

项目说明文档

操作系统课程设计

——xv6 及 Labs 课程项目

作者姓名：_____杜天乐_____

学 号：_____2251310_____

指导教师：_____王冬青_____

学院专业：_____软件学院 软件工程_____

同济大学

Tongji University

目录

项目说明文档.....	1
操作系统课程设计.....	1
Tongji University.....	1
0. 实验准备.....	6
0.1 安装 VMWare 虚拟机.....	6
0.2 在 VMWare 中安装 Ubuntu 系统.....	6
0.3 在 Ubuntu 系统中配置 xv6 环境.....	6
1. Lab1: Xv6 and Unix utilities.....	8
1.1 Boot xv6.....	8
1.1.1 实验目的.....	8
1.1.2 实验步骤.....	8
1.1.3 实验中遇到的问题和解决办法.....	9
1.1.4 实验心得.....	9
1.2 sleep.....	10
1.2.1 实验目的.....	10
1.2.2 实验步骤.....	10
1.2.3 实验中遇到的问题和解决办法.....	12
1.2.4 实验心得.....	12
1.3 pingpong.....	12
1.3.1 实验目的.....	12
1.3.2 实验步骤.....	13
1.3.3 实验中遇到的问题和解决办法.....	14
1.3.4 实验心得.....	15
1.4 primes.....	15
1.4.1 实验目的.....	15
1.4.2 实验步骤.....	15
1.4.3 实验中遇到的问题和解决办法.....	18
1.4.4 实验心得.....	19
1.5 find.....	19
1.5.1 实验目的.....	19
1.5.2 实验步骤.....	19
1.5.3 实验中遇到的问题和解决办法.....	24
1.5.4 实验心得.....	24
1.6 xargs.....	24
1.6.1 实验目的.....	24
1.6.2 实验步骤.....	25
1.6.3 实验中遇到的问题和解决办法.....	27
1.6.4 实验心得.....	27
1.7 Lab1 实验成绩.....	27
2. Lab2: system calls.....	29
2.1 System call tracing.....	29
2.1.1 实验目的.....	29
2.1.2 实验步骤.....	29

2.1.3 实验中遇到的问题和解决办法	33
2.1.4 实验心得	34
2.2 Sysinfo	34
2.2.1 实验目的	34
2.2.2 实验步骤	34
2.2.3 实验中遇到的问题和解决办法	39
2.2.4 实验心得	40
2.3 Lab2 实验成绩	40
3. Lab3: page tables	41
3.1 Speed up system calls	41
3.1.1 实验目的	41
3.1.2 实验步骤	41
3.1.3 实验中遇到的问题和解决办法	43
3.1.4 实验心得	43
3.2 Print a page table	44
3.2.1 实验目的	44
3.2.2 实验步骤	44
3.2.3 实验中遇到的问题和解决办法	46
3.2.4 实验心得	47
3.3 Detecting which pages have been accessed	47
3.3.1 实验目的	47
3.3.2 实验步骤	47
3.3.3 实验中遇到的问题和解决办法	49
3.3.4 实验心得	50
3.4 Lab3 实验成绩	50
4. Lab4: page tables	51
4.1 RISC-V assembly	51
4.1.1 实验目的	51
4.1.2 实验步骤	51
4.1.3 实验中遇到的问题和解决办法	55
4.1.4 实验心得	55
4.2 Backtrace (moderate)	55
4.2.1 实验目的	55
4.2.2 实验步骤	55
4.2.3 实验中遇到的问题和解决办法	57
4.2.4 实验心得	58
4.3 Alarm	58
4.3.1 实验目的	58
4.3.2 实验步骤	58
4.3.3 实验中遇到的问题和解决办法	63
4.3.4 实验心得	64
4.4 Lab4 实验成绩	64
5. Lab5: Copy-on-Write Fork for xv6	65
5.1 Implement copy-on write	65

