&ideLock, "ide");
3856     plicenable(IRQ_IDE);
3857     ioapicenable(IRQ_IDE, ncpu - 1);
3858     idewait(0);
3859
3860     // Check if disk 1 is present
3861     outb(0x1f6, 0xe0 | (1<<4));
3862     for(i=0; i<1000; i++){
3863         if(inb(0x1f7) != 0){
3864             havedisk1 = 1;
3865             break;
3866         }
3867     }
3868
3869     // Switch back to disk 0.
3870     outb(0x1f6, 0xe0 | (0<<4));
3871 }
3872
3873 // Start the request for b. Caller must hold ideLock.
3874 static void
3875 idestart(struct buf *b)
3876 {
3877     if(b == 0)
3878         panic("idestart");
3879
3880     idewait(0);
3881     outb(0x3f6, 0); // generate interrupt
3882     outb(0x1f2, 1); // number of sectors
3883     outb(0x1f3, b->sector & 0xff);
3884     outb(0x1f4, (b->sector >> 8) & 0xff);
3885     outb(0x1f5, (b->sector >> 16) & 0xff);
3886     outb(0x1f6, 0xe0 | ((b->dev&1)<<4) | ((b->sector>>24)&0x0f));
3887     if(b->flags & B_DIRTY){
3888         outb(0x1f7, IDE_CMD_WRITE);
3889         outsl(0x1f0, b->data, 512/4);
3890     } else {
3891         outb(0x1f7, IDE_CMD_READ);
3892     }
3893 }
3894
3895
3896
3897
3898
3899

```

```

3900 // Interrupt handler.
3901 void
3902 ideintr(void)
3903 {
3904     struct buf *b;
3905
3906     // First queued buffer is the active request.
3907     acquire(&idelock);
3908     if((b = idequeue) == 0){
3909         release(&idelock);
3910         // printf("spurious IDE interrupt\n");
3911         return;
3912     }
3913     idequeue = b->qnext;
3914
3915     // Read data if needed.
3916     if(! (b->flags & B_DIRTY) && idewait(1) >= 0)
3917         insl(0x1f0, b->data, 512/4);
3918
3919     // Wake process waiting for this buf.
3920     b->flags |= B_VALID;
3921     b->flags &= ~B_DIRTY;
3922     wakeup(b);
3923
3924     // Start disk on next buf in queue.
3925     if(idequeue != 0)
3926         idestart(idequeue);
3927
3928     release(&idelock);
3929 }
3930
3931
3932
3933
3934
3935
3936
3937
3938
3939
3940
3941
3942
3943
3944
3945
3946
3947
3948
3949

```

```

3950 // Sync buf with disk.
3951 // If B_DIRTY is set, write buf to disk, clear B_DIRTY, set B_VALID.
3952 // Else if B_VALID is not set, read buf from disk, set B_VALID.
3953 void
3954 iderw(struct buf *b)
3955 {
3956     struct buf **pp;
3957
3958     if(! (b->flags & B_BUSY))
3959         panic("iderw: buf not busy");
3960     if((b->flags & (B_VALID|B_DIRTY)) == B_VALID)
3961         panic("iderw: nothing to do");
3962     if(b->dev != 0 && !havedisk1)
3963         panic("iderw: ide disk 1 not present");
3964
3965     acquire(&idelock);
3966
3967     // Append b to idequeue.
3968     b->qnext = 0;
3969     for(pp=&idequeue; *pp; pp=&(*pp)->qnext)
3970         ;
3971     *pp = b;
3972
3973     // Start disk if necessary.
3974     if(idequeue == b)
3975         idestart(b);
3976
3977     // Wait for request to finish.
3978     while((b->flags & (B_VALID|B_DIRTY)) != B_VALID){
3979         sleep(b, &idelock);
3980     }
3981
3982     release(&idelock);
3983 }
3984
3985
3986
3987
3988
3989
3990
3991
3992
3993
3994
3995
3996
3997
3998
3999

```

```

4000 // Buffer cache.
4001 //
4002 // The buffer cache is a linked list of buf structures holding
4003 // cached copies of disk block contents. Caching disk blocks
4004 // in memory reduces the number of disk reads and also provides
4005 // a synchronization point for disk blocks used by multiple processes.
4006 //
4007 // Interface:
4008 // * To get a buffer for a particular disk block, call bread.
4009 // * After changing buffer data, call bwrite to write it to disk.
4010 // * When done with the buffer, call brelse.
4011 // * Do not use the buffer after calling brelse.
4012 // * Only one process at a time can use a buffer,
4013 //   so do not keep them longer than necessary.
4014 //
4015 // The implementation uses three state flags internally:
4016 // * B_BUSY: the block has been returned from bread
4017 //   and has not been passed back to brelse.
4018 // * B_VALID: the buffer data has been read from the disk.
4019 // * B_DIRTY: the buffer data has been modified
4020 //   and needs to be written to disk.
4021
4022 #include "types.h"
4023 #include "defs.h"
4024 #include "param.h"
4025 #include "spinlock.h"
4026 #include "buf.h"
4027
4028 struct {
4029   struct spinlock lock;
4030   struct buf buf[NBUF];
4031 }
4032 // Linked list of all buffers, through prev/next.
4033 // head.next is most recently used.
4034 struct buf head;
4035 } bcache;
4036
4037 void
4038 binit(void)
4039 {
4040   struct buf *b;
4041
4042   initlock(&bcache.lock, "bcache");
4043
4044
4045
4046
4047
4048
4049

```

```

4050 // Create linked list of buffers
4051 bcache.head.prev = &bcache.head;
4052 bcache.head.next = &bcache.head;
4053 for(b = bcache.buf; b < bcache.buf+NBUF; b++){
4054   b->next = bcache.head.next;
4055   b->prev = &bcache.head;
4056   b->dev = -1;
4057   bcache.head.next->prev = b;
4058   bcache.head.next = b;
4059 }
4060 }
4061
4062 // Look through buffer cache for sector on device dev.
4063 // If not found, allocate fresh block.
4064 // In either case, return B_BUSY buffer.
4065 static struct buf*
4066 bget(uint dev, uint sector)
4067 {
4068   struct buf *b;
4069
4070   acquire(&bcache.lock);
4071
4072   loop:
4073   // Is the sector already cached?
4074   for(b = bcache.head.next; b != &bcache.head; b = b->next){
4075     if(b->dev == dev && b->sector == sector){
4076       if(!(b->flags & B_BUSY)){
4077         b->flags |= B_BUSY;
4078         release(&bcache.lock);
4079         return b;
4080       }
4081       sleep(b, &bcache.lock);
4082       goto loop;
4083     }
4084   }
4085
4086   // Not cached; recycle some non-busy and clean buffer.
4087   for(b = bcache.head.prev; b != &bcache.head; b = b->prev){
4088     if((b->flags & B_BUSY) == 0 && (b->flags & B_DIRTY) == 0){
4089       b->dev = dev;
4090       b->sector = sector;
4091       b->flags = B_BUSY;
4092       release(&bcache.lock);
4093       return b;
4094     }
4095   }
4096   panic("bget: no buffers");
4097 }
4098
4099

```



```

4100 // Return a B_BUSY buf with the contents of the indicated disk sector.
4101 struct buf*
4102 bread(uint dev, uint sector)
4103 {
4104     struct buf *b;
4105
4106     b = bget(dev, sector);
4107     if(!b->flags & B_VALID)
4108         iderw(b);
4109     return b;
4110 }
4111
4112 // Write b's contents to disk. Must be B_BUSY.
4113 void
4114 bwrite(struct buf *b)
4115 {
4116     if((b->flags & B_BUSY) == 0)
4117         panic("bwrite");
4118     b->flags |= B_DIRTY;
4119     iderw(b);
4120 }
4121
4122 // Release a B_BUSY buffer.
4123 // Move to the head of the MRU list.
4124 void
4125 brelse(struct buf *b)
4126 {
4127     if((b->flags & B_BUSY) == 0)
4128         panic("brelse");
4129
4130     acquire(&bcache.lock);
4131
4132     b->next->prev = b->prev;
4133     b->prev->next = b->next;
4134     b->next = bcache.head.next;
4135     b->prev = &bcache.head;
4136     bcache.head.next->prev = b;
4137     bcache.head.next = b;
4138
4139     b->flags &= ~B_BUSY;
4140     wakeup(b);
4141
4142     release(&bcache.lock);
4143 }
4144
4145
4146
4147
4148
4149

```

```

4150 #include "types.h"
4151 #include "defs.h"
4152 #include "param.h"
4153 #include "spinlock.h"
4154 #include "fs.h"
4155 #include "buf.h"
4156
4157 // Simple logging. Each system call that might write the file system
4158 // should be surrounded with begin_trans() and commit_trans() calls.
4159 //
4160 // The log holds at most one transaction at a time. Commit forces
4161 // the log (with commit record) to disk, then installs the affected
4162 // blocks to disk, then erases the log. begin_trans() ensures that
4163 // only one system call can be in a transaction; others must wait.
4164 //
4165 // Allowing only one transaction at a time means that the file
4166 // system code doesn't have to worry about the possibility of
4167 // one transaction reading a block that another one has modified,
4168 // for example an i-node block.
4169 //
4170 // Read-only system calls don't need to use transactions, though
4171 // this means that they may observe uncommitted data. I-node and
4172 // buffer locks prevent read-only calls from seeing inconsistent data.
4173 //
4174 // The log is a physical re-do log containing disk blocks.
4175 // The on-disk log format:
4176 //   header block, containing sector #s for block A, B, C, ...
4177 //   block A
4178 //   block B
4179 //   block C
4180 //   ...
4181 // Log appends are synchronous.
4182
4183 // Contents of the header block, used for both the on-disk header block
4184 // and to keep track in memory of logged sector #s before commit.
4185 struct logheader {
4186     int n;
4187     int sector[LOGSIZE];
4188 };
4189
4190 struct log {
4191     struct spinlock lock;
4192     int start;
4193     int size;
4194     int busy; // a transaction is active
4195     int dev;
4196     struct logheader lh;
4197 };
4198
4199

```

```

4200 struct log log;
4201
4202 static void recover_from_log(void);
4203
4204 void
4205 init_log(void)
4206 {
4207     if (sizeof(struct logheader) >= BSIZE)
4208         panic("initlog: too big logheader");
4209
4210     struct superblock sb;
4211     initlock(&log.lock, "log");
4212     readsb(ROOTDEV, &sb);
4213     log.start = sb.size - sb.nlog;
4214     log.size = sb.nlog;
4215     log.dev = ROOTDEV;
4216     recover_from_log();
4217 }
4218
4219 // Copy committed blocks from log to their home location
4220 static void
4221 install_trans(void)
4222 {
4223     int tail;
4224
4225     for (tail = 0; tail < log.lh.n; tail++) {
4226         struct buf *lbuf = bread(log.dev, log.start+tail+1); // read log block
4227         struct buf *dbuf = bread(log.dev, log.lh.sector[tail]); // read dst
4228         memmove(dbuf->data, lbuf->data, BSIZE); // copy block to dst
4229         bwrite(dbuf); // write dst to disk
4230         brelse(lbuf);
4231         brelse(dbuf);
4232     }
4233 }
4234
4235 // Read the log header from disk into the in-memory log header
4236 static void
4237 read_head(void)
4238 {
4239     struct buf *buf = bread(log.dev, log.start);
4240     struct logheader *lh = (struct logheader *) (buf->data);
4241     int i;
4242     log.lh.n = lh->n;
4243     for (i = 0; i < log.lh.n; i++) {
4244         log.lh.sector[i] = lh->sector[i];
4245     }
4246     brelse(buf);
4247 }
4248
4249

```

```

4250 // Write in-memory log header to disk.
4251 // This is the true point at which the
4252 // current transaction commits.
4253 static void
4254 write_head(void)
4255 {
4256     struct buf *buf = bread(log.dev, log.start);
4257     struct logheader *hb = (struct logheader *) (buf->data);
4258     int i;
4259     hb->n = log.lh.n;
4260     for (i = 0; i < log.lh.n; i++) {
4261         hb->sector[i] = log.lh.sector[i];
4262     }
4263     bwrite(buf);
4264     brelse(buf);
4265 }
4266
4267 static void
4268 recover_from_log(void)
4269 {
4270     read_head();
4271     install_trans(); // if committed, copy from log to disk
4272     log.lh.n = 0;
4273     write_head(); // clear the log
4274 }
4275
4276 void
4277 begin_trans(void)
4278 {
4279     acquire(&log.lock);
4280     while (log.busy) {
4281         sleep(&log, &log.lock);
4282     }
4283     log.busy = 1;
4284     release(&log.lock);
4285 }
4286
4287
4288
4289
4290
4291
4292
4293
4294
4295
4296
4297
4298
4299

```

```

4300 void
4301 commit_trans(void)
4302 {
4303     if (log.lh.n > 0) {
4304         write_head(); // Write header to disk -- the real commit
4305         install_trans(); // Now install writes to home locations
4306         log.lh.n = 0;
4307         write_head(); // Erase the transaction from the log
4308     }
4309
4310     acquire(&log.lock);
4311     log.busy = 0;
4312     wakeup(&log);
4313     release(&log.lock);
4314 }
4315
4316 // Caller has modified b->data and is done with the buffer.
4317 // Append the block to the log and record the block number,
4318 // but don't write the log header (which would commit the write).
4319 // log_write() replaces bwrite(); a typical use is:
4320 // bp = bread(...);
4321 // modify bp->data[]
4322 // log_write(bp)
4323 // brelse(bp)
4324 void
4325 log_write(struct buf *b)
4326 {
4327     int i;
4328
4329     if (log.lh.n >= LOGSIZE || log.lh.n >= log.size - 1)
4330         panic("too big a transaction");
4331     if (!log.busy)
4332         panic("write outside of trans");
4333
4334     for (i = 0; i < log.lh.n; i++) {
4335         if (log.lh.sector[i] == b->sector) // log absorption?
4336             break;
4337     }
4338     log.lh.sector[i] = b->sector;
4339     struct buf *lbuf = bread(b->dev, log.start+i+1);
4340     memmove(lbuf->data, b->data, BSIZE);
4341     bwrite(lbuf);
4342     brelse(lbuf);
4343     if (i == log.lh.n)
4344         log.lh.n++;
4345     b->flags |= B_DIRTY; // XXX prevent eviction
4346 }
4347
4348
4349

```

```

4350 // Blank page.
4351
4352
4353
4354
4355
4356
4357
4358
4359
4360
4361
4362
4363
4364
4365
4366
4367
4368
4369
4370
4371
4372
4373
4374
4375
4376
4377
4378
4379
4380
4381
4382
4383
4384
4385
4386
4387
4388
4389
4390
4391
4392
4393
4394
4395
4396
4397
4398
4399

```

```

4400 // File system implementation. Five layers:
4401 // + Blocks: allocator for raw disk blocks.
4402 // + Log: crash recovery for multi-step updates.
4403 // + Files: inode allocator, reading, writing, metadata.
4404 // + Directories: inode with special contents (list of other inodes!)
4405 // + Names: paths like /usr/rm/xv6/fs.c for convenient naming.
4406 //
4407 // This file contains the low-level file system manipulation
4408 // routines. The (higher-level) system call implementations
4409 // are in sysfile.c.
4410
4411 #include "types.h"
4412 #include "defs.h"
4413 #include "param.h"
4414 #include "stat.h"
4415 #include "mmu.h"
4416 #include "proc.h"
4417 #include "spinlock.h"
4418 #include "buf.h"
4419 #include "fs.h"
4420 #include "file.h"
4421
4422 #define min(a, b) ((a) < (b) ? (a) : (b))
4423 static void itrunc(struct inode*);
4424
4425 // Read the super block.
4426 void
4427 readsb(int dev, struct superblock *sb)
4428 {
4429     struct buf *bp;
4430
4431     bp = bread(dev, 1);
4432     memmove(sb, bp->data, sizeof(*sb));
4433     brelse(bp);
4434 }
4435
4436 // Zero a block.
4437 static void
4438 bzero(int dev, int bno)
4439 {
4440     struct buf *bp;
4441
4442     bp = bread(dev, bno);
4443     memset(bp->data, 0, BSIZE);
4444     log_write(bp);
4445     brelse(bp);
4446 }
4447
4448
4449

```