5.1.1 实验目的	65
5.1.2 实验步骤	65
5.1.3 实验中遇到的问题和解决办法	73
5.1.4 实验心得	73
5.2 Lab5 实验成绩	73
6. Lab6: Multithreading.....	75
6.1 Uthread: switching between threads	75
6.1.1 实验目的	75
6.1.2 实验步骤	75
6.1.3 实验中遇到的问题和解决办法	78
6.1.4 实验心得	79
6.2 Using threads.....	79
6.2.1 实验目的	79
6.2.2 实验步骤	79
6.2.3 实验中遇到的问题和解决办法	81
6.2.4 实验心得	82
6.3 Barrier.....	82
6.3.1 实验目的	82
6.3.2 实验步骤	82
6.3.3 实验中遇到的问题和解决办法	84
6.3.4 实验心得	85
6.4 Lab6 实验成绩	85
7. Lab7: Networking.....	87
7.1 Your Job	87
7.1.1 实验目的	87
7.1.2 实验步骤	87
7.1.3 实验中遇到的问题和解决办法	91
7.1.4 实验心得	92
7.2 Lab7 实验成绩	92
8. Lab8: Locks.....	94
8.1 Memory allocator.....	94
8.1.1 实验目的	94
8.1.2 实验步骤	94
8.1.3 实验中遇到的问题和解决办法	100
8.1.4 实验心得	101
8.2 Buffer cache.....	101
8.2.1 实验目的	101
8.2.2 实验步骤	101
8.2.3 实验中遇到的问题和解决办法	106
8.2.4 实验心得	107
8.3 Lab8 实验成绩	107
9. Lab9: File system.....	108
9.1 Large files.....	108
9.1.1 实验目的	108

9.1.2 实验步骤	108
9.1.3 实验中遇到的问题和解决办法	115
9.1.4 实验心得	116
9.2 Symbolic links	116
9.2.1 实验目的	116
9.2.2 实验步骤	116
9.2.3 实验中遇到的问题和解决办法	120
9.2.4 实验心得	121
9.3 Lab9 实验成绩	121
10. Lab10: mmap	122
10.1 mmap	122
10.1.1 实验目的	122
10.1.2 实验步骤	122
10.1.3 实验中遇到的问题和解决办法	134
10.1.4 实验心得	134
10.2 Lab10 实验成绩	135