```

4450 // Blocks.
4451
4452 // Allocate a zeroed disk block.
4453 static uint
4454 balloc(uint dev)
4455 {
4456     int b, bi, m;
4457     struct buf *bp;
4458     struct superblock sb;
4459
4460     bp = 0;
4461     readsb(dev, &sb);
4462     for(b = 0; b < sb.size; b += BPB){
4463         bp = bread(dev, BBLOCK(b, sb.ninodes));
4464         for(bi = 0; bi < BPB && b + bi < sb.size; bi++){
4465             m = 1 << (bi % 8);
4466             if((bp->data[bi/8] & m) == 0){ // Is block free?
4467                 bp->data[bi/8] |= m; // Mark block in use.
4468                 log_write(bp);
4469                 brelse(bp);
4470                 bzero(dev, b + bi);
4471                 return b + bi;
4472             }
4473         }
4474         brelse(bp);
4475     }
4476     panic("balloc: out of blocks");
4477 }
4478
4479 // Free a disk block.
4480 static void
4481 bfree(int dev, uint b)
4482 {
4483     struct buf *bp;
4484     struct superblock sb;
4485     int bi, m;
4486
4487     readsb(dev, &sb);
4488     bp = bread(dev, BBLOCK(b, sb.ninodes));
4489     bi = b % BPB;
4490     m = 1 << (bi % 8);
4491     if((bp->data[bi/8] & m) == 0)
4492         panic("freeing free block");
4493     bp->data[bi/8] &= ~m;
4494     log_write(bp);
4495     brelse(bp);
4496 }
4497
4498
4499

```

```

4500 // Inodes.
4501 //
4502 // An inode describes a single unnamed file.
4503 // The inode disk structure holds metadata: the file's type,
4504 // its size, the number of links referring to it, and the
4505 // list of blocks holding the file's content.
4506 //
4507 // The inodes are laid out sequentially on disk immediately after
4508 // the superblock. Each inode has a number, indicating its
4509 // position on the disk.
4510 //
4511 // The kernel keeps a cache of in-use inodes in memory
4512 // to provide a place for synchronizing access
4513 // to inodes used by multiple processes. The cached
4514 // inodes include book-keeping information that is
4515 // not stored on disk: ip->ref and ip->flags.
4516 //
4517 // An inode and its in-memory representative go through a
4518 // sequence of states before they can be used by the
4519 // rest of the file system code.
4520 //
4521 // * Allocation: an inode is allocated if its type (on disk)
4522 // is non-zero. ialloc() allocates, iput() frees if
4523 // the link count has fallen to zero.
4524 //
4525 // * Referencing in cache: an entry in the inode cache
4526 // is free if ip->ref is zero. Otherwise ip->ref tracks
4527 // the number of in-memory pointers to the entry (open
4528 // files and current directories). iget() to find or
4529 // create a cache entry and increment its ref, iput()
4530 // to decrement ref.
4531 //
4532 // * Valid: the information (type, size, &c) in an inode
4533 // cache entry is only correct when the I_VALID bit
4534 // is set in ip->flags. ilock() reads the inode from
4535 // the disk and sets I_VALID, while iput() clears
4536 // I_VALID if ip->ref has fallen to zero.
4537 //
4538 // * Locked: file system code may only examine and modify
4539 // the information in an inode and its content if it
4540 // has first locked the inode. The I_BUSY flag indicates
4541 // that the inode is locked. ilock() sets I_BUSY,
4542 // while iunlock clears it.
4543 //
4544 // Thus a typical sequence is:
4545 // ip = iget(dev, inum)
4546 // ilock(ip)
4547 // ... examine and modify ip->xxx ...
4548 // iunlock(ip)
4549 // iput(ip)

```

```

4550 //
4551 // ilock() is separate from iget() so that system calls can
4552 // get a long-term reference to an inode (as for an open file)
4553 // and only lock it for short periods (e.g., in read()).
4554 // The separation also helps avoid deadlock and races during
4555 // pathname lookup. iget() increments ip->ref so that the inode
4556 // stays cached and pointers to it remain valid.
4557 //
4558 // Many internal file system functions expect the caller to
4559 // have locked the inodes involved; this lets callers create
4560 // multi-step atomic operations.
4561 //
4562 struct {
4563   struct spinlock lock;
4564   struct inode inode[NINODE];
4565 } icache;
4566
4567 void
4568 init(void)
4569 {
4570   initlock(&icache.lock, "icache");
4571 }
4572
4573 static struct inode* iget(uint dev, uint inum);
4574
4575
4576
4577
4578
4579
4580
4581
4582
4583
4584
4585
4586
4587
4588
4589
4590
4591
4592
4593
4594
4595
4596
4597
4598
4599

```

```

4600 // Allocate a new inode with the given type on device dev.
4601 // A free inode has a type of zero.
4602 struct inode*
4603 ialloc(uint dev, short type)
4604 {
4605     int inum;
4606     struct buf *bp;
4607     struct dinode *dip;
4608     struct superblock sb;
4609
4610     readsb(dev, &sb);
4611
4612     for(inum = 1; inum < sb.ninodes; inum++){
4613         bp = bread(dev, IBLOCK(inum));
4614         dip = (struct dinode*)bp->data + inum%IPB;
4615         if(dip->type == 0){ // a free inode
4616             memset(dip, 0, sizeof(*dip));
4617             dip->type = type;
4618             log_write(bp); // mark it allocated on the disk
4619             brelse(bp);
4620             return iget(dev, inum);
4621         }
4622         brelse(bp);
4623     }
4624     panic("ialloc: no inodes");
4625 }
4626
4627 // Copy a modified in-memory inode to disk.
4628 void
4629 iupdate(struct inode *ip)
4630 {
4631     struct buf *bp;
4632     struct dinode *dip;
4633
4634     bp = bread(ip->dev, IBLOCK(ip->inum));
4635     dip = (struct dinode*)bp->data + ip->inum%IPB;
4636     dip->type = ip->type;
4637     dip->major = ip->major;
4638     dip->minor = ip->minor;
4639     dip->nlink = ip->nlink;
4640     dip->size = ip->size;
4641     memmove(dip->addrs, ip->addrs, sizeof(ip->addrs));
4642     log_write(bp);
4643     brelse(bp);
4644 }
4645
4646
4647
4648
4649

```

```

4650 // Find the inode with number inum on device dev
4651 // and return the in-memory copy. Does not lock
4652 // the inode and does not read it from disk.
4653 static struct inode*
4654 iget(uint dev, uint inum)
4655 {
4656     struct inode *ip, *empty;
4657
4658     acquire(&icache.lock);
4659
4660     // Is the inode already cached?
4661     empty = 0;
4662     for(ip = &icache.inode[0]; ip < &icache.inode[NINODE]; ip++){
4663         if(ip->ref > 0 && ip->dev == dev && ip->inum == inum){
4664             ip->ref++;
4665             release(&icache.lock);
4666             return ip;
4667         }
4668         if(empty == 0 && ip->ref == 0) // Remember empty slot.
4669             empty = ip;
4670     }
4671
4672     // Recycle an inode cache entry.
4673     if(empty == 0)
4674         panic("iget: no inodes");
4675
4676     ip = empty;
4677     ip->dev = dev;
4678     ip->inum = inum;
4679     ip->ref = 1;
4680     ip->flags = 0;
4681     release(&icache.lock);
4682
4683     return ip;
4684 }
4685
4686 // Increment reference count for ip.
4687 // Returns ip to enable ip = idup(ip1) idiom.
4688 struct inode*
4689 idup(struct inode *ip)
4690 {
4691     acquire(&icache.lock);
4692     ip->ref++;
4693     release(&icache.lock);
4694     return ip;
4695 }
4696
4697
4698
4699

```

```

4700 // Lock the given inode.
4701 // Reads the inode from disk if necessary.
4702 void
4703 ilock(struct inode *ip)
4704 {
4705     struct buf *bp;
4706     struct dinode *dip;
4707
4708     if(ip == 0 || ip->ref < 1)
4709         panic("ilock");
4710
4711     acquire(&icache.lock);
4712     while(ip->flags & I_BUSY)
4713         sleep(ip, &icache.lock);
4714     ip->flags |= I_BUSY;
4715     release(&icache.lock);
4716
4717     if(!(ip->flags & I_INVALID)){
4718         bp = bread(ip->dev, IBLOCK(ip->inum));
4719         dip = (struct dinode*)bp->data + ip->inum%IPB;
4720         ip->type = dip->type;
4721         ip->major = dip->major;
4722         ip->minor = dip->minor;
4723         ip->nlink = dip->nlink;
4724         ip->size = dip->size;
4725         memmove(ip->addrs, dip->addrs, sizeof(ip->addrs));
4726         brelse(bp);
4727         ip->flags |= I_VALID;
4728         if(ip->type == 0)
4729             panic("ilock: no type");
4730     }
4731 }
4732
4733 // Unlock the given inode.
4734 void
4735 iunlock(struct inode *ip)
4736 {
4737     if(ip == 0 || !(ip->flags & I_BUSY) || ip->ref < 1)
4738         panic("iunlock");
4739
4740     acquire(&icache.lock);
4741     ip->flags &= ~I_BUSY;
4742     wakeup(ip);
4743     release(&icache.lock);
4744 }
4745
4746
4747
4748
4749

```

```

4750 // Drop a reference to an in-memory inode.
4751 // If that was the last reference, the inode cache entry can
4752 // be recycled.
4753 // If that was the last reference and the inode has no links
4754 // to it, free the inode (and its content) on disk.
4755 void
4756 iput(struct inode *ip)
4757 {
4758     acquire(&icache.lock);
4759     if(ip->ref == 1 && (ip->flags & I_INVALID) && ip->nlink == 0){
4760         // inode has no links: truncate and free inode.
4761         if(ip->flags & I_BUSY)
4762             panic("iput busy");
4763         ip->flags |= I_BUSY;
4764         release(&icache.lock);
4765         itrunc(ip);
4766         ip->type = 0;
4767         iupdate(ip);
4768         acquire(&icache.lock);
4769         ip->flags = 0;
4770         wakeup(ip);
4771     }
4772     ip->ref--;
4773     release(&icache.lock);
4774 }
4775
4776 // Common idiom: unlock, then put.
4777 void
4778 iunlockput(struct inode *ip)
4779 {
4780     iunlock(ip);
4781     iput(ip);
4782 }
4783
4784
4785
4786
4787
4788
4789
4790
4791
4792
4793
4794
4795
4796
4797
4798
4799

```



```

4800 // Inode content
4801 //
4802 // The content (data) associated with each inode is stored
4803 // in blocks on the disk. The first NDIRECT block numbers
4804 // are listed in ip->addrs[]. The next NINDIRECT blocks are
4805 // listed in block ip->addrs[NDIRECT].
4806
4807 // Return the disk block address of the nth block in inode ip.
4808 // If there is no such block, bmap allocates one.
4809 static uint
4810 bmap(struct inode *ip, uint bn)
4811 {
4812     uint addr, *a;
4813     struct buf *bp;
4814
4815     if(bn < NDIRECT){
4816         if((addr = ip->addrs[bn]) == 0)
4817             ip->addrs[bn] = addr = balloc(ip->dev);
4818         return addr;
4819     }
4820     bn -= NDIRECT;
4821
4822     if(bn < NINDIRECT){
4823         // Load indirect block, allocating if necessary.
4824         if((addr = ip->addrs[NDIRECT]) == 0)
4825             ip->addrs[NDIRECT] = addr = balloc(ip->dev);
4826         bp = bread(ip->dev, addr);
4827         a = (uint*)bp->data;
4828         if((addr = a[bn]) == 0){
4829             a[bn] = addr = balloc(ip->dev);
4830             log_write(bp);
4831         }
4832         brelse(bp);
4833         return addr;
4834     }
4835
4836     panic("bmap: out of range");
4837 }
4838
4839
4840
4841
4842
4843
4844
4845
4846
4847
4848
4849

```

```

4850 // Truncate inode (discard contents).
4851 // Only called when the inode has no links
4852 // to it (no directory entries referring to it)
4853 // and has no in-memory reference to it (is
4854 // not an open file or current directory).
4855 static void
4856 itrunc(struct inode *ip)
4857 {
4858     int i, j;
4859     struct buf *bp;
4860     uint *a;
4861
4862     for(i = 0; i < NDIRECT; i++){
4863         if(ip->addrs[i]){
4864             bfree(ip->dev, ip->addrs[i]);
4865             ip->addrs[i] = 0;
4866         }
4867     }
4868
4869     if(ip->addrs[NDIRECT]){
4870         bp = bread(ip->dev, ip->addrs[NDIRECT]);
4871         a = (uint*)bp->data;
4872         for(j = 0; j < NINDIRECT; j++){
4873             if(a[j])
4874                 bfree(ip->dev, a[j]);
4875         }
4876         brelse(bp);
4877         bfree(ip->dev, ip->addrs[NDIRECT]);
4878         ip->addrs[NDIRECT] = 0;
4879     }
4880
4881     ip->size = 0;
4882     update(ip);
4883 }
4884
4885 // Copy stat information from inode.
4886 void
4887 stati(struct inode *ip, struct stat *st)
4888 {
4889     st->dev = ip->dev;
4890     st->ino = ip->inum;
4891     st->type = ip->type;
4892     st->nlink = ip->nlink;
4893     st->size = ip->size;
4894 }
4895
4896
4897
4898
4899

```

```

4900 // Read data from inode.
4901 int
4902 readi(struct inode *ip, char *dst, uint off, uint n)
4903 {
4904     uint tot, m;
4905     struct buf *bp;
4906
4907     if(ip->type == T_DEV){
4908         if(ip->major < 0 || ip->major >= NDEV || !devsw[ip->major].read)
4909             return -1;
4910         return devsw[ip->major].read(ip, dst, n);
4911     }
4912
4913     if(off > ip->size || off + n < off)
4914         return -1;
4915     if(off + n > ip->size)
4916         n = ip->size - off;
4917
4918     for(tot=0; tot<n; tot+=m, off+=m, dst+=m){
4919         bp = bread(ip->dev, bmap(ip, off/BSIZE));
4920         m = min(n - tot, BSIZE - off%BSIZE);
4921         memmove(dst, bp->data + off%BSIZE, m);
4922         brelse(bp);
4923     }
4924     return n;
4925 }
4926
4927
4928
4929
4930
4931
4932
4933
4934
4935
4936
4937
4938
4939
4940
4941
4942
4943
4944
4945
4946
4947
4948
4949

```

```

4950 // Write data to inode.
4951 int
4952 writei(struct inode *ip, char *src, uint off, uint n)
4953 {
4954     uint tot, m;
4955     struct buf *bp;
4956
4957     if(ip->type == T_DEV){
4958         if(ip->major < 0 || ip->major >= NDEV || !devsw[ip->major].write)
4959             return -1;
4960         return devsw[ip->major].write(ip, src, n);
4961     }
4962
4963     if(off > ip->size || off + n < off)
4964         return -1;
4965     if(off + n > MAXFILE*BSIZE)
4966         return -1;
4967
4968     for(tot=0; tot<n; tot+=m, off+=m, src+=m){
4969         bp = bread(ip->dev, bmap(ip, off/BSIZE));
4970         m = min(n - tot, BSIZE - off%BSIZE);
4971         memmove(bp->data + off%BSIZE, src, m);
4972         log_write(bp);
4973         brelse(bp);
4974     }
4975
4976     if(n > 0 && off > ip->size){
4977         ip->size = off;
4978         iupdate(ip);
4979     }
4980     return n;
4981 }
4982
4983
4984
4985
4986
4987
4988
4989
4990
4991
4992
4993
4994
4995
4996
4997
4998
4999

```

```

5000 // Directories
5001
5002 int
5003 namecmp(const char *s, const char *t)
5004 {
5005     return strcmp(s, t, DIRSIZ);
5006 }
5007
5008 // Look for a directory entry in a directory.
5009 // If found, set *poff to byte offset of entry.
5010 struct inode*
5011 dirlookup(struct inode *dp, char *name, uint *poff)
5012 {
5013     uint off, inum;
5014     struct dirent de;
5015
5016     if(dp->type != T_DIR)
5017         panic("dirlookup not DIR");
5018
5019     for(off = 0; off < dp->size; off += sizeof(de)){
5020         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5021             panic("dirlink read");
5022         if(de.inum == 0)
5023             continue;
5024         if(namecmp(name, de.name) == 0){
5025             // entry matches path element
5026             if(poff)
5027                 *poff = off;
5028             inum = de.inum;
5029             return iget(dp->dev, inum);
5030         }
5031     }
5032
5033     return 0;
5034 }
5035
5036
5037
5038
5039
5040
5041
5042
5043
5044
5045
5046
5047
5048
5049

```

```

5050 // Write a new directory entry (name, inum) into the directory dp.
5051 int
5052 dirlink(struct inode *dp, char *name, uint inum)
5053 {
5054     int off;
5055     struct dirent de;
5056     struct inode *ip;
5057
5058     // Check that name is not present.
5059     if((ip = dirlookup(dp, name, 0)) != 0){
5060         iput(ip);
5061         return -1;
5062     }
5063
5064     // Look for an empty dirent.
5065     for(off = 0; off < dp->size; off += sizeof(de)){
5066         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5067             panic("dirlink read");
5068         if(de.inum == 0)
5069             break;
5070     }
5071
5072     strncpy(de.name, name, DIRSIZ);
5073     de.inum = inum;
5074     if(writei(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5075         panic("dirlink");
5076
5077     return 0;
5078 }
5079
5080
5081
5082
5083
5084
5085
5086
5087
5088
5089
5090
5091
5092
5093
5094
5095
5096
5097
5098
5099

```

```

5100 // Paths
5101
5102 // Copy the next path element from path into name.
5103 // Return a pointer to the element following the copied one.
5104 // The returned path has no leading slashes,
5105 // so the caller can check *path=='\0' to see if the name is the last one.
5106 // If no name to remove, return 0.
5107 //
5108 // Examples:
5109 // skipel("a/bb/c", name) = "bb/c", setting name = "a"
5110 // skipel("///a//bb", name) = "bb", setting name = "a"
5111 // skipel("a", name) = "", setting name = "a"
5112 // skipel("", name) = skipel("///", name) = 0
5113 //
5114 static char*
5115 skipel(char *path, char *name)
5116 {
5117   char *s;
5118   int len;
5119
5120   while(*path == '/')
5121     path++;
5122   if(*path == 0)
5123     return 0;
5124   s = path;
5125   while(*path != '/' && *path != 0)
5126     path++;
5127   len = path - s;
5128   if(len >= DIRSIZ)
5129     memmove(name, s, DIRSIZ);
5130   else {
5131     memmove(name, s, len);
5132     name[len] = 0;
5133   }
5134   while(*path == '/')
5135     path++;
5136   return path;
5137 }
5138
5139
5140
5141
5142
5143
5144
5145
5146
5147
5148
5149

```

```

5150 // Look up and return the inode for a path name.
5151 // If parent != 0, return the inode for the parent and copy the final
5152 // path element into name, which must have room for DIRSIZ bytes.
5153 static struct inode*
5154 namex(char *path, int nameparent, char *name)
5155 {
5156   struct inode *ip, *next;
5157
5158   if(*path == '/')
5159     ip = iget(ROOTDEV, ROOTINO);
5160   else
5161     ip = idup(proc->cwd);
5162
5163   while((path = skipel(path, name)) != 0){
5164     ilock(ip);
5165     if(ip->type != T_DIR){
5166       iunlockput(ip);
5167       return 0;
5168     }
5169     if(nameparent && *path == '\0'){
5170       // Stop one level early.
5171       iunlock(ip);
5172       return ip;
5173     }
5174     if((next = dirlookup(ip, name, 0)) == 0){
5175       iunlockput(ip);
5176       return 0;
5177     }
5178     iunlockput(ip);
5179     ip = next;
5180   }
5181   if(nameparent){
5182     iput(ip);
5183     return 0;
5184   }
5185   return ip;
5186 }
5187
5188 struct inode*
5189 namei(char *path)
5190 {
5191   char name[DIRSIZ];
5192   return namex(path, 0, name);
5193 }
5194
5195 struct inode*
5196 nameiparent(char *path, char *name)
5197 {
5198   return namex(path, 1, name);
5199 }

```

```

5200 //
5201 // File descriptors
5202 //
5203
5204 #include "types.h"
5205 #include "defs.h"
5206 #include "param.h"
5207 #include "fs.h"
5208 #include "file.h"
5209 #include "spinlock.h"
5210
5211 struct devsw devsw[NDEV];
5212 struct {
5213   struct spinlock lock;
5214   struct file file[NFILE];
5215 } ftable;
5216
5217 void
5218 fileinit(void)
5219 {
5220   initlock(&ftable.lock, "ftable");
5221 }
5222
5223 // Allocate a file structure.
5224 struct file*
5225 filealloc(void)
5226 {
5227   struct file *f;
5228
5229   acquire(&ftable.lock);
5230   for(f = ftable.file; f < ftable.file + NFILE; f++){
5231     if(f->ref == 0){
5232       f->ref = 1;
5233       release(&ftable.lock);
5234       return f;
5235     }
5236   }
5237   release(&ftable.lock);
5238   return 0;
5239 }
5240
5241
5242
5243
5244
5245
5246
5247
5248
5249

```

```

5250 // Increment ref count for file f.
5251 struct file*
5252 filedup(struct file *f)
5253 {
5254   acquire(&ftable.lock);
5255   if(f->ref < 1)
5256     panic("filedup");
5257   f->ref++;
5258   release(&ftable.lock);
5259   return f;
5260 }
5261
5262 // Close file f. (Decrement ref count, close when reaches 0.)
5263 void
5264 fileclose(struct file *f)
5265 {
5266   struct file ff;
5267
5268   acquire(&ftable.lock);
5269   if(f->ref < 1)
5270     panic("fileclose");
5271   if(--f->ref > 0){
5272     release(&ftable.lock);
5273     return;
5274   }
5275   ff = *f;
5276   f->ref = 0;
5277   f->type = FD_NONE;
5278   release(&ftable.lock);
5279
5280   if(ff.type == FD_PIPE)
5281     pipeclose(ff.pipe, ff.writable);
5282   else if(ff.type == FD_INODE){
5283     begin_trans();
5284     iput(ff.ip);
5285     commit_trans();
5286   }
5287 }
5288
5289
5290
5291
5292
5293
5294
5295
5296
5297
5298
5299

```

```

5300 // Get metadata about file f.
5301 int
5302 filestat(struct file *f, struct stat *st)
5303 {
5304     if(f->type == FD_INODE){
5305         ilock(f->ip);
5306         stati(f->ip, st);
5307         iunlock(f->ip);
5308         return 0;
5309     }
5310     return -1;
5311 }
5312
5313 // Read from file f.
5314 int
5315 fileread(struct file *f, char *addr, int n)
5316 {
5317     int r;
5318
5319     if(f->readable == 0)
5320         return -1;
5321     if(f->type == FD_PIPE)
5322         return piperead(f->pipe, addr, n);
5323     if(f->type == FD_INODE){
5324         ilock(f->ip);
5325         if((r = readi(f->ip, addr, f->off, n)) > 0)
5326             f->off += r;
5327         iunlock(f->ip);
5328         return r;
5329     }
5330     panic("fileread");
5331 }
5332
5333
5334
5335
5336
5337
5338
5339
5340
5341
5342
5343
5344
5345
5346
5347
5348
5349

```

```

5350 // Write to file f.
5351 int
5352 filewrite(struct file *f, char *addr, int n)
5353 {
5354     int r;
5355
5356     if(f->writable == 0)
5357         return -1;
5358     if(f->type == FD_PIPE)
5359         return pipewrite(f->pipe, addr, n);
5360     if(f->type == FD_INODE){
5361         // write a few blocks at a time to avoid exceeding
5362         // the maximum log transaction size, including
5363         // i-node, indirect block, allocation blocks,
5364         // and 2 blocks of slop for non-aligned writes.
5365         // this really belongs lower down, since writei()
5366         // might be writing a device like the console.
5367         int max = ((LOGSIZE-1-2) / 2) * 512;
5368         int i = 0;
5369         while(i < n){
5370             int n1 = n - i;
5371             if(n1 > max)
5372                 n1 = max;
5373
5374             begin_trans();
5375             ilock(f->ip);
5376             if ((r = writei(f->ip, addr + i, f->off, n1)) > 0)
5377                 f->off += r;
5378             iunlock(f->ip);
5379             commit_trans();
5380
5381             if(r < 0)
5382                 break;
5383             if(r != n1)
5384                 panic("short filewrite");
5385             i += r;
5386         }
5387         return i == n ? n : -1;
5388     }
5389     panic("filewrite");
5390 }
5391
5392
5393
5394
5395
5396
5397
5398
5399

```

```

5400 //
5401 // File-system system calls.
5402 // Mostly argument checking, since we don't trust
5403 // user code, and calls into file.c and fs.c.
5404 //
5405
5406 #include "types.h"
5407 #include "defs.h"
5408 #include "param.h"
5409 #include "stat.h"
5410 #include "mmu.h"
5411 #include "proc.h"
5412 #include "fs.h"
5413 #include "file.h"
5414 #include "fcntl.h"
5415
5416 // Fetch the nth word-sized system call argument as a file descriptor
5417 // and return both the descriptor and the corresponding struct file.
5418 static int
5419 argfd(int n, int *pfd, struct file **pf)
5420 {
5421     int fd;
5422     struct file *f;
5423
5424     if(argint(n, &fd) < 0)
5425         return -1;
5426     if(fd < 0 || fd >= NOFILE || (f=proc->ofile[fd]) == 0)
5427         return -1;
5428     if(pfd)
5429         *pfd = fd;
5430     if(pf)
5431         *pf = f;
5432     return 0;
5433 }
5434
5435 // Allocate a file descriptor for the given file.
5436 // Takes over file reference from caller on success.
5437 static int
5438 fdalloc(struct file *f)
5439 {
5440     int fd;
5441
5442     for(fd = 0; fd < NOFILE; fd++){
5443         if(proc->ofile[fd] == 0){
5444             proc->ofile[fd] = f;
5445             return fd;
5446         }
5447     }
5448     return -1;
5449 }

```

```

5450 int
5451 sys_dup(void)
5452 {
5453     struct file *f;
5454     int fd;
5455
5456     if(argfd(0, 0, &f) < 0)
5457         return -1;
5458     if((fd=fdalloc(f)) < 0)
5459         return -1;
5460     fildup(f);
5461     return fd;
5462 }
5463
5464 int
5465 sys_read(void)
5466 {
5467     struct file *f;
5468     int n;
5469     char *p;
5470
5471     if(argfd(0, 0, &f) < 0 || argint(2, &n) < 0 || argptr(1, &p, n) < 0)
5472         return -1;
5473     return fileread(f, p, n);
5474 }
5475
5476 int
5477 sys_write(void)
5478 {
5479     struct file *f;
5480     int n;
5481     char *p;
5482
5483     if(argfd(0, 0, &f) < 0 || argint(2, &n) < 0 || argptr(1, &p, n) < 0)
5484         return -1;
5485     return filewrite(f, p, n);
5486 }
5487
5488 int
5489 sys_close(void)
5490 {
5491     int fd;
5492     struct file *f;
5493
5494     if(argfd(0, &fd, &f) < 0)
5495         return -1;
5496     proc->ofile[fd] = 0;
5497     fileclose(f);
5498     return 0;
5499 }

```



```

5500 int
5501 sys_fstat(void)
5502 {
5503     struct file *f;
5504     struct stat *st;
5505
5506     if(argfd(0, 0, &f) < 0 || argptr(1, (void*)&st, sizeof(*st)) < 0)
5507         return -1;
5508     return filestat(f, st);
5509 }
5510
5511 // Create the path new as a link to the same inode as old.
5512 int
5513 sys_link(void)
5514 {
5515     char name[DIRSIZ], *new, *old;
5516     struct inode *dp, *ip;
5517
5518     if(argstr(0, &old) < 0 || argstr(1, &new) < 0)
5519         return -1;
5520     if((ip = namei(old)) == 0)
5521         return -1;
5522
5523     begin_trans();
5524
5525     ilock(ip);
5526     if(ip->type == T_DIR) {
5527         iunlockput(ip);
5528         commit_trans();
5529         return -1;
5530     }
5531
5532     ip->nlink++;
5533     iupdate(ip);
5534     iunlock(ip);
5535
5536     if((dp = nameiparent(new, name)) == 0)
5537         goto bad;
5538     ilock(dp);
5539     if(dp->dev != ip->dev || dirlink(dp, name, ip->inum) < 0) {
5540         iunlockput(dp);
5541         goto bad;
5542     }
5543     iunlockput(dp);
5544     iput(ip);
5545
5546     commit_trans();
5547
5548     return 0;
5549

```

```

5550 bad:
5551     ilock(ip);
5552     ip->nlink--;
5553     iupdate(ip);
5554     iunlockput(ip);
5555     commit_trans();
5556     return -1;
5557 }
5558
5559 // Is the directory dp empty except for "." and ".." ?
5560 static int
5561 isdirempty(struct inode *dp)
5562 {
5563     int off;
5564     struct dirent de;
5565
5566     for(off=2*sizeof(de); off<dp->size; off+=sizeof(de)){
5567         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5568             panic("isdirempty: readi");
5569         if(de.inum != 0)
5570             return 0;
5571     }
5572     return 1;
5573 }
5574
5575
5576
5577
5578
5579
5580
5581
5582
5583
5584
5585
5586
5587
5588
5589
5590
5591
5592
5593
5594
5595
5596
5597
5598
5599

```

```

5600 int
5601 sys_unlink(void)
5602 {
5603     struct inode *ip, *dp;
5604     struct dirent de;
5605     char name[DIRSIZ], *path;
5606     uint off;
5607
5608     if(argstr(0, &path) < 0)
5609         return -1;
5610     if((dp = nameiparent(path, name)) == 0)
5611         return -1;
5612
5613     begin_trans();
5614
5615     ilock(dp);
5616
5617     // Cannot unlink "." or "..".
5618     if(namecmp(name, ".") == 0 || namecmp(name, "..") == 0)
5619         goto bad;
5620
5621     if((ip = dirlookup(dp, name, &off)) == 0)
5622         goto bad;
5623     ilock(ip);
5624
5625     if(ip->nlink < 1)
5626         panic("unlink: nlink < 1");
5627     if(ip->type == T_DIR && !isdirempty(ip)){
5628         unlinkput(ip);
5629         goto bad;
5630     }
5631
5632     memset(&de, 0, sizeof(de));
5633     if(writei(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5634         panic("unlink: write");
5635     if(ip->type == T_DIR) {
5636         dp->nlink--;
5637         update(dp);
5638     }
5639     unlinkput(dp);
5640
5641     ip->nlink--;
5642     update(ip);
5643     unlinkput(ip);
5644
5645     commit_trans();
5646
5647     return 0;
5648
5649

```

```

5650 bad:
5651     unlinkput(dp);
5652     commit_trans();
5653     return -1;
5654 }
5655
5656 static struct inode*
5657 create(char *path, short type, short major, short minor)
5658 {
5659     uint off;
5660     struct inode *ip, *dp;
5661     char name[DIRSIZ];
5662
5663     if((dp = nameiparent(path, name)) == 0)
5664         return 0;
5665     ilock(dp);
5666
5667     if((ip = dirlookup(dp, name, &off)) != 0){
5668         unlinkput(dp);
5669         ilock(ip);
5670         if(type == T_FILE && ip->type == T_FILE)
5671             return ip;
5672         unlinkput(ip);
5673         return 0;
5674     }
5675
5676     if((ip = ialloc(dp->dev, type)) == 0)
5677         panic("create: ialloc");
5678
5679     ilock(ip);
5680     ip->major = major;
5681     ip->minor = minor;
5682     ip->nlink = 1;
5683     update(ip);
5684
5685     if(type == T_DIR){ // Create . and .. entries.
5686         dp->nlink++; // for "."
5687         update(dp);
5688         // No ip->nlink++ for ".": avoid cyclic ref count.
5689         if(dirlink(ip, ".", ip->inum) < 0 || dirlink(ip, "..", dp->inum) < 0)
5690             panic("create dots");
5691     }
5692
5693     if(dirlink(dp, name, ip->inum) < 0)
5694         panic("create: dirlink");
5695
5696     unlinkput(dp);
5697
5698     return ip;
5699 }

```