0. 实验准备

0.1 安装 VMWare 虚拟机

访问 (<https://www.virtualbox.org/wiki/Downloads>)，使用默认设置进行安装，并重启计算机。

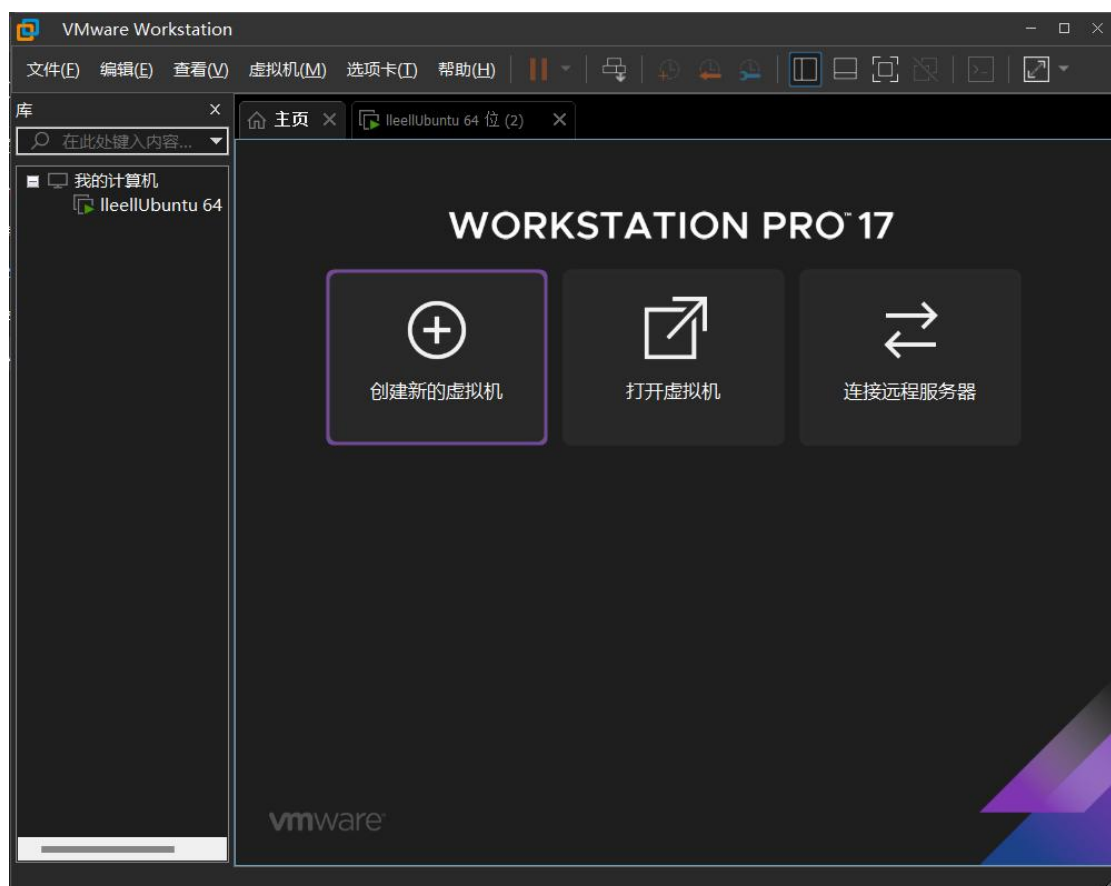
0.2 在 VMWare 中安装 Ubuntu 系统

(1) 下载 Ubuntu 镜像：

从 Ubuntu 官网或者从国内一些高校的镜像站中进行下载。我使用的是清华大学镜像站 (<https://mirrors.tuna.tsinghua.edu.cn/ubuntu-releases/>)。选择合适的版本，我使用的是 20.04 的版本，不同版本的 Ubuntu 对后续实验有影响（我第一次选择的是 24 版的，需要进行额外步骤，换成 20.04 的版本即可）

(2) 在 VMWare 中安装 Ubuntu 系统：

打开 VMWare，在主页点击“创建新的虚拟机”：



根据指示和自己的具体需求，结合网上教程，配置虚拟机，我参考的是 CSDN 中的文章 (https://blog.csdn.net/sherry__yuzu/article/details/140409245)

0.3 在 Ubuntu 系统中配置 xv6 环境

(1) 首先需要认识：

①sudo：系统级权限，不带这个很多命令执行不了

②apt-get：安装软件、更新操作等都通过改命令来执行

(2) 打开 Ubuntu 系统的终端，安装本项目所需的所有软件，运行：

```
$ sudo apt-get update && sudo apt-get upgrade
```

```
$ sudo apt-get install git build-essential gdb-multiarch qemu-system-misc  
gcc-riscv64-linux-gnu binutils-riscv64-linux-gnu
```

(3) 测试安装：

```
$ qemu-system-riscv64 --version
```

```
$ riscv64-linux-gnu-gcc --version
```

得到测试结果如下图所示，表示安装完成：

```
l1eell@l1eell-VMware-Virtual-Platform:~/Desktop$ qemu-system-riscv64 --version  
QEMU emulator version 8.2.2 (Debian 1:8.2.2+ds-0ubuntu1)  
Copyright (c) 2003-2023 Fabrice Bellard and the QEMU Project developers  
l1eell@l1eell-VMware-Virtual-Platform:~/Desktop$ riscv64-linux-gnu-gcc --version  
riscv64-linux-gnu-gcc (Ubuntu 13.2.0-23ubuntu4) 13.2.0  
Copyright (C) 2023 Free Software Foundation, Inc.  
This is free software; see the source for copying conditions. There is NO  
warranty; not even for MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.
```

(4) 获取实验室的 xv6 源代码（注意设置 VPN，挂梯子）：

```
$ git clone git://g.csail.mit.edu/xv6-labs-2021
```

```
l1eell@l1eell-VMware-Virtual-Platform:~/桌面$ git clone git://g.csail.mit.edu/xv  
6-labs-2021  
正克隆到 'xv6-labs-2021'...  
remote: Enumerating objects: 7051, done.  
remote: Counting objects: 100% (7051/7051), done.  
remote: Compressing objects: 100% (3423/3423), done.  
remote: Total 7051 (delta 3702), reused 6830 (delta 3600), pack-reused 0  
接收对象中: 100% (7051/7051), 17.20 MiB | 1.97 MiB/s, 完成。  
处理 delta 中: 100% (3702/3702), 完成。  
warning: remote HEAD refers to nonexistent ref, unable to checkout
```