```

5700 int
5701 sys_open(void)
5702 {
5703     char *path;
5704     int fd, omode;
5705     struct file *f;
5706     struct inode *ip;
5707
5708     if(argstr(0, &path) < 0 || argint(1, &omode) < 0)
5709         return -1;
5710     if(omode & O_CREATE){
5711         begin_trans();
5712         ip = create(path, T_FILE, 0, 0);
5713         commit_trans();
5714         if(ip == 0)
5715             return -1;
5716     } else {
5717         if((ip = namei(path)) == 0)
5718             return -1;
5719         ilock(ip);
5720         if(ip->type == T_DIR && omode != O_RDONLY){
5721             unlockput(ip);
5722             return -1;
5723         }
5724     }
5725
5726     if((f = filealloc()) == 0 || (fd = fdalloc(f)) < 0){
5727         if(f)
5728             fileclose(f);
5729         unlockput(ip);
5730         return -1;
5731     }
5732     unlock(ip);
5733
5734     f->type = FD_INODE;
5735     f->ip = ip;
5736     f->off = 0;
5737     f->readable = !(omode & O_WRONLY);
5738     f->writable = (omode & O_WRONLY) || (omode & O_RDWR);
5739     return fd;
5740 }
5741
5742
5743
5744
5745
5746
5747
5748
5749

```

```

5750 int
5751 sys_mkdir(void)
5752 {
5753     char *path;
5754     struct inode *ip;
5755
5756     begin_trans();
5757     if(argstr(0, &path) < 0 || (ip = create(path, T_DIR, 0, 0)) == 0){
5758         commit_trans();
5759         return -1;
5760     }
5761     unlockput(ip);
5762     commit_trans();
5763     return 0;
5764 }
5765
5766 int
5767 sys_mknod(void)
5768 {
5769     struct inode *ip;
5770     char *path;
5771     int len;
5772     int major, minor;
5773
5774     begin_trans();
5775     if((len=argstr(0, &path)) < 0 ||
5776        argint(1, &major) < 0 ||
5777        argint(2, &minor) < 0 ||
5778        (ip = create(path, T_DEV, major, minor)) == 0){
5779         commit_trans();
5780         return -1;
5781     }
5782     unlockput(ip);
5783     commit_trans();
5784     return 0;
5785 }
5786
5787
5788
5789
5790
5791
5792
5793
5794
5795
5796
5797
5798
5799

```

```

5800 int
5801 sys_chdir(void)
5802 {
5803     char *path;
5804     struct inode *ip;
5805
5806     if(argstr(0, &path) < 0 || (ip = namei(path)) == 0)
5807         return -1;
5808     ilock(ip);
5809     if(ip->type != T_DIR){
5810         iunlockput(ip);
5811         return -1;
5812     }
5813     iunlock(ip);
5814     iput(proc->cwd);
5815     proc->cwd = ip;
5816     return 0;
5817 }
5818
5819 int
5820 sys_exec(void)
5821 {
5822     char *path, *argv[MAXARG];
5823     int i;
5824     uint uargv, uarg;
5825
5826     if(argstr(0, &path) < 0 || argint(1, (int*)&uargv) < 0){
5827         return -1;
5828     }
5829     memset(argv, 0, sizeof(argv));
5830     for(i=0; i++;){
5831         if(i >= NELEM(argv))
5832             return -1;
5833         if(fetchint(uargv+4*i, (int*)&uarg) < 0)
5834             return -1;
5835         if(uarg == 0){
5836             argv[i] = 0;
5837             break;
5838         }
5839         if(fetchstr(uarg, &argv[i]) < 0)
5840             return -1;
5841     }
5842     return exec(path, argv);
5843 }
5844
5845
5846
5847
5848
5849

```

```

5850 int
5851 sys_pipe(void)
5852 {
5853     int *fd;
5854     struct file *rf, *wf;
5855     int fd0, fd1;
5856
5857     if(argptr(0, (void*)&fd, 2*sizeof(fd[0])) < 0)
5858         return -1;
5859     if(pipealloc(&rf, &wf) < 0)
5860         return -1;
5861     fd0 = -1;
5862     if((fd0 = fdalloc(rf)) < 0 || (fd1 = fdalloc(wf)) < 0){
5863         if(fd0 >= 0)
5864             proc->ofile[fd0] = 0;
5865         fileclose(rf);
5866         fileclose(wf);
5867         return -1;
5868     }
5869     fd[0] = fd0;
5870     fd[1] = fd1;
5871     return 0;
5872 }
5873
5874
5875
5876
5877
5878
5879
5880
5881
5882
5883
5884
5885
5886
5887
5888
5889
5890
5891
5892
5893
5894
5895
5896
5897
5898
5899

```

```

5900 #include "types.h"
5901 #include "param.h"
5902 #include "memlayout.h"
5903 #include "mmu.h"
5904 #include "proc.h"
5905 #include "defs.h"
5906 #include "x86.h"
5907 #include "elf.h"
5908
5909 int
5910 exec(char *path, char **argv)
5911 {
5912     char *s, *last;
5913     int i, off;
5914     uint argc, sz, sp, ustack[3+MAXARG+1];
5915     struct elfhdr elf;
5916     struct inode *ip;
5917     struct proghdr ph;
5918     pde_t *pgdir, *oldpgdir;
5919
5920     if((ip = namei(path)) == 0)
5921         return -1;
5922     ilock(ip);
5923     pgdir = 0;
5924
5925     // Check ELF header
5926     if(readi(ip, (char*)&elf, 0, sizeof(elf)) < sizeof(elf))
5927         goto bad;
5928     if(elf.magic != ELF_MAGIC)
5929         goto bad;
5930
5931     if((pgdir = setupkvm()) == 0)
5932         goto bad;
5933
5934     // Load program into memory.
5935     sz = 0;
5936     for(i=0, off=elf.phoff; i<elf.phnum; i++, off+=sizeof(ph)){
5937         if(readi(ip, (char*)&ph, off, sizeof(ph)) != sizeof(ph))
5938             goto bad;
5939         if(ph.type != ELF_PROG_LOAD)
5940             continue;
5941         if(ph.memsz < ph.filesz)
5942             goto bad;
5943         if((sz = allocvm(pgdir, sz, ph.vaddr + ph.memsz)) == 0)
5944             goto bad;
5945         if(loadvm(pgdir, (char*)&ph.vaddr, ip, ph.off, ph.filesz) < 0)
5946             goto bad;
5947     }
5948     iunlockput(ip);
5949     ip = 0;

```

```

5950 // Allocate two pages at the next page boundary.
5951 // Make the first inaccessible. Use the second as the user stack.
5952 sz = PGROUNDUP(sz);
5953 if((sz = allocvm(pgdir, sz, sz + 2*PGSIZE)) == 0)
5954     goto bad;
5955 clearpteu(pgdir, (char*)(sz - 2*PGSIZE));
5956 sp = sz;
5957
5958 // Push argument strings, prepare rest of stack in ustack.
5959 for(argc = 0; argv[argc]; argc++) {
5960     if(argc >= MAXARG)
5961         goto bad;
5962     sp = (sp - (strlen(argv[argc]) + 1)) & ~3;
5963     if(copyout(pgdir, sp, argv[argc], strlen(argv[argc]) + 1) < 0)
5964         goto bad;
5965     ustack[3+argc] = sp;
5966 }
5967 ustack[3+argc] = 0;
5968
5969 ustack[0] = 0xffffffff; // fake return PC
5970 ustack[1] = argc;
5971 ustack[2] = sp - (argc+1)*4; // argv pointer
5972
5973 sp -= (3+argc+1) * 4;
5974 if(copyout(pgdir, sp, ustack, (3+argc+1)*4) < 0)
5975     goto bad;
5976
5977 // Save program name for debugging.
5978 for(last=s=path; *s; s++)
5979     if(*s == '/')
5980         last = s+1;
5981 safestrcpy(proc->name, last, sizeof(proc->name));
5982
5983 // Commit to the user image.
5984 oldpgdir = proc->pgdir;
5985 proc->pgdir = pgdir;
5986 proc->sz = sz;
5987 proc->tf->eip = elf.entry; // main
5988 proc->tf->esp = sp;
5989 switchvm(proc);
5990 freevm(oldpgdir);
5991 return 0;
5992
5993 bad:
5994 if(pgdir)
5995     freevm(pgdir);
5996 if(ip)
5997     iunlockput(ip);
5998 return -1;
5999 }

```

```

6000 #include "types.h"
6001 #include "defs.h"
6002 #include "param.h"
6003 #include "mmu.h"
6004 #include "proc.h"
6005 #include "fs.h"
6006 #include "file.h"
6007 #include "spinlock.h"
6008
6009 #define PIPESIZE 512
6010
6011 struct pipe {
6012   struct spinlock lock;
6013   char data[PIPESIZE];
6014   uint nread; // number of bytes read
6015   uint nwrite; // number of bytes written
6016   int readopen; // read fd is still open
6017   int writeopen; // write fd is still open
6018 };
6019
6020 int
6021 pipealloc(struct file **f0, struct file **f1)
6022 {
6023   struct pipe *p;
6024
6025   p = 0;
6026   *f0 = *f1 = 0;
6027   if((*f0 = filealloc()) == 0 || (*f1 = filealloc()) == 0)
6028     goto bad;
6029   if((p = (struct pipe*)kalloc()) == 0)
6030     goto bad;
6031   p->readopen = 1;
6032   p->writeopen = 1;
6033   p->nwrite = 0;
6034   p->nread = 0;
6035   initlock(&p->lock, "pipe");
6036   (*f0)->type = FD_PIPE;
6037   (*f0)->readable = 1;
6038   (*f0)->writable = 0;
6039   (*f0)->pipe = p;
6040   (*f1)->type = FD_PIPE;
6041   (*f1)->readable = 0;
6042   (*f1)->writable = 1;
6043   (*f1)->pipe = p;
6044   return 0;
6045
6046
6047
6048
6049

```

```

6050   bad:
6051     if(p)
6052       kfree((char*)p);
6053     if(*f0)
6054       fileclose(*f0);
6055     if(*f1)
6056       fileclose(*f1);
6057     return -1;
6058   }
6059
6060 void
6061 pipeclose(struct pipe *p, int writable)
6062 {
6063   acquire(&p->lock);
6064   if(writable){
6065     p->writeopen = 0;
6066     wakeup(&p->nread);
6067   } else {
6068     p->readopen = 0;
6069     wakeup(&p->nwrite);
6070   }
6071   if(p->readopen == 0 && p->writeopen == 0){
6072     release(&p->lock);
6073     kfree((char*)p);
6074   } else
6075     release(&p->lock);
6076   }
6077
6078
6079 int
6080 pipewrite(struct pipe *p, char *addr, int n)
6081 {
6082   int i;
6083
6084   acquire(&p->lock);
6085   for(i = 0; i < n; i++){
6086     while(p->nwrite == p->nread + PIPESIZE){
6087       if(p->readopen == 0 || proc->killed) {
6088         release(&p->lock);
6089         return -1;
6090       }
6091       wakeup(&p->nread);
6092       sleep(&p->nwrite, &p->lock);
6093     }
6094     p->data[p->nwrite++ % PIPESIZE] = addr[i];
6095   }
6096   wakeup(&p->nread);
6097   release(&p->lock);
6098   return n;
6099 }

```

```

6100 int
6101 piperead(struct pipe *p, char *addr, int n)
6102 {
6103     int i;
6104
6105     acquire(&p->lock);
6106     while(p->nread == p->nwrite && p->writeopen){
6107         if(proc->killed){
6108             release(&p->lock);
6109             return -1;
6110         }
6111         sleep(&p->nread, &p->lock);
6112     }
6113     for(i = 0; i < n; i++){
6114         if(p->nread == p->nwrite)
6115             break;
6116         addr[i] = p->data[p->nread++ % PIPESIZE];
6117     }
6118     wakeup(&p->nwrite);
6119     release(&p->lock);
6120     return i;
6121 }
6122
6123
6124
6125
6126
6127
6128
6129
6130
6131
6132
6133
6134
6135
6136
6137
6138
6139
6140
6141
6142
6143
6144
6145
6146
6147
6148
6149

```

```

6150 #include "types.h"
6151 #include "x86.h"
6152
6153 void*
6154 memset(void *dst, int c, uint n)
6155 {
6156     if ((int)dst%4 == 0 && n%4 == 0){
6157         c &= 0xFF;
6158         stosl(dst, (c<<24)|(c<<16)|(c<<8)|c, n/4);
6159     } else
6160         stosb(dst, c, n);
6161     return dst;
6162 }
6163
6164 int
6165 memcmp(const void *v1, const void *v2, uint n)
6166 {
6167     const uchar *s1, *s2;
6168
6169     s1 = v1;
6170     s2 = v2;
6171     while(n-- > 0){
6172         if(*s1 != *s2)
6173             return *s1 - *s2;
6174         s1++, s2++;
6175     }
6176     return 0;
6177 }
6178
6179
6180 void*
6181 memmove(void *dst, const void *src, uint n)
6182 {
6183     const char *s;
6184     char *d;
6185
6186     s = src;
6187     d = dst;
6188     if(s < d && s + n > d){
6189         s += n;
6190         d += n;
6191         while(n-- > 0)
6192             *--d = *--s;
6193     } else
6194         while(n-- > 0)
6195             *d++ = *s++;
6196
6197     return dst;
6198 }
6199

```



```
6200 // memcpy exists to placate GCC. Use memmove.
6201 void*
6202 memcpy(void *dst, const void *src, uint n)
6203 {
6204     return memmove(dst, src, n);
6205 }
6206
6207 int
6208 strncmp(const char *p, const char *q, uint n)
6209 {
6210     while(n > 0 && *p && *p == *q)
6211         n--, p++, q++;
6212     if(n == 0)
6213         return 0;
6214     return (uchar)*p - (uchar)*q;
6215 }
6216
6217 char*
6218 strncpy(char *s, const char *t, int n)
6219 {
6220     char *os;
6221
6222     os = s;
6223     while(n-- > 0 && (*s++ = *t++) != 0)
6224         ;
6225     while(n-- > 0)
6226         *s++ = 0;
6227     return os;
6228 }
6229
6230 // Like strncpy but guaranteed to NUL-terminate.
6231 char*
6232 safestrcpy(char *s, const char *t, int n)
6233 {
6234     char *os;
6235
6236     os = s;
6237     if(n <= 0)
6238         return os;
6239     while(--n > 0 && (*s++ = *t++) != 0)
6240         ;
6241     *s = 0;
6242     return os;
6243 }
6244
6245
6246
6247
6248
6249
```

```
6250 int
6251 strlen(const char *s)
6252 {
6253     int n;
6254     for(n = 0; s[n]; n++)
6255         ;
6256     return n;
6257 }
6258
6259
6260
6261
6262
6263
6264
6265
6266
6267
6268
6269
6270
6271
6272
6273
6274
6275
6276
6277
6278
6279
6280
6281
6282
6283
6284
6285
6286
6287
6288
6289
6290
6291
6292
6293
6294
6295
6296
6297
6298
6299
```

```

6300 // See MultiProcessor Specification Version 1.1[14]
6301
6302 struct mp {          // floating pointer
6303     uchar signature[4];
6304     void *physaddr;
6305     uchar length;
6306     uchar specrev;
6307     uchar checksum;
6308     uchar type;
6309     uchar imcrp;
6310     uchar reserved[3];
6311 };
6312
6313 struct mpconf {
6314     uchar signature[4];
6315     ushort length;
6316     uchar version;
6317     uchar checksum;
6318     uchar product[20];
6319     uint *oemtable;
6320     ushort oemlength;
6321     ushort entry;
6322     uint *lapicaddr;
6323     ushort xlength;
6324     uchar xchecksum;
6325     uchar reserved;
6326 };
6327
6328 struct mpproc {
6329     uchar type;
6330     uchar apicid;
6331     uchar version;
6332     uchar flags;
6333     #define MPB00T 0x02
6334     uchar signature[4];
6335     uint feature;
6336     uchar reserved[8];
6337 };
6338
6339 struct mpioapic {
6340     uchar type;
6341     uchar apicno;
6342     uchar version;
6343     uchar flags;
6344     uint *addr;
6345 };
6346
6347
6348
6349
6350 // Table entry types
6351 #define MPPROC 0x00 // One per processor
6352 #define MPBUS 0x01 // One per bus
6353 #define MPIOAPIC 0x02 // One per I/O APIC
6354 #define MPIOINTR 0x03 // One per bus interrupt source
6355 #define MPLINTR 0x04 // One per system interrupt source
6356
6357
6358
6359
6360
6361
6362
6363
6364
6365
6366
6367
6368
6369
6370
6371
6372
6373
6374
6375
6376
6377
6378
6379
6380
6381
6382
6383
6384
6385
6386
6387
6388
6389
6390
6391
6392
6393
6394
6395
6396
6397
6398
6399

```

```

6400 // Multiprocessor support
6401 // Search memory for MP description structures.
6402 // http://developer.intel.com/design/pentium/datashts/24201606.pdf
6403
6404 #include "types.h"
6405 #include "defs.h"
6406 #include "param.h"
6407 #include "memlayout.h"
6408 #include "mp.h"
6409 #include "x86.h"
6410 #include "mmu.h"
6411 #include "proc.h"
6412
6413 struct cpu cpus[NCPU];
6414 static struct cpu *bcpu;
6415 int ismp;
6416 int ncpu;
6417 uchar ioapicid;
6418
6419 int
6420 mpbcpu(void)
6421 {
6422     return bcpu->cpu;
6423 }
6424
6425 static uchar
6426 sum(uchar *addr, int len)
6427 {
6428     int i, sum;
6429
6430     sum = 0;
6431     for(i=0; i<len; i++)
6432         sum += addr[i];
6433     return sum;
6434 }
6435
6436 // Look for an MP structure in the len bytes at addr.
6437 static struct mp*
6438 mpsearch1(uint a, int len)
6439 {
6440     uchar *e, *p, *addr;
6441
6442     addr = p2v(a);
6443     e = addr+len;
6444     for(p = addr; p < e; p += sizeof(struct mp))
6445         if(memcmp(p, "MP_", 4) == 0 && sum(p, sizeof(struct mp)) == 0)
6446             return (struct mp*)p;
6447     return 0;
6448 }
6449

```

```

6450 // Search for the MP Floating Pointer Structure, which according to the
6451 // spec is in one of the following three locations:
6452 // 1) in the first KB of the EBDA;
6453 // 2) in the last KB of system base memory;
6454 // 3) in the BIOS ROM between 0xE0000 and 0xFFFFF.
6455 static struct mp*
6456 mpsearch(void)
6457 {
6458     uchar *bda;
6459     uint p;
6460     struct mp *mp;
6461
6462     bda = (uchar *) P2V(0x400);
6463     if((p = ((bda[0x0F] < 8) | bda[0x0E] << 4))){
6464         if((mp = mpsearch1(p, 1024))
6465            return mp;
6466     } else {
6467         p = ((bda[0x14] < 8) | bda[0x13] * 1024);
6468         if((mp = mpsearch1(p-1024, 1024))
6469            return mp;
6470     }
6471     return mpsearch1(0xF0000, 0x10000);
6472 }
6473
6474 // Search for an MP configuration table. For now,
6475 // don't accept the default configurations (physaddr == 0).
6476 // Check for correct signature, calculate the checksum and,
6477 // if correct, check the version.
6478 // To do: check extended table checksum.
6479 static struct mpconf*
6480 mpconfig(struct mp **pmp)
6481 {
6482     struct mpconf *conf;
6483     struct mp *mp;
6484
6485     if((mp = mpsearch()) == 0 || mp->physaddr == 0)
6486         return 0;
6487     conf = (struct mpconf*) p2v((uint) mp->physaddr);
6488     if(memcmp(conf, "PCMP", 4) != 0)
6489         return 0;
6490     if(conf->version != 1 && conf->version != 4)
6491         return 0;
6492     if(sum((uchar*)conf, conf->length) != 0)
6493         return 0;
6494     *pmp = mp;
6495     return conf;
6496 }
6497
6498
6499

```

```

6500 void
6501 mpinit(void)
6502 {
6503     uchar *p, *e;
6504     struct mp *mp;
6505     struct mpconf *conf;
6506     struct mpproc *proc;
6507     struct mpioapic *ioapic;
6508
6509     bcpu = &cpus[0];
6510     if((conf = mpconfig(&mp)) == 0)
6511         return;
6512     ismp = 1;
6513     lapic = (uint*)conf->lapicaddr;
6514     for(p=(uchar*)(conf+1), e=(uchar*)conf+conf->length; p<e; ){
6515         switch(*p){
6516         case MPPROC:
6517             proc = (struct mpproc*)p;
6518             if(ncpu != proc->apicid){
6519                 cprintf("mpinit: ncpu=%d apicid=%d\n", ncpu, proc->apicid);
6520                 ismp = 0;
6521             }
6522             if(proc->flags & MPBOOT)
6523                 bcpu = &cpus[ncpu];
6524             cpus[ncpu].id = ncpu;
6525             ncpu++;
6526             p += sizeof(struct mpproc);
6527             continue;
6528         case MPIOAPIC:
6529             ioapic = (struct mpioapic*)p;
6530             ioapicid = ioapic->apicno;
6531             p += sizeof(struct mpioapic);
6532             continue;
6533         case MPBUS:
6534         case MPIOINTR:
6535         case MPLINTR:
6536             p += 8;
6537             continue;
6538         default:
6539             cprintf("mpinit: unknown config type %x\n", *p);
6540             ismp = 0;
6541         }
6542     }
6543     if(!ismp){
6544         // Didn't like what we found; fall back to no MP.
6545         ncpu = 1;
6546         lapic = 0;
6547         ioapicid = 0;
6548         return;
6549     }

```

```

6550     if(mp->imcrp){
6551         // Bochs doesn't support IMCR, so this doesn't run on Bochs.
6552         // But it would on real hardware.
6553         outb(0x22, 0x70); // Select IMCR
6554         outb(0x23, inb(0x23) | 1); // Mask external interrupts.
6555     }
6556 }
6557
6558
6559
6560
6561
6562
6563
6564
6565
6566
6567
6568
6569
6570
6571
6572
6573
6574
6575
6576
6577
6578
6579
6580
6581
6582
6583
6584
6585
6586
6587
6588
6589
6590
6591
6592
6593
6594
6595
6596
6597
6598
6599

```

```

6600 // The local APIC manages internal (non-I/O) interrupts.
6601 // See Chapter 8 & Appendix C of Intel processor manual volume 3.
6602
6603 #include "types.h"
6604 #include "defs.h"
6605 #include "memlayout.h"
6606 #include "traps.h"
6607 #include "mmu.h"
6608 #include "x86.h"
6609
6610 // Local APIC registers, divided by 4 for use as uint[] indices.
6611 #define ID (0x0020/4) // ID
6612 #define VER (0x0030/4) // Version
6613 #define TPR (0x0080/4) // Task Priority
6614 #define EOI (0x00B0/4) // EOI
6615 #define SVR (0x00F0/4) // Spurious Interrupt Vector
6616 #define ENABLE 0x00000100 // Unit Enable
6617 #define ESR (0x0280/4) // Error Status
6618 #define ICRLO (0x0300/4) // Interrupt Command
6619 #define INIT 0x00000500 // INIT/RESET
6620 #define STARTUP 0x00000600 // Startup IPI
6621 #define DELIVS 0x00001000 // Delivery status
6622 #define ASSERT 0x00004000 // Assert interrupt (vs deassert)
6623 #define DEASSERT 0x00000000
6624 #define LEVEL 0x00008000 // Level triggered
6625 #define BCST 0x00080000 // Send to all APICs, including self.
6626 #define BUSY 0x00001000
6627 #define FIXED 0x00000000
6628 #define ICRHI (0x0310/4) // Interrupt Command [63:32]
6629 #define TIMER (0x0320/4) // Local Vector Table 0 (TIMER)
6630 #define XI 0x00000008 // divide counts by 1
6631 #define PERIODIC 0x00020000 // Periodic
6632 #define PCINT (0x0340/4) // Performance Counter LVT
6633 #define LINT0 (0x0350/4) // Local Vector Table 1 (LINT0)
6634 #define LINT1 (0x0360/4) // Local Vector Table 2 (LINT1)
6635 #define ERROR (0x0370/4) // Local Vector Table 3 (ERROR)
6636 #define MASKED 0x00010000 // Interrupt masked
6637 #define TICC (0x0380/4) // Timer Initial Count
6638 #define TCCR (0x0390/4) // Timer Current Count
6639 #define TDCR (0x03E0/4) // Timer Divide Configuration
6640
6641 volatile uint *lapic; // Initialized in mp.c
6642
6643 static void
6644 lapicw(int index, int value)
6645 {
6646     lapic[index] = value;
6647     lapic[ID]; // wait for write to finish, by reading
6648 }
6649

```

```

6650 void
6651 lapicinit(void)
6652 {
6653     if(!lapic)
6654         return;
6655
6656     // Enable local APIC; set spurious interrupt vector.
6657     lapicw(SVR, ENABLE | (T_IRQ0 + IRQ_SPURIOUS));
6658
6659     // The timer repeatedly counts down at bus frequency
6660     // from lapic[TICR] and then issues an interrupt.
6661     // If xv6 cared more about precise timekeeping,
6662     // TICR would be calibrated using an external time source.
6663     lapicw(TDCR, XI);
6664     lapicw(TIMER, PERIODIC | (T_IRQ0 + IRQ_TIMER));
6665     lapicw(TICR, 10000000);
6666
6667     // Disable logical interrupt lines.
6668     lapicw(LINT0, MASKED);
6669     lapicw(LINT1, MASKED);
6670
6671     // Disable performance counter overflow interrupts
6672     // on machines that provide that interrupt entry.
6673     if((lapic[VER]>>16) & 0xFF) >= 4)
6674         lapicw(PCINT, MASKED);
6675
6676     // Map error interrupt to IRQ_ERROR.
6677     lapicw(ERROR, T_IRQ0 + IRQ_ERROR);
6678
6679     // Clear error status register (requires back-to-back writes).
6680     lapicw(ESR, 0);
6681     lapicw(ESR, 0);
6682
6683     // Ack any outstanding interrupts.
6684     lapicw(EOI, 0);
6685
6686     // Send an Init Level De-Assert to synchronise arbitration ID's.
6687     lapicw(ICRHI, 0);
6688     lapicw(ICRLO, BCST | INIT | LEVEL);
6689     while(lapic[ICRLO] & DELIVS)
6690         ;
6691
6692     // Enable interrupts on the APIC (but not on the processor).
6693     lapicw(TPR, 0);
6694 }
6695
6696
6697
6698
6699

```

```

6700 int
6701 cpunum(void)
6702 {
6703     // Cannot call cpu when interrupts are enabled:
6704     // result not guaranteed to last long enough to be used!
6705     // Would prefer to panic but even printing is chancy here:
6706     // almost everything, including printf and panic, calls cpu,
6707     // often indirectly through acquire and release.
6708     if(readeflags() & FL_IF) {
6709         static int n;
6710         if(n++ == 0)
6711             printf("cpu called from %x with interrupts enabled\n",
6712                    __builtin_return_address(0));
6713     }
6714
6715     if(lapic)
6716         return lapic[ID]>>24;
6717     return 0;
6718 }
6719
6720 // Acknowledge interrupt.
6721 void
6722 lapiceoi(void)
6723 {
6724     if(lapic)
6725         lapicw(EOI, 0);
6726 }
6727
6728 // Spin for a given number of microseconds.
6729 // On real hardware would want to tune this dynamically.
6730 void
6731 microdelay(int us)
6732 {
6733 }
6734
6735 #define IO_RTC 0x70
6736
6737 // Start additional processor running entry code at addr.
6738 // See Appendix B of MultiProcessor Specification.
6739 void
6740 lapicstartap(uchar apicid, uint addr)
6741 {
6742     int i;
6743     ushort *wrv;
6744
6745     // The BSP must initialize CMOS shutdown code to 0AH
6746     // and the warm reset vector (DWORD based at 40:67) to point at
6747     // the AP startup code prior to the [universal startup algorithm].
6748     outb(IO_RTC, 0xF); // offset 0xF is shutdown code
6749     outb(IO_RTC+1, 0x0A);

```

```

6750     wrv = (ushort*)P2V((0x40<<4 | 0x67)); // Warm reset vector
6751     wrv[0] = 0;
6752     wrv[1] = addr >> 4;
6753
6754     // "Universal startup algorithm."
6755     // Send INIT (level-triggered) interrupt to reset other CPU.
6756     lapicw(ICRHI, apicid<<24);
6757     lapicw(ICRLO, INIT | LEVEL | ASSERT);
6758     microdelay(200);
6759     lapicw(ICRLO, INIT | LEVEL);
6760     microdelay(100); // should be 10ms, but too slow in Bochs!
6761
6762     // Send startup IPI (twice!) to enter code.
6763     // Regular hardware is supposed to only accept a STARTUP
6764     // when it is in the halted state due to an INIT. So the second
6765     // should be ignored, but it is part of the official Intel algorithm.
6766     // Bochs complains about the second one. Too bad for Bochs.
6767     for(i = 0; i < 2; i++) {
6768         lapicw(ICRHI, apicid<<24);
6769         lapicw(ICRLO, STARTUP | (addr>>12));
6770         microdelay(200);
6771     }
6772 }
6773
6774
6775
6776
6777
6778
6779
6780
6781
6782
6783
6784
6785
6786
6787
6788
6789
6790
6791
6792
6793
6794
6795
6796
6797
6798
6799

```

```

6800 // The I/O APIC manages hardware interrupts for an SMP system.
6801 // http://www.intel.com/design/chipsets/datashts/29056601.pdf
6802 // See also picirq.c.
6803
6804 #include "types.h"
6805 #include "defs.h"
6806 #include "traps.h"
6807
6808 #define IOAPIC 0xFEC00000 // Default physical address of IO APIC
6809
6810 #define REG_ID 0x00 // Register index: ID
6811 #define REG_VER 0x01 // Register index: version
6812 #define REG_TABLE 0x10 // Redirection table base
6813
6814 // The redirection table starts at REG_TABLE and uses
6815 // two registers to configure each interrupt.
6816 // The first (low) register in a pair contains configuration bits.
6817 // The second (high) register contains a bitmask telling which
6818 // CPUs can serve that interrupt.
6819 #define INT_DISABLED 0x00010000 // Interrupt disabled
6820 #define INT_LEVEL 0x00008000 // Level-triggered (vs edge-)
6821 #define INT_ACTIVELOW 0x00002000 // Active low (vs high)
6822 #define INT_LOGICAL 0x00000800 // Destination is CPU id (vs APIC ID)
6823
6824 volatile struct ioapic *ioapic;
6825
6826 // IO APIC MMIO structure: write reg, then read or write data.
6827 struct ioapic {
6828     uint reg;
6829     uint pad[3];
6830     uint data;
6831 };
6832
6833 static uint
6834 ioapicread(int reg)
6835 {
6836     ioapic->reg = reg;
6837     return ioapic->data;
6838 }
6839
6840 static void
6841 ioapicwrite(int reg, uint data)
6842 {
6843     ioapic->reg = reg;
6844     ioapic->data = data;
6845 }
6846
6847
6848
6849

```

```

6850 void
6851 ioapicinit(void)
6852 {
6853     int i, id, maxintr;
6854
6855     if(!ismp)
6856         return;
6857
6858     ioapic = (volatile struct ioapic*)IOAPIC;
6859     maxintr = (ioapicread(REG_VER) >> 16) & 0xFF;
6860     id = ioapicread(REG_ID) >> 24;
6861     if(id != ioapicid)
6862         cprintf("ioapicinit: id isn't equal to ioapicid; not a MP\n");
6863
6864     // Mark all interrupts edge-triggered, active high, disabled,
6865     // and not routed to any CPUs.
6866     for(i = 0; i <= maxintr; i++){
6867         ioapicwrite(REG_TABLE+2*i, INT_DISABLED | (T_IRQ0 + i));
6868         ioapicwrite(REG_TABLE+2*i+1, 0);
6869     }
6870 }
6871
6872 void
6873 ioapicenable(int irq, int cpunum)
6874 {
6875     if(!ismp)
6876         return;
6877
6878     // Mark interrupt edge-triggered, active high,
6879     // enabled, and routed to the given cpunum,
6880     // which happens to be that cpu's APIC ID.
6881     ioapicwrite(REG_TABLE+2*irq, T_IRQ0 + irq);
6882     ioapicwrite(REG_TABLE+2*irq+1, cpunum << 24);
6883 }
6884
6885
6886
6887
6888
6889
6890
6891
6892
6893
6894
6895
6896
6897
6898
6899

```