Github 地址：<https://github.com/l1eell104/OS-Xv6-Lab-2024.git>

1. Lab1: Xv6 and Unix utilities

1.1 Boot xv6

1.1.1 实验目的

启动 xv6 操作系统内核，以初步掌握其基本功能与使用方法。通过本实验，了解 xv6 操作系统的工作原理，以及如何使用简单的命令进行交互。

1.1.2 实验步骤

(1) 切换到 util 分支

进入克隆的仓库目录，并切换到 util 分支。

```
cd xv6-labs-2021  
  
git checkout util  
  
make clean
```

```
lleell@lleell-VMware-Virtual-Platform:~/桌面$ cd xv6-labs-2021  
lleell@lleell-VMware-Virtual-Platform:~/桌面/xv6-labs-2021$ git checkout util  
分支 'util' 设置为跟踪 'origin/util'。  
切换到一个新分支 'util'
```

(2) 构建并运行 xv6

在确保 Ubuntu 版本为 20.04 的情况下，使用 `make qemu` 命令构建并运行 xv6。xv6 启动后，init 进程会启动一个 shell 等待用户的命令。

```
make qemu
```

```
xv6 kernel is booting  
  
hart 2 starting  
hart 1 starting  
init: starting sh  
$ █
```

(3) 验证 xv6 运行状态

输入 `ls` 命令以验证 xv6 是否正确加载初始文件系统，并显示其中的内容。


```

$ ls
.          1 1 1024
..         1 1 1024
README    2 2 2226
xargstest.sh 2 3 93
cat        2 4 23904
echo       2 5 22736
forktest   2 6 13096
grep       2 7 27264
init       2 8 23840
kill       2 9 22704
ln         2 10 22664
ls         2 11 26136
mkdir      2 12 22808
rm         2 13 22800
sh         2 14 41672
stressfs   2 15 23808
usertests  2 16 156024
grind      2 17 37976
wc         2 18 25048
zombie     2 19 22200
console    3 20 0
$

```

(4) 检查进程信息

由于 xv6 没有内置的 `ps` 命令, 使用快捷键 `Ctrl-p` 来显示当前运行的进程信息。通常可以看到两个进程: `init` 和 `sh` (shell)。

(5) 退出 qemu 环境

输入 `ctrl+A`, 再输入 `X`, 退出 qemu 环境

```

$ QEMU: Terminated
lleeell@ubuntu:~/Desktop/xv6-labs-2021$

```

1.1.3 实验中遇到的问题和解决办法

(1) 问题 1: 安装 QEMU 模拟器的时候出现错误。

解决办法: 通过查阅官方文档或在线教程了解 QEMU 的安装过程。询问同学或教师有关启动项目的具体步骤。

(2) 问题 2: Ubuntu 版本不兼容。

解决办法: 更换到支持的 Ubuntu 版本, 之前 Ubuntu 的版本是 24, 无法正常实验, 更换至 20.04 版本。

1.1.4 实验心得

通过本次实验, 我对 xv6 操作系统有了初步的认识。不仅了解了如何使用 Git 管理源代码, 还掌握了如何在虚拟环境中构建和运行 xv6 操作系统。此外, 通过实际操作, 我对操作系统的基本概念有了更直观的理解, 比如文件系统的加载、进程的概念等。此次实验也让我意识到了理论知识与实践之间的差距, 未来还需要继续努力学习, 提高自己的动手能力。

1.2 sleep

1.2.1 实验目的

本实验旨在实现一个 Unix 风格的 sleep 工具，该工具能够根据用户指定的时间单位 (tick) 来暂停进程。通过这个实验，我们学习了如何获取和解析命令行参数、使用系统调用来暂停进程、在用户空间中编写和调试程序以及了解如何使自编写的程序能够被 xv6 操作系统编译和运行。此外，实验还将加深我们对系统调用机制的理解，并提高我们在有限资源环境下进行编程的能力。