```

6900 // Intel 8259A programmable interrupt controllers.
6901
6902 #include "types.h"
6903 #include "x86.h"
6904 #include "traps.h"
6905
6906 // I/O Addresses of the two programmable interrupt controllers
6907 #define IO_PIC1 0x20 // Master (IRQs 0-7)
6908 #define IO_PIC2 0xA0 // Slave (IRQs 8-15)
6909
6910 #define IRQ_SLAVE 2 // IRQ at which slave connects to master
6911
6912 // Current IRQ mask.
6913 // Initial IRQ mask has interrupt 2 enabled (for slave 8259A).
6914 static ushort irqmask = 0xFFFF & ~(1<<IRQ_SLAVE);
6915
6916 static void
6917 picsetmask(ushort mask)
6918 {
6919     irqmask = mask;
6920     outb(IO_PIC1+1, mask);
6921     outb(IO_PIC2+1, mask >> 8);
6922 }
6923
6924 void
6925 picenable(int irq)
6926 {
6927     picsetmask(irqmask & ~(1<<irq));
6928 }
6929
6930 // Initialize the 8259A interrupt controllers.
6931 void
6932 picinit(void)
6933 {
6934     // mask all interrupts
6935     outb(IO_PIC1+1, 0xFF);
6936     outb(IO_PIC2+1, 0xFF);
6937
6938     // Set up master (8259A-1)
6939
6940     // ICW1: 0001g0hi
6941     // g: 0 = edge triggering, 1 = level triggering
6942     // h: 0 = cascaded PICs, 1 = master only
6943     // i: 0 = no ICW4, 1 = ICW4 required
6944     outb(IO_PIC1, 0x11);
6945
6946     // ICW2: Vector offset
6947     outb(IO_PIC1+1, T_IRQ0);
6948
6949

```

```

6950 // ICW3: (master PIC) bit mask of IR lines connected to slaves
6951 // (slave PIC) 3-bit # of slave's connection to master
6952 outb(IO_PIC1+1, 1<<IRQ_SLAVE);
6953
6954 // ICW4: 000nbmap
6955 // n: 1 = special fully nested mode
6956 // b: 1 = buffered mode
6957 // m: 0 = slave PIC, 1 = master PIC
6958 // (ignored when b is 0, as the master/slave role
6959 // can be hardwired).
6960 // a: 1 = Automatic EOI mode
6961 // p: 0 = MCS-80/85 mode, 1 = intel x86 mode
6962 outb(IO_PIC1+1, 0x3);
6963
6964 // Set up slave (8259A-2)
6965 outb(IO_PIC2, 0x11); // ICW1
6966 outb(IO_PIC2+1, T_IRQ0 + 8); // ICW2
6967 outb(IO_PIC2+1, IRQ_SLAVE); // ICW3
6968 // NB Automatic EOI mode doesn't tend to work on the slave.
6969 // Linux source code says it's "to be investigated".
6970 outb(IO_PIC2+1, 0x3); // ICW4
6971
6972 // OCW3: 0ef01prs
6973 // ef: 0x = NOP, 10 = clear specific mask, 11 = set specific mask
6974 // p: 0 = no polling, 1 = polling mode
6975 // rs: 0x = NOP, 10 = read IRR, 11 = read ISR
6976 outb(IO_PIC1, 0x68); // clear specific mask
6977 outb(IO_PIC1, 0x0a); // read IRR by default
6978
6979 outb(IO_PIC2, 0x68); // OCW3
6980 outb(IO_PIC2, 0x0a); // OCW3
6981
6982 if(irqmask != 0xFFFF)
6983     picsetmask(irqmask);
6984 }
6985
6986
6987
6988
6989
6990
6991
6992
6993
6994
6995
6996
6997
6998
6999

```

```

7000 // PC keyboard interface constants
7001
7002 #define KBSTATP 0x64 // kbd controller status port(I)
7003 #define KBD_DIB 0x01 // kbd data in buffer
7004 #define KBDATAP 0x60 // kbd data port(I)
7005
7006 #define NO 0
7007
7008 #define SHIFT (1<<0)
7009 #define CTL (1<<1)
7010 #define ALT (1<<2)
7011
7012 #define CAPSLOCK (1<<3)
7013 #define NUMLOCK (1<<4)
7014 #define SCROLLLOCK (1<<5)
7015
7016 #define E0ESC (1<<6)
7017
7018 // Special keycodes
7019 #define KEY_HOME 0xE0
7020 #define KEY_END 0xE1
7021 #define KEY_UP 0xE2
7022 #define KEY_DN 0xE3
7023 #define KEY_LF 0xE4
7024 #define KEY_RT 0xE5
7025 #define KEY_PGUP 0xE6
7026 #define KEY_PGDN 0xE7
7027 #define KEY_INS 0xE8
7028 #define KEY_DEL 0xE9
7029
7030 // CC('A') == Control-A
7031 #define C(x) (x - '@')
7032
7033 static uchar shiftcode[256] =
7034 {
7035   [0x1D] CTL,
7036   [0x2A] SHIFT,
7037   [0x36] SHIFT,
7038   [0x38] ALT,
7039   [0x9D] CTL,
7040   [0xB8] ALT
7041 };
7042
7043 static uchar togglecode[256] =
7044 {
7045   [0x3A] CAPSLOCK,
7046   [0x45] NUMLOCK,
7047   [0x46] SCROLLLOCK
7048 };
7049

```

```

7050 static uchar normalmap[256] =
7051 {
7052   NO, 0x1B, '1', '2', '3', '4', '5', '6', // 0x00
7053   '7', '8', '9', '0', '=', '\b', '\t', // 0x10
7054   'q', 'w', 'e', 'r', 't', 'y', '\n', // 0x20
7055   'o', 'p', '[', ']', '\n', 'a', 's', // 0x30
7056   'd', 'f', 'g', 'h', 'j', 'k', 'l', // 0x40
7057   '\'', ',', NO, '\\', 'z', 'x', 'c', 'v', // 0x50
7058   'b', 'n', 'm', '.', '/', NO, '*', // 0x60
7059   NO, NO, NO, NO, NO, NO, NO, NO, // 0x70
7060   NO, NO, NO, NO, NO, NO, NO, NO, // 0x80
7061   '8', '9', '-', '4', '5', '6', '+', // 0x90
7062   '2', '3', '0', '.', NO, NO, NO, NO, // 0xA0
7063   [0x9C] '\n', // KP_Enter
7064   [0x85] '/', // KP_Div
7065   [0xC8] KEY_UP, [0xD0] KEY_DN,
7066   [0xC9] KEY_PGUP, [0xD1] KEY_PGDN,
7067   [0xCB] KEY_LF, [0xCD] KEY_RT,
7068   [0x97] KEY_HOME, [0xCF] KEY_END,
7069   [0xD2] KEY_INS, [0xD3] KEY_DEL
7070 };
7071
7072 static uchar shiftmap[256] =
7073 {
7074   NO, 0x33, '!', '@', '#', '$', '%', '^', // 0x00
7075   '&', '*', '(', ')', '-', '=', '\b', '\t', // 0x10
7076   'Q', 'W', 'E', 'R', 'T', 'Y', '\n', // 0x20
7077   'O', 'P', '[', ']', '\n', 'A', 'S', // 0x30
7078   'D', 'F', 'G', 'H', 'J', 'K', 'L', // 0x40
7079   'I', '~', NO, '|', 'Z', 'X', 'C', 'V', // 0x50
7080   'B', 'N', '-', 'M', '<', '>', '?', '*', // 0x60
7081   NO, NO, NO, NO, NO, NO, NO, NO, // 0x70
7082   NO, NO, NO, NO, NO, NO, NO, NO, // 0x80
7083   '8', '9', '-', '4', '5', '6', '+', // 0x90
7084   '2', '3', '0', '.', NO, NO, NO, NO, // 0xA0
7085   [0x9C] '\n', // KP_Enter
7086   [0x85] '/', // KP_Div
7087   [0xC8] KEY_UP, [0xD0] KEY_DN,
7088   [0xC9] KEY_PGUP, [0xD1] KEY_PGDN,
7089   [0xCB] KEY_LF, [0xCD] KEY_RT,
7090   [0x97] KEY_HOME, [0xCF] KEY_END,
7091   [0xD2] KEY_INS, [0xD3] KEY_DEL
7092 };
7093
7094
7095
7096
7097
7098
7099

```

```

7100 static uchar ctlmap[256] =
7101 {
7102     NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO,
7103     NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO, NO,
7104     C('Q'), C('W'), C('E'), C('R'), C('T'), C('U'), C('I'),
7105     C('O'), C('P'), NO, NO, '\r', C('A'), C('S'),
7106     C('D'), C('F'), C('G'), C('H'), C('J'), C('K'), NO,
7107     NO, NO, C('\\'), C('Z'), C('X'), C('C'), C('V'),
7108     C('B'), C('N'), C('M'), NO, C('/'), NO,
7109     [0x9C] '\r', // KP_Enter
7110     [0xB5] C('/'), // KP_Div
7111     [0xC8] KEY_UP, [0xD0] KEY_DN,
7112     [0xC9] KEY_PGUP, [0xD1] KEY_PGDN,
7113     [0xCB] KEY_LF, [0xCD] KEY_RT,
7114     [0x97] KEY_HOME, [0xCF] KEY_END,
7115     [0xD2] KEY_INS, [0xD3] KEY_DEL
7116 };
7117
7118
7119
7120
7121
7122
7123
7124
7125
7126
7127
7128
7129
7130
7131
7132
7133
7134
7135
7136
7137
7138
7139
7140
7141
7142
7143
7144
7145
7146
7147
7148
7149

```

```

7150 #include "types.h"
7151 #include "x86.h"
7152 #include "defs.h"
7153 #include "kbd.h"
7154
7155 int
7156 kbdgetc(void)
7157 {
7158     static uint shift;
7159     static uchar *charcode[4] = {
7160         normalmap, shiftmap, ctlmap, ctlmap
7161     };
7162     uint st, data, c;
7163
7164     st = inb(KBSTATP);
7165     if((st & KBS_DIB) == 0)
7166         return -1;
7167     data = inb(KBDATAP);
7168
7169     if(data == 0xE0){
7170         shift |= E0ESC;
7171         return 0;
7172     } else if(data & 0x80){
7173         // Key released
7174         data = (shift & E0ESC ? data : data & 0x7F);
7175         shift &= ~(shiftcode[data] | E0ESC);
7176         return 0;
7177     } else if(shift & E0ESC){
7178         // Last character was an E0 escape; or with 0x80
7179         data |= 0x80;
7180         shift &= ~E0ESC;
7181     }
7182
7183     shift |= shiftcode[data];
7184     shift ^= togglecode[data];
7185     c = charcode[shift & (CTL | SHIFT)][data];
7186     if(shift & CAPSLOCK){
7187         if('a' <= c && c <= 'z')
7188             c += 'A' - 'a';
7189         else if('A' <= c && c <= 'Z')
7190             c += 'a' - 'A';
7191     }
7192     return c;
7193 }
7194
7195 void
7196 kbdtintr(void)
7197 {
7198     consoleintr(kbdgetc);
7199 }

```

```

7200 // Console input and output.
7201 // Input is from the keyboard or serial port.
7202 // Output is written to the screen and serial port.
7203
7204 #include "types.h"
7205 #include "defs.h"
7206 #include "param.h"
7207 #include "traps.h"
7208 #include "spinlock.h"
7209 #include "fs.h"
7210 #include "file.h"
7211 #include "memlayout.h"
7212 #include "mmu.h"
7213 #include "proc.h"
7214 #include "x86.h"
7215
7216 static void consputc(int);
7217
7218 static int panicked = 0;
7219
7220 static struct {
7221   struct spinlock lock;
7222   int locking;
7223 } cons;
7224
7225 static void
7226 printint(int xx, int base, int sign)
7227 {
7228   static char digits[] = "0123456789abcdef";
7229   char buf[16];
7230   int i;
7231   uint x;
7232
7233   if(sign && (sign == xx < 0))
7234     x = -xx;
7235   else
7236     x = xx;
7237
7238   i = 0;
7239   do{
7240     buf[i++] = digits[x % base];
7241   }while((x /= base) != 0);
7242
7243   if(sign)
7244     buf[i++] = '-';
7245
7246   while(--i >= 0)
7247     consputc(buf[i]);
7248 }
7249

```

```

7250 // Print to the console. only understands %d, %x, %p, %s.
7251 void
7252 cprintf(char *fmt, ...)
7253 {
7254   int i, c, locking;
7255   uint *argp;
7256   char *s;
7257
7258   locking = cons.locking;
7259   if(locking)
7260     acquire(&cons.lock);
7261
7262   if (fmt == 0)
7263     panic("null fmt");
7264
7265   argp = (uint*)(void*)&fmt + 1;
7266   for(i = 0; (c = fmt[i] & 0xff) != 0; i++){
7267     if(c != '%'){
7268       consputc(c);
7269       continue;
7270     }
7271     c = fmt[++i] & 0xff;
7272     if(c == 0)
7273       break;
7274     switch(c){
7275     case 'd':
7276       printint(*argp++, 10, 1);
7277       break;
7278     case 'x':
7279     case 'p':
7280       printint(*argp++, 16, 0);
7281       break;
7282     case 's':
7283       if((s = (char*)*argp++) == 0)
7284         s = "(null)";
7285       for(; *s; s++)
7286         consputc(*s);
7287       break;
7288     case '%':
7289       consputc('%');
7290       break;
7291     default:
7292       // Print unknown % sequence to draw attention.
7293       consputc('%');
7294       consputc(c);
7295       break;
7296     }
7297   }
7298
7299

```

```

7300 if(locking)
7301     release(&cons.lock);
7302 }
7303
7304 void
7305 panic(char *s)
7306 {
7307     int i;
7308     uint pcs[10];
7309
7310     cli();
7311     cons.locking = 0;
7312     cprintf("cpu%d: panic: ", cpu->id);
7313     cprintf(s);
7314     cprintf("\n");
7315     getcallerpcs(&s, pcs);
7316     for(i=0; i<10; i++)
7317         cprintf("%p", pcs[i]);
7318     panicked = 1; // freeze other CPU
7319     for(;;)
7320         ;
7321 }
7322
7323
7324
7325
7326
7327
7328
7329
7330
7331
7332
7333
7334
7335
7336
7337
7338
7339
7340
7341
7342
7343
7344
7345
7346
7347
7348
7349

```

```

7350 #define BACKSPACE 0x100
7351 #define CRTPORT 0x3d4
7352 static ushort *crt = (ushort*)P2V(0xb8000); // CGA memory
7353
7354 static void
7355 cputc(int c)
7356 {
7357     int pos;
7358
7359     // Cursor position: col + 80*row.
7360     outb(CRTPORT, 14);
7361     pos = inb(CRTPORT+1) << 8;
7362     outb(CRTPORT, 15);
7363     pos |= inb(CRTPORT+1);
7364
7365     if(c == '\n')
7366         pos += 80 - pos%80;
7367     else if(c == BACKSPACE){
7368         if(pos > 0) --pos;
7369     } else
7370         crt[pos++] = (c&0xff) | 0x0700; // black on white
7371
7372     if((pos/80) >= 24){ // Scroll up.
7373         memmove(crt, crt+80, sizeof(crt[0])*23*80);
7374         pos -= 80;
7375         memset(crt+pos, 0, sizeof(crt[0])*(24*80 - pos));
7376     }
7377
7378     outb(CRTPORT, 14);
7379     outb(CRTPORT+1, pos>>8);
7380     outb(CRTPORT, 15);
7381     outb(CRTPORT+1, pos);
7382     crt[pos] = ' ' | 0x0700;
7383 }
7384
7385 void
7386 consputc(int c)
7387 {
7388     if(panicked){
7389         cli();
7390         for(;;)
7391             ;
7392     }
7393
7394     if(c == BACKSPACE){
7395         uartputc('\b'); uartputc(' '); uartputc('\b');
7396     } else
7397         uartputc(c);
7398     cputc(c);
7399 }

```

```

7400 #define INPUT_BUF 128
7401 struct {
7402   struct spinlock lock;
7403   char buf[INPUT_BUF];
7404   uint r; // Read index
7405   uint w; // Write index
7406   uint e; // Edit index
7407 } input;
7408
7409 #define C(x) ((x) - '@') // Control-x
7410
7411 void
7412 consoleintr(int (*getc)(void))
7413 {
7414   int c;
7415
7416   acquire(&input.lock);
7417   while((c = getc()) >= 0){
7418     switch(c){
7419     case C('P'): // Process listing.
7420       procdump();
7421       break;
7422     case C('U'): // Kill line.
7423       while(input.e != input.w &&
7424             input.buf[(input.e-1) % INPUT_BUF] != '\n'){
7425         input.e--;
7426         consputc(BACKSPACE);
7427       }
7428       break;
7429     case C('H'): case '\x7f': // Backspace
7430       if(input.e != input.w){
7431         input.e--;
7432         consputc(BACKSPACE);
7433       }
7434       break;
7435     default:
7436       if(c != 0 && input.e - input.r < INPUT_BUF){
7437         c = (c == '\r') ? '\n' : c;
7438         input.buf[input.e++] = INPUT_BUF;
7439         consputc(c);
7440         if(c == '\n' || c == C('D') || input.e == input.r + INPUT_BUF){
7441           input.w = input.e;
7442           wakeup(&input.r);
7443         }
7444       }
7445       break;
7446     }
7447   }
7448   release(&input.lock);
7449 }

```

```

7450 int
7451 consoleread(struct inode *ip, char *dst, int n)
7452 {
7453   uint target;
7454   int c;
7455
7456   iunlock(ip);
7457   target = n;
7458   acquire(&input.lock);
7459   while(n > 0){
7460     while(input.r == input.w){
7461       if(proc->killed){
7462         release(&input.lock);
7463         ilock(ip);
7464         return -1;
7465       }
7466       sleep(&input.r, &input.lock);
7467     }
7468     c = input.buf[input.r++] % INPUT_BUF;
7469     if(c == C('D')){ // EOF
7470       if(n < target){
7471         // Save ^D for next time, to make sure
7472         // caller gets a 0-byte result.
7473         input.r--;
7474       }
7475       break;
7476     }
7477     *dst++ = c;
7478     --n;
7479     if(c == '\n')
7480       break;
7481   }
7482   release(&input.lock);
7483   ilock(ip);
7484
7485   return target - n;
7486 }
7487
7488
7489
7490
7491
7492
7493
7494
7495
7496
7497
7498
7499

```

```

7500 int
7501 consolewrite(struct inode *ip, char *buf, int n)
7502 {
7503     int i;
7504
7505     iunlock(ip);
7506     acquire(&cons.lock);
7507     for(i = 0; i < n; i++)
7508         consputc(buf[i] & 0xff);
7509     release(&cons.lock);
7510     ilock(ip);
7511
7512     return n;
7513 }
7514
7515 void
7516 consoleinit(void)
7517 {
7518     initlock(&cons.lock, "console");
7519     initlock(&input.lock, "input");
7520
7521     devsw[CONSOLE].write = consolewrite;
7522     devsw[CONSOLE].read = consoleread;
7523     cons.locking = 1;
7524
7525     piconable(IRQ_KBD);
7526     ioapicenable(IRQ_KBD, 0);
7527 }
7528
7529
7530
7531
7532
7533
7534
7535
7536
7537
7538
7539
7540
7541
7542
7543
7544
7545
7546
7547
7548
7549

```

```

7550 // Intel 8253/8254/82C54 Programmable Interval Timer (PIT).
7551 // Only used on uniprocessors;
7552 // SMP machines use the local APIC timer.
7553
7554 #include "types.h"
7555 #include "defs.h"
7556 #include "traps.h"
7557 #include "x86.h"
7558
7559 #define IO_TIMER1      0x040      // 8253 Timer #1
7560
7561 // Frequency of all three count-down timers;
7562 // (TIMER_FREQ/freq) is the appropriate count
7563 // to generate a frequency of freq Hz.
7564
7565 #define TIMER_FREQ      1193182
7566 #define TIMER_DIV(x)    ((TIMER_FREQ+(x)/2)/(x))
7567
7568 #define TIMER_MODE      (IO_TIMER1 + 3) // timer mode port
7569 #define TIMER_SELO      0x00          // select counter 0
7570 #define TIMER_RATEGEN    0x04          // mode 2, rate generator
7571 #define TIMER_16BIT      0x30          // r/w counter 16 bits, LSB first
7572
7573 void
7574 timerinit(void)
7575 {
7576     // Interrupt 100 times/sec.
7577     outb(TIMER_MODE, TIMER_SELO | TIMER_RATEGEN | TIMER_16BIT);
7578     outb(IO_TIMER1, TIMER_DIV(100) % 256);
7579     outb(IO_TIMER1, TIMER_DIV(100) / 256);
7580     piconable(IRQ_TIMER);
7581 }
7582
7583
7584
7585
7586
7587
7588
7589
7590
7591
7592
7593
7594
7595
7596
7597
7598
7599

```

```

7600 // Intel 8250 serial port (UART).
7601
7602 #include "types.h"
7603 #include "defs.h"
7604 #include "param.h"
7605 #include "traps.h"
7606 #include "spinlock.h"
7607 #include "fs.h"
7608 #include "file.h"
7609 #include "mmu.h"
7610 #include "proc.h"
7611 #include "x86.h"
7612
7613 #define COM1    0x3f8
7614
7615 static int uart;    // is there a uart?
7616
7617 void
7618 uartinit(void)
7619 {
7620     char *p;
7621
7622     // Turn off the FIFO
7623     outb(COM1+2, 0);
7624
7625     // 9600 baud, 8 data bits, 1 stop bit, parity off.
7626     outb(COM1+3, 0x80);    // Unlock divisor
7627     outb(COM1+0, 115200/9600);
7628     outb(COM1+1, 0);
7629     outb(COM1+3, 0x03);    // Lock divisor, 8 data bits.
7630     outb(COM1+4, 0);
7631     outb(COM1+1, 0x01);    // Enable receive interrupts.
7632
7633     // If status is 0xFF, no serial port.
7634     if(inb(COM1+5) == 0xFF)
7635         return;
7636     uart = 1;
7637
7638     // Acknowledge pre-existing interrupt conditions;
7639     // enable interrupts.
7640     inb(COM1+2);
7641     inb(COM1+0);
7642     picenable(IRQ_COM1);
7643     ioapicenable(IRQ_COM1, 0);
7644
7645     // Announce that we're here.
7646     for(p="xv6...\n"; *p; p++)
7647         uartputc(*p);
7648 }
7649

```

```

7650 void
7651 uartputc(int c)
7652 {
7653     int i;
7654
7655     if(!uart)
7656         return;
7657     for(i = 0; i < 128 && !(inb(COM1+5) & 0x20); i++)
7658         microdelay(10);
7659     outb(COM1+0, c);
7660 }
7661
7662 static int
7663 uartgetc(void)
7664 {
7665     if(!uart)
7666         return -1;
7667     if(!(inb(COM1+5) & 0x01))
7668         return -1;
7669     return inb(COM1+0);
7670 }
7671
7672 void
7673 uartintr(void)
7674 {
7675     consoleintr(uartgetc);
7676 }
7677
7678
7679
7680
7681
7682
7683
7684
7685
7686
7687
7688
7689
7690
7691
7692
7693
7694
7695
7696
7697
7698
7699

```



```
7700 # Initial process execs /init.
7701
7702 #include "syscall.h"
7703 #include "traps.h"
7704
7705
7706 # exec(init, argv)
7707 .globl start
7708 start:
7709     pushl $argv
7710     pushl $init
7711     pushl $0 // where caller pc would be
7712     movl $SYS_exec, %eax
7713     int $T_SYSCALL
7714
7715 # for(;;) exit();
7716 exit:
7717     movl $SYS_exit, %eax
7718     int $T_SYSCALL
7719     jmp exit
7720
7721 # char init[] = "/init\0";
7722 init:
7723     .string "/init\0"
7724
7725 # char *argv[] = { init, 0 };
7726 .p2align 2
7727 argv:
7728     .long init
7729     .long 0
7730
7731
7732
7733
7734
7735
7736
7737
7738
7739
7740
7741
7742
7743
7744
7745
7746
7747
7748
7749
```

```
7750 #include "syscall.h"
7751 #include "traps.h"
7752
7753 #define SYSCALL(name) \
7754     .globl name; \
7755     name: \
7756         movl $SYS_ ## name, %eax; \
7757         int $T_SYSCALL; \
7758         ret
7759
7760 SYSCALL(fork)
7761 SYSCALL(exit)
7762 SYSCALL(wait)
7763 SYSCALL(pipe)
7764 SYSCALL(read)
7765 SYSCALL(write)
7766 SYSCALL(close)
7767 SYSCALL(kill)
7768 SYSCALL(exec)
7769 SYSCALL(open)
7770 SYSCALL(mknod)
7771 SYSCALL(unlink)
7772 SYSCALL(fstat)
7773 SYSCALL(link)
7774 SYSCALL(mkdir)
7775 SYSCALL(chdir)
7776 SYSCALL(dup)
7777 SYSCALL(getpid)
7778 SYSCALL(sbrk)
7779 SYSCALL(sleep)
7780 SYSCALL(uptime)
7781
7782
7783
7784
7785
7786
7787
7788
7789
7790
7791
7792
7793
7794
7795
7796
7797
7798
7799
```

```

7800 // init: The initial user-level program
7801
7802 #include "types.h"
7803 #include "stat.h"
7804 #include "user.h"
7805 #include "fcntl.h"
7806
7807 char *argv[] = { "sh", 0 };
7808
7809 int
7810 main(void)
7811 {
7812     int pid, wpid;
7813
7814     if(open("console", O_RDWR) < 0){
7815         mknode("console", 1, 1);
7816         open("console", O_RDWR);
7817     }
7818     dup(0); // stdout
7819     dup(0); // stderr
7820
7821     for(;;){
7822         printf(1, "init: starting sh\n");
7823         pid = fork();
7824         if(pid < 0){
7825             printf(1, "init: fork failed\n");
7826             exit();
7827         }
7828         if(pid == 0){
7829             exec("sh", argv);
7830             printf(1, "init: exec sh failed\n");
7831             exit();
7832         }
7833         while((wpid=wait()) >= 0 && wpid != pid)
7834             printf(1, "zombie!\n");
7835     }
7836 }
7837
7838
7839
7840
7841
7842
7843
7844
7845
7846
7847
7848
7849

```

```

7850 // Shell.
7851
7852 #include "types.h"
7853 #include "user.h"
7854 #include "fcntl.h"
7855
7856 // Parsed command representation
7857 #define EXEC 1
7858 #define REDIR 2
7859 #define PIPE 3
7860 #define LIST 4
7861 #define BACK 5
7862
7863 #define MAXARGS 10
7864
7865 struct cmd {
7866     int type;
7867 };
7868
7869 struct execcmd {
7870     int type;
7871     char *argv[MAXARGS];
7872     char *eargv[MAXARGS];
7873 };
7874
7875 struct redircmd {
7876     int type;
7877     struct cmd *cmd;
7878     char *file;
7879     char *efile;
7880     int mode;
7881     int fd;
7882 };
7883
7884 struct pipecmd {
7885     int type;
7886     struct cmd *left;
7887     struct cmd *right;
7888 };
7889
7890 struct listcmd {
7891     int type;
7892     struct cmd *left;
7893     struct cmd *right;
7894 };
7895
7896 struct backcmd {
7897     int type;
7898     struct cmd *cmd;
7899 };

```

```

7900 int fork1(void); // Fork but panics on failure.
7901 void panic(char*);
7902 struct cmd *parsecmd(char*);
7903
7904 // Execute cmd. Never returns.
7905 void
7906 runcmd(struct cmd *cmd)
7907 {
7908     int p[2];
7909     struct backcmd *bcmd;
7910     struct execcmd *ecmd;
7911     struct listcmd *lcmd;
7912     struct pipecmd *pcmd;
7913     struct redircmd *rcmd;
7914
7915     if(cmd == 0)
7916         exit();
7917
7918     switch(cmd->type){
7919     default:
7920         panic("runcmd");
7921
7922     case EXEC:
7923         ecmd = (struct execcmd*)cmd;
7924         if(ecmd->argv[0] == 0)
7925             exit();
7926         exec(ecmd->argv[0], ecmd->argv);
7927         printf(2, "exec %s failed\n", ecmd->argv[0]);
7928         break;
7929
7930     case REDIR:
7931         rcmd = (struct redircmd*)cmd;
7932         close(rcmd->fd);
7933         if(open(rcmd->file, rcmd->mode) < 0){
7934             printf(2, "open %s failed\n", rcmd->file);
7935             exit();
7936         }
7937         runcmd(rcmd->cmd);
7938         break;
7939
7940     case LIST:
7941         lcmd = (struct listcmd*)cmd;
7942         if(fork1() == 0)
7943             runcmd(lcmd->left);
7944         wait();
7945         runcmd(lcmd->right);
7946         break;
7947
7948
7949

```

```

7950 case PIPE:
7951     pcmd = (struct pipecmd*)cmd;
7952     if(pipe(p) < 0)
7953         panic("pipe");
7954     if(fork1() == 0){
7955         close(1);
7956         dup(p[1]);
7957         close(p[0]);
7958         close(p[1]);
7959         runcmd(pcmd->left);
7960     }
7961     if(fork1() == 0){
7962         close(0);
7963         dup(p[0]);
7964         close(p[0]);
7965         close(p[1]);
7966         runcmd(pcmd->right);
7967     }
7968     close(p[0]);
7969     close(p[1]);
7970     wait();
7971     wait();
7972     break;
7973
7974 case BACK:
7975     bcmd = (struct backcmd*)cmd;
7976     if(fork1() == 0)
7977         runcmd(bcmd->cmd);
7978     break;
7979 }
7980 exit();
7981 }
7982
7983 int
7984 getcmd(char *buf, int nbuf)
7985 {
7986     printf(2, "$ ");
7987     memset(buf, 0, nbuf);
7988     gets(buf, nbuf);
7989     if(buf[0] == 0) // EOF
7990         return -1;
7991     return 0;
7992 }
7993
7994
7995
7996
7997
7998
7999

```

```

8000 int
8001 main(void)
8002 {
8003     static char buf[100];
8004     int fd;
8005
8006     // Assumes three file descriptors open.
8007     while((fd = open("console", 0_RDWR)) >= 0){
8008         if(fd >= 3){
8009             close(fd);
8010             break;
8011         }
8012     }
8013
8014     // Read and run input commands.
8015     while(getcmd(buf, sizeof(buf)) >= 0){
8016         if(buf[0] == 'c' && buf[1] == 'd' && buf[2] == ' '){
8017             // Clumsy but will have to do for now.
8018             // Chdir has no effect on the parent if run in the child.
8019             buf[strlen(buf)-1] = 0; // chop \n
8020             if(chdir(buf+3) < 0)
8021                 printf(2, "cannot cd %s\n", buf+3);
8022             continue;
8023         }
8024         if(fork1() == 0)
8025             runcmd(parsecmd(buf));
8026         wait();
8027     }
8028     exit();
8029 }
8030
8031 void
8032 panic(char *s)
8033 {
8034     printf(2, "%s\n", s);
8035     exit();
8036 }
8037
8038 int
8039 fork1(void)
8040 {
8041     int pid;
8042
8043     pid = fork();
8044     if(pid == -1)
8045         panic("fork");
8046     return pid;
8047 }
8048
8049

```