1.2.2 实验步骤

(1) 设计 sleep 程序

观察 user 中的文件，熟悉 xv6 用户程序的编程风格，注意到它与标准 C 语言程序的不同之处，例如使用的头文件和程序退出方式。

在 user/user.h 中定义了 `int sleep(int)` sleep 接受参数的个数为 1 个。

在 user/ 目录下创建 user/sleep.c。



sleep.c

设计一个 sleep.c 的命令行接口，以确保程序能够正确地读取和解释用户输入的时间值。

使用 `atoi()` 函数将命令行参数转换为整数。

调用 `sleep()` 系统调用来暂停进程指定的时间单位。设计具体代码如下：

```
#include "kernel/types.h" // 包含内核定义的基本数据类型

#include "kernel/stat.h" // 包含文件状态相关的定义

#include "user/user.h" // 包含用户态程序所需的各种系统调用和函数原型 int

// main 函数，程序的入口点，argc 表示命令行参数的数量，argv 是一个指向字符串数组的指针

main(int argc, char *argv[])

{

    // 检查传递给程序的参数数量是否少于 2。如果少于 2，说明用户没有提供需要的时间参数

    if(argc < 2){
```

```
fprintf(2,"Usage:sleep [time]\n");// 向标准错误输出（文件描述符 2）打印一条使用说明信息
```

```
exit(1);// 状态码 1 表示错误
```

```
}
```

```
int time = atoi(argv[1]);// 将传递的第一个参数（argv[1]）转换为整数，并存储在变量 time 中
```

```
sleep(time);// 调用 sleep 函数，使程序暂停执行 time 秒
```

```
exit(0);// 状态码 0 表示正确
```

```
}
```

(2) 加入 Makefile

打开 Makefile 并找到 UPROGS 变量。

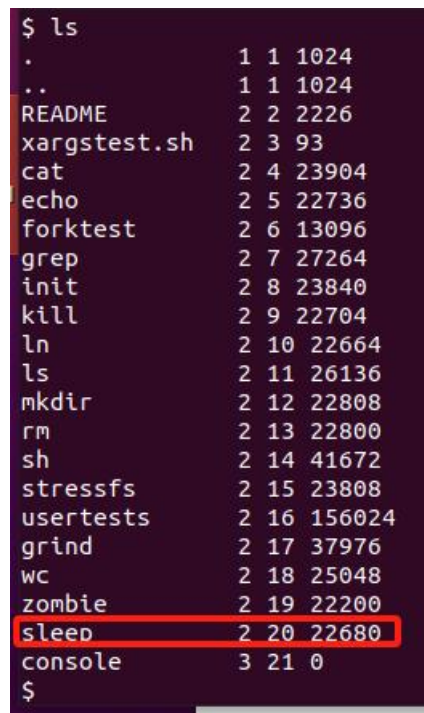
添加 sleep 到 UPROGS 变量中，确保 xv6 构建系统能够编译它。

```
$U/_sleep\
```

(3) 构建和测试

在终端中执行 make qemu 命令，进入 xv6 的 QEMU 环境。

进入 xv6 后输入 ls，出现下图表示程序已经在文件系统的根目录下了。



```
$ ls
.          1 1 1024
..         1 1 1024
README    2 2 2226
xargstest.sh 2 3 93
cat       2 4 23904
echo      2 5 22736
forktest  2 6 13096
grep      2 7 27264
init      2 8 23840
kill      2 9 22704
ln        2 10 22664
ls        2 11 26136
mkdir     2 12 22808
rm        2 13 22800
sh        2 14 41672
stressfs  2 15 23808
usertests 2 16 156024
grind     2 17 37976
wc        2 18 25048
zombie    2 19 22200
sleep     2 20 22680
console   3 21 0
$
```

(4) 性能评估

使用 xv6 自带的评估脚本来验证程序是否正确工作，在终端里输入 ./grade-lab-util sleep

评估结果如下，测试通过，程序的行为符合预期。

```
llee11@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-util sleep
make: 'kernel/kernel' is up to date.
== Test sleep, no arguments == sleep, no arguments: OK (1.6s)
== Test sleep, returns == sleep, returns: OK (0.9s)
== Test sleep, makes syscall == sleep, makes syscall: OK (1.1s)
```

1.2.3 实验中遇到的问题和解决办法

(1) 问题 1: 编译过程中出现了未定义的符号错误。

解决方法: 为了排除这个问题, 我检查了所有头文件的包含指令, 并确认 `Makefile` 中正确地包含了 `sleep` 程序。通过这样的步骤, 我解决了链接阶段的问题, 并确保程序能够成功编译。

(2) 问题 2: 如何将字符串参数转换为整数。

解决方法: 使用了 `user/user.h` 中提供的 `atoi()` 函数来进行转换。

(3) 问题 3: 如何正确退出程序。

解决方法: 根据 `user/echo.c` 的示例, 我发现 `xv6` 用户程序通常使用 `exit()` 而不是 `return` 来结束进程。因此, 在我的 `sleep` 程序中, 我也采用了同样的方式来退出。

1.2.4 实验心得

通过本次实验, 我对操作系统的工作原理有了更深入的理解。特别是在用户空间与内核空间之间的交互方面, 我学会了如何利用系统调用来实现特定的功能。此外, 我也熟悉了 `xv6` 的构建流程和调试技巧。这些经验对于未来学习更复杂的操作系统概念非常有帮助。具体表现为:

系统调用的重要性: 系统调用是用户程序与操作系统之间沟通的桥梁, 理解它们的工作机制对于编写高效的应用程序至关重要。

C 语言编程的灵活性: 虽然 `xv6` 的用户程序与标准 C 语言略有不同, 但这种差异增强了我对 C 语言特性的理解和掌握。

调试技巧: 学会如何有效地调试 `xv6` 应用程序对于快速解决问题非常有用。

1.3 pingpong

1.3.1 实验目的

本实验旨在通过编写一个 `pingpong` 程序来验证 `xv6` 操作系统中进程间通信机制的有效性。具体而言, 程序需要创建一个子进程, 并使用管道与子进程进行通信。父进程先向子进程发送一个字节的数据 (通常为 "ping"), 子进程接收后输出 `<pid>: received ping` (其中 `<pid>` 是子进程的进程号), 随后子进程向父进程发送一个字节的数据 (通常是 "pong"), 父进程接收后输出 `<pid>: received pong` (其中 `<pid>` 是父进程的进程号)。

1.3.2 实验步骤

(1) 创建 pingpong 文件

在 user/ 目录下创建 pingpong.c 文件。



在 Makefile 中添加 pingpong 到 UPROGS 变量中，确保它会被编译。

```
$U/_pingpong\
```

(2) 设计 pingpong 程序

使用 pipe() 系统调用来创建两个管道，一个用于父进程到子进程的通信，另一个用于子进程到父进程的通信。

使用 fork() 创建子进程。fork 创建子进程时，会复制父进程的寄存器、内存空间和进程控制块，并且将子进程的进程控制块中的父进程指针指向原父进程。此时，父子进程几乎完全相同。为了区分父进程和子进程，fork 给它们返回不同的值：父进程返回子进程的 PID，而子进程返回 0。

在子进程中，使用 read() 从管道中读取数据，并使用 printf() 打印接收信息。

在子进程中，使用 write() 向管道中写入数据。

在父进程中，使用 write() 向管道中写入数据。

在父进程中，使用 wait() 等待子进程结束，并使用 read() 从管道中读取数据。

使用 getpid() 获取当前进程的进程号。

设计代码如下：

```
#include "kernel/types.h"// 包含内核类型定义
#include "kernel/stat.h"// 包含内核类型定义
#include "user/user.h" // 包含用户空间函数声明

int
main(int argc, char *argv[])
{
    int fd1[2]; // 定义管道 fd1
    int fd2[2]; // 定义管道 fd2
    pipe(fd1); // 创建管道 fd1
```

```

pipe(fd2); // 创建管道 fd2

char mes[81];

if(fork()== 0){
    // 子进程

    read(fd1[0],mes,4); // 从管道 fd1 读取最多 4 字节数据
    printf("%d:received %s\n",getpid(),mes); // 打印接收到的数据
    write(fd2[1],"pong",strlen("pong")); // 向管道 fd2 写入字符串 "pong"
}
else {
    // 父进程

    write(fd1[1],"ping", strlen("ping")); // 向管道 fd1 写入字符串 "ping"
    wait(0); // 等待子进程结束

    read(fd2[0],mes,4); // 从管道 fd2 读取最多 4 字节数据
    printf("%d: received %s\n",getpid(),mes); // 打印接收到的数据
}

exit(0); // 结束进程
}

```

(3) 编译与运行

执行 `make qemu` 编译并启动 xv6 模拟器。在 xv6 的 shell 中输入 `pingpong` 来运行程序。观察程序输出，确保符合预期。

```

$ pingpong
4: received ping
3: received pong
$