```

8050 // Constructors
8051
8052 struct cmd*
8053 execcmd(void)
8054 {
8055     struct execcmd *cmd;
8056
8057     cmd = malloc(sizeof(*cmd));
8058     memset(cmd, 0, sizeof(*cmd));
8059     cmd->type = EXEC;
8060     return (struct cmd*)cmd;
8061 }
8062
8063 struct cmd*
8064 redircmd(struct cmd *subcmd, char *file, char *efile, int mode, int fd)
8065 {
8066     struct redircmd *cmd;
8067
8068     cmd = malloc(sizeof(*cmd));
8069     memset(cmd, 0, sizeof(*cmd));
8070     cmd->type = REDIR;
8071     cmd->cmd = subcmd;
8072     cmd->file = file;
8073     cmd->efile = efile;
8074     cmd->mode = mode;
8075     cmd->fd = fd;
8076     return (struct cmd*)cmd;
8077 }
8078
8079 struct cmd*
8080 pipecmd(struct cmd *left, struct cmd *right)
8081 {
8082     struct pipecmd *cmd;
8083
8084     cmd = malloc(sizeof(*cmd));
8085     memset(cmd, 0, sizeof(*cmd));
8086     cmd->type = PIPE;
8087     cmd->left = left;
8088     cmd->right = right;
8089     return (struct cmd*)cmd;
8090 }
8091
8092
8093
8094
8095
8096
8097
8098
8099

```

```

8100 struct cmd*
8101 listcmd(struct cmd *left, struct cmd *right)
8102 {
8103     struct listcmd *cmd;
8104
8105     cmd = malloc(sizeof(*cmd));
8106     memset(cmd, 0, sizeof(*cmd));
8107     cmd->type = LIST;
8108     cmd->left = left;
8109     cmd->right = right;
8110     return (struct cmd*)cmd;
8111 }
8112
8113 struct cmd*
8114 backcmd(struct cmd *subcmd)
8115 {
8116     struct backcmd *cmd;
8117
8118     cmd = malloc(sizeof(*cmd));
8119     memset(cmd, 0, sizeof(*cmd));
8120     cmd->type = BACK;
8121     cmd->cmd = subcmd;
8122     return (struct cmd*)cmd;
8123 }
8124
8125
8126
8127
8128
8129
8130
8131
8132
8133
8134
8135
8136
8137
8138
8139
8140
8141
8142
8143
8144
8145
8146
8147
8148
8149

```

```

8150 // Parsing
8151
8152 char whitespace[] = "\t\r\n\v";
8153 char symbols[] = "<|>&;O";
8154
8155 int
8156 gettoken(char **ps, char *es, char **q, char **eq)
8157 {
8158     char *s;
8159     int ret;
8160
8161     s = *ps;
8162     while(s < es && strchr(whitespace, *s))
8163         s++;
8164     if(q)
8165         *q = s;
8166     ret = *s;
8167     switch(*s){
8168     case 0:
8169         break;
8170     case '|':
8171     case '(':
8172     case ')':
8173     case ',':
8174     case '&':
8175     case '<':
8176         s++;
8177         break;
8178     case '>':
8179         s++;
8180         if(*s == '>'){
8181             ret = '+';
8182             s++;
8183         }
8184         break;
8185     default:
8186         ret = 'a';
8187         while(s < es && !strchr(whitespace, *s) && !strchr(symbols, *s))
8188             s++;
8189         break;
8190     }
8191     if(eq)
8192         *eq = s;
8193
8194     while(s < es && strchr(whitespace, *s))
8195         s++;
8196     *ps = s;
8197     return ret;
8198 }
8199

```

```

8200 int
8201 peek(char **ps, char *es, char *toks)
8202 {
8203     char *s;
8204
8205     s = *ps;
8206     while(s < es && strchr(whitespace, *s))
8207         s++;
8208     *ps = s;
8209     return *s && strchr(toks, *s);
8210 }
8211
8212 struct cmd *parseline(char**, char*);
8213 struct cmd *parsepipe(char**, char*);
8214 struct cmd *parseexec(char**, char*);
8215 struct cmd *nulterminate(struct cmd*);
8216
8217 struct cmd*
8218 parsecmd(char *s)
8219 {
8220     char *es;
8221     struct cmd *cmd;
8222
8223     es = s + strlen(s);
8224     cmd = parseline(&s, es);
8225     peek(&s, es, "");
8226     if(s != es){
8227         printf(2, "leftovers: %s\n", s);
8228         panic("syntax");
8229     }
8230     nulterminate(cmd);
8231     return cmd;
8232 }
8233
8234 struct cmd*
8235 parseline(char **ps, char *es)
8236 {
8237     struct cmd *cmd;
8238
8239     cmd = parsepipe(ps, es);
8240     while(peek(ps, es, "&")){
8241         gettoken(ps, es, 0, 0);
8242         cmd = backcmd(cmd);
8243     }
8244     if(peek(ps, es, ";")){
8245         gettoken(ps, es, 0, 0);
8246         cmd = listcmd(cmd, parseline(ps, es));
8247     }
8248     return cmd;
8249 }

```

```

8250 struct cmd*
8251 parsepipe(char **ps, char *es)
8252 {
8253     struct cmd *cmd;
8254
8255     cmd = parseexec(ps, es);
8256     if(peek(ps, es, "|")){
8257         gettoken(ps, es, 0, 0);
8258         cmd = pipecmd(cmd, parsepipe(ps, es));
8259     }
8260     return cmd;
8261 }
8262
8263 struct cmd*
8264 parseredirs(struct cmd *cmd, char **ps, char *es)
8265 {
8266     int tok;
8267     char *q, *eq;
8268
8269     while(peek(ps, es, "<>")){
8270         tok = gettoken(ps, es, 0, 0);
8271         if(gettoken(ps, es, &q, &eq) != 'a')
8272             panic("missing file for redirection");
8273         switch(tok){
8274             case '<':
8275                 cmd = redircmd(cmd, q, eq, O_RDONLY, 0);
8276                 break;
8277             case '>':
8278                 cmd = redircmd(cmd, q, eq, O_WRONLY|O_CREATE, 1);
8279                 break;
8280             case '+': // >>
8281                 cmd = redircmd(cmd, q, eq, O_WRONLY|O_CREATE, 1);
8282                 break;
8283         }
8284     }
8285     return cmd;
8286 }
8287
8288
8289
8290
8291
8292
8293
8294
8295
8296
8297
8298
8299

```

```

8300 struct cmd*
8301 parseblock(char **ps, char *es)
8302 {
8303     struct cmd *cmd;
8304
8305     if(!peek(ps, es, "("))
8306         panic("parseblock");
8307     gettoken(ps, es, 0, 0);
8308     cmd = parseline(ps, es);
8309     if(!peek(ps, es, ")"))
8310         panic("syntax - missing )");
8311     gettoken(ps, es, 0, 0);
8312     cmd = parseredirs(cmd, ps, es);
8313     return cmd;
8314 }
8315
8316 struct cmd*
8317 parseexec(char **ps, char *es)
8318 {
8319     char *q, *eq;
8320     int tok, argc;
8321     struct execcmd *cmd;
8322     struct cmd *ret;
8323
8324     if(peek(ps, es, "("))
8325         return parseblock(ps, es);
8326
8327     ret = execcmd();
8328     cmd = (struct execcmd*)ret;
8329
8330     argc = 0;
8331     ret = parseredirs(ret, ps, es);
8332     while(!peek(ps, es, "|)&")){
8333         if((tok=gettoken(ps, es, &q, &eq)) == 0)
8334             break;
8335         if(tok != 'a')
8336             panic("syntax");
8337         cmd->argv[argc] = q;
8338         cmd->eargv[argc] = eq;
8339         argc++;
8340     }
8341     if(argc >= MAXARGS)
8342         panic("too many args");
8343     ret = parseredirs(ret, ps, es);
8344     cmd->argv[argc] = 0;
8345     cmd->eargv[argc] = 0;
8346     return ret;
8347 }
8348
8349

```

```

8350 // NUL-terminate all the counted strings.
8351 struct cmd*
8352 nulterminate(struct cmd *cmd)
8353 {
8354     int i;
8355     struct backcmd *bcmd;
8356     struct execcmd *ecmd;
8357     struct listcmd *lcmd;
8358     struct pipecmd *pcmd;
8359     struct redircmd *rcmd;
8360
8361     if(cmd == 0)
8362         return 0;
8363
8364     switch(cmd->type){
8365     case EXEC:
8366         ecmd = (struct execcmd*)cmd;
8367         for(i=0; ecmd->argv[i]; i++)
8368             *ecmd->eargv[i] = 0;
8369         break;
8370
8371     case REDIR:
8372         rcmd = (struct redircmd*)cmd;
8373         nulterminate(rcmd->cmd);
8374         *rcmd->efile = 0;
8375         break;
8376
8377     case PIPE:
8378         pcmd = (struct pipecmd*)cmd;
8379         nulterminate(pcmd->left);
8380         nulterminate(pcmd->right);
8381         break;
8382
8383     case LIST:
8384         lcmd = (struct listcmd*)cmd;
8385         nulterminate(lcmd->left);
8386         nulterminate(lcmd->right);
8387         break;
8388
8389     case BACK:
8390         bcmd = (struct backcmd*)cmd;
8391         nulterminate(bcmd->cmd);
8392         break;
8393     }
8394     return cmd;
8395 }
8396
8397
8398
8399

```

```

8400 #include "asm.h"
8401 #include "memlayout.h"
8402 #include "mmu.h"
8403
8404 # Start the first CPU: switch to 32-bit protected mode, jump into C.
8405 # The BIOS loads this code from the first sector of the hard disk into
8406 # memory at physical address 0x7c00 and starts executing in real mode
8407 # with %cs=0 %ip=7c00.
8408
8409 .code16                # Assemble for 16-bit mode
8410 .globl start
8411 start:
8412 cli
8413
8414 # Zero data segment registers DS, ES, and SS.
8415 xorw %ax,%ax          # Set %ax to zero
8416 movw %ax,%ds          # -> Data Segment
8417 movw %ax,%es          # -> Extra Segment
8418 movw %ax,%ss          # -> Stack Segment
8419
8420 # Physical address line A20 is tied to zero so that the first PCs
8421 # with 2 MB would run software that assumed 1 MB. Undo that.
8422 seta20.1:
8423 inb $0x64,%al          # Wait for not busy
8424 testb $0x2,%al
8425 jnz seta20.1
8426
8427 movb $0xd1,%al         # 0xd1 -> port 0x64
8428 outb %al,$0x64
8429
8430 seta20.2:
8431 inb $0x64,%al          # Wait for not busy
8432 testb $0x2,%al
8433 jnz seta20.2
8434
8435 movb $0xdf,%al         # 0xdf -> port 0x60
8436 outb %al,$0x60
8437
8438 # Switch from real to protected mode. Use a bootstrap GDT that makes
8439 # virtual addresses map directly to physical addresses so that the
8440 # effective memory map doesn't change during the transition.
8441 lgdt gdtdesc
8442 movl %cr0, %eax
8443 orl $CR0_PE, %eax
8444 movl %eax, %cr0
8445
8446
8447
8448
8449

```

```

8450 # Complete transition to 32-bit protected mode by using long jmp
8451 # to reload %cs and %eip. The segment descriptors are set up with no
8452 # translation, so that the mapping is still the identity mapping.
8453 jmp $(SEG_KCODE<<3), $start32
8454
8455 .code32 # Tell assembler to generate 32-bit code now.
8456 start32:
8457 # Set up the protected-mode data segment registers
8458 movw $(SEG_KDATA<<3), %ax # Our data segment selector
8459 movw %ax,%ds             # -> DS: Data Segment
8460 movw %ax,%es             # -> ES: Extra Segment
8461 movw %ax,%ss             # -> SS: Stack Segment
8462 movw $0, %ax             # Zero segments not ready for use
8463 movw %ax,%fs             # -> FS
8464 movw %ax,%gs             # -> GS
8465
8466 # Set up the stack pointer and call into C.
8467 movl $start, %esp
8468 call bootmain
8469
8470 # If bootmain returns (it shouldn't), trigger a Bochs
8471 # breakpoint if running under Bochs, then loop.
8472 movw $0x8a00, %ax        # 0x8a00 -> port 0x8a00
8473 movw %ax,%dx
8474 outw %ax,%dx
8475 movw $0x8ae0, %ax       # 0x8ae0 -> port 0x8a00
8476 outw %ax,%dx
8477 spin:
8478 jmp spin
8479
8480 # Bootstrap GDT
8481 .p2align 2              # force 4 byte alignment
8482 gdt:
8483 SEG_NULLASM            # null seg
8484 SEG_ASM(STA_X|STA_R, 0x0, 0xffffffff) # code seg
8485 SEG_ASM(STA_W, 0x0, 0xffffffff)      # data seg
8486
8487 gdtdesc:
8488 .word (gdtdesc - gdt - 1)           # sizeof(gdt) - 1
8489 .long gdt                             # address gdt
8490
8491
8492
8493
8494
8495
8496
8497
8498
8499

```



```

8500 // Boot loader.
8501 //
8502 // Part of the boot sector, along with bootasm.S, which calls bootmain().
8503 // bootasm.S has put the processor into protected 32-bit mode.
8504 // bootmain() loads an ELF kernel image from the disk starting at
8505 // sector 1 and then jumps to the kernel entry routine.
8506
8507 #include "types.h"
8508 #include "elf.h"
8509 #include "x86.h"
8510 #include "memlayout.h"
8511
8512 #define SECTSIZE 512
8513
8514 void readseg(uchar*, uint, uint);
8515
8516 void
8517 bootmain(void)
8518 {
8519     struct elfhdr *elf;
8520     struct proghdr *ph, *eph;
8521     void (*entry)(void);
8522     uchar* pa;
8523
8524     elf = (struct elfhdr*)0x10000; // scratch space
8525
8526     // Read 1st page off disk
8527     readseg((uchar*)elf, 4096, 0);
8528
8529     // Is this an ELF executable?
8530     if(elf->magic != ELF_MAGIC)
8531         return; // let bootasm.S handle error
8532
8533     // Load each program segment (ignores ph flags).
8534     ph = (struct proghdr*)((uchar*)elf + elf->phoff);
8535     eph = ph + elf->phnum;
8536     for(; ph < eph; ph++){
8537         pa = (uchar*)ph->paddr;
8538         readseg(pa, ph->filesz, ph->off);
8539         if(ph->memsz > ph->filesz)
8540             stosb(pa + ph->filesz, 0, ph->memsz - ph->filesz);
8541     }
8542
8543     // Call the entry point from the ELF header.
8544     // Does not return!
8545     entry = (void(*) (void))(elf->entry);
8546     entry();
8547 }
8548
8549

```

```

8550 void
8551 waitdisk(void)
8552 {
8553     // Wait for disk ready.
8554     while((inb(0x1F7) & 0xC0) != 0x40)
8555         ;
8556 }
8557
8558 // Read a single sector at offset into dst.
8559 void
8560 readsect(void *dst, uint offset)
8561 {
8562     // Issue command.
8563     waitdisk();
8564     outb(0x1F2, 1); // count = 1
8565     outb(0x1F3, offset);
8566     outb(0x1F4, offset >> 8);
8567     outb(0x1F5, offset >> 16);
8568     outb(0x1F6, (offset >> 24) | 0xE0);
8569     outb(0x1F7, 0x20); // cmd 0x20 - read sectors
8570
8571     // Read data.
8572     waitdisk();
8573     insl(0x1F0, dst, SECTSIZE/4);
8574 }
8575
8576 // Read 'count' bytes at 'offset' from kernel into physical address 'pa'.
8577 // Might copy more than asked.
8578 void
8579 readseg(uchar* pa, uint count, uint offset)
8580 {
8581     uchar* epa;
8582
8583     epa = pa + count;
8584
8585     // Round down to sector boundary.
8586     pa -= offset % SECTSIZE;
8587
8588     // Translate from bytes to sectors; kernel starts at sector 1.
8589     offset = (offset / SECTSIZE) + 1;
8590
8591     // If this is too slow, we could read lots of sectors at a time.
8592     // We'd write more to memory than asked, but it doesn't matter --
8593     // we load in increasing order.
8594     for(; pa < epa; pa += SECTSIZE, offset++)
8595         readsect(pa, offset);
8596 }
8597
8598
8599

```