```

(4) 评估与测试

使用实验自带的测评工具 `./grade-lab-util pingpong` 对程序进行自动测试。评估结果如下，测试通过，程序的行为符合预期。

```

llee11@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-util pingpong
make: 'kernel/kernel' is up to date.
== Test pingpong == pingpong: OK (1.7s)

```

1.3.3 实验中遇到的问题和解决办法

(1) 问题 1：父进程在子进程未完成之前就提前退出。

解决方法：在父进程中调用 `wait()` 系统调用来等待子进程的结束，这可以避免父进程过早退出，保证子进程能够正常完成其任务。

1.3.4 实验心得

通过这次实验，我对进程间通信有了更深刻的理解，特别是利用管道来进行父子进程间的数据交换。实验中，我不仅熟悉了 `pipe()`, `fork()`, `read()`, `write()` 和 `getpid()` 等系统调用的具体应用，还学会了如何通过这些调用来构建一个 pingpong 程序。在调试过程中，我也逐渐掌握了如何有效地使用 `printf()` 来输出调试信息，这对于定位程序中的错误非常有帮助。此外，我还学会了如何设计测试案例，并使用自动化测试工具来确保程序的正确性和健壮性。整个实验过程不仅提高了我的 C 语言编程能力，也增强了我对操作系统底层机制的认识。

1.4 primes

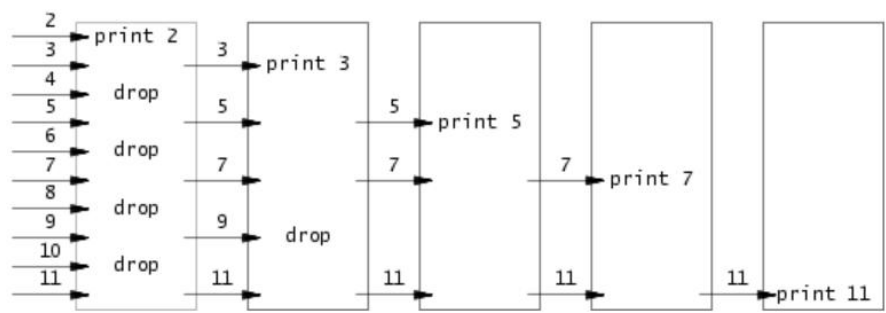
1.4.1 实验目的

本实验旨在通过使用管道和进程间的通信机制来筛选出指定范围内的所有素数。利用 xv6 操作系统环境,通过 `fork` 和 `pipe` 系统调用来创建一系列子进程，每个子进程负责筛选出特定因子的素数。实验的主要目标是理解进程间通信的机制，并学习如何高效地利用这些机制来解决问题。

1.4.2 实验步骤

(1) 理解质数筛算法：

主要思想是让每个进程负责检查一个质数是否为输入的因子。当发现一个新质数时，会创建一个新的进程来检验该质因子。进程之间通过管道传递数据。最初的父进程生成从 2 到 `maxn` 的整数，然后生成一个子进程来检查质因子 2。第一个无法被质因子整除的数用于创建下一个子进程，并输出该数，而质数筛的特性保证该数是质数。



完成数据传递或更新时，需要及时关闭一个进程不需要的文件描述符（防止程序在父进程到达 35 之前耗尽 xv6 的资源）父进程需要等待子进程的结束，并

回收共享的资源和数据等，即一旦第一个进程到达 35，它应该等待直到整个管道终止，包括所有子进程、孙进程等。

(2) 创建 primes 文件

在 user 目录下创建名为 primes.c 的新文件。



将 prime 程序添加到 MAKEFILE 文件的 UPROGS 部分。

```
$U/_primes\
```

(3) 设计 primes 程序

编写 prime 函数，该函数接收一个管道读端作为参数，并从中读取数字。如果读取到的是素数，则输出该素数，并创建一个新的管道和子进程继续筛选素数。

在 main 函数中，创建一个管道，并通过 fork 创建子进程。父进程负责向管道写入从 2 到 35 的所有数字，而子进程则调用 prime 函数开始筛选素数的过程。

在适当的位置关闭管道的读端或写端，以避免资源泄漏。

```
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"

void prime(int rd) {
    int n;
    read(rd, &n, 4); // 把读取的内容放进 n
    printf("prime %d\n", n); // 输出当前的 n
    int created = 0; // 是否已创建管道
    int p[2];
    int num;

    while (read(rd, &num, 4) != 0) { // 用 num 存管道里在 n 之后的数
        if (created == 0) {
            pipe(p); // 创建对应于 n 的管道
            created = 1;
            int pid = fork(); // 创建子进程
            if (pid == 0) { // 子进程
```



```

        close(p[1]);
        prime(p[0]); // 递归，判断子管道里的内容是不是质数
        return;
    } else { // 当前进程
        close(p[0]); // 关闭读取，允许写入
    }
}

if (num % n != 0) { // 如果 num 不是 n 的倍数，则它有可能是质数
    write(p[1], &num, 4); // 把它写入子管道
}
}

close(rd);
if (created) {
    close(p[1]);
    wait(0); // 等待子进程结束
}
}

int main(int argc, char *argv[]) {
    int p[2];
    pipe(p); // 创建管道 p
    int pid = fork();
    if (pid != 0) { // 父进程
        close(p[0]); // 关闭 p 的读取。只有读取端关闭，才能进行写入
        for (int i = 2; i <= 35; i++) {
            write(p[1], &i, 4); // 写入
        }
        close(p[1]); // 关闭 p 的写入
        wait(0); // 等待子进程结束
    }
}

```

```

    } else { // 子进程

        close(p[1]);

        prime(p[0]);

        close(p[0]);

    }

    exit(0);

}

```

(4) 编译与运行

使用 `make qemu` 编译并运行程序。

进入 `xv6` 后，输入 `primes` 命令来执行程序。

```

$ primes
prime 2
prime 3
prime 5
prime 7
prime 11
prime 13
prime 17
prime 19
prime 23
prime 29
prime 31

```

(5) 评估与测试

使用 `xv6` 自带的测评工具 `./grade-lab-util primes` 来进行自动测试。评估结果如下，测试通过，程序的行为符合预期。

```

lleell@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-util primes
make: 'kernel/kernel' is up to date.
== Test primes == primes: OK (1.1s)
lleell@ubuntu:~/Desktop/xv6-labs-2021$

```

(6) 分析结果

观察程序输出，确保所有素数都被正确筛选出来。检查程序是否能够正确处理进程间的同步问题，并且所有进程都能顺利退出。

1.4.3 实验中遇到的问题和解决办法

(1) 问题 1：在实验初期，由于没有及时关闭进程不需要的文件描述符，导致程序在第一个进程达到 35 之前就因资源不足而失败。

解决方案：在每个进程不再需要管道的一端时立即关闭相应的文件描述符，例如，在父进程中关闭管道的读端，在子进程中关闭管道的写端。

(2) 问题 2：在程序设计中，如何确保父进程在所有子进程都完成任务之后才退出，成为了一个难点。

解决方案：使用 `wait` 系统调用来等待所有子进程的结束，这确保了主进程在所有输出打印完毕并且所有子进程退出之后才退出。

1.4.4 实验心得

通过本次实验，我对进程间通信有了更深入的理解。特别是对 `fork` 和 `pipe` 系统调用的使用有了实践经验。我发现合理地管理进程资源非常重要，不仅能够提高程序的性能，还能避免资源泄露等问题。此外，我还学会了如何在复杂的程序结构中进行调试和测试，这对于确保程序的稳定性和可靠性至关重要。最后，实验中遇到的问题让我意识到，在实际编程中考虑各种边界情况是非常重要的，这样可以避免很多潜在的问题。

1.5 find

1.5.1 实验目的

本实验旨在实现一个简化版的 `find` 工具，该工具能够在指定的目录树中查找具有特定名称的所有文件。通过这个实验，我们能够更好地理解文件系统的操作，如文件路径的处理、递归搜索目录结构以及字符串匹配算法。

1.5.2 实验步骤

(1) 分析需求, 创建 `find` 文件:

在 `user` 目录下创建名为 `find.c` 的新文件。



将 `find` 程序添加到 `MAKEFILE` 文件的 `UPROGS` 部分。

```
$U/_find\
```

(2) 设计 `find` 程序:

主函数 (`main`): 负责接收用户输入的路径和文件名, 验证参数的有效性, 并启动查找过程。

查找函数 (`find`): 递归地搜索目录树中的文件, 调用 `match` 函数来判断文件名是否匹配。打开目录, 获取文件状态, 判断是文件还是目录, 如果是文件则调用 `match` 函数; 如果是目录, 则递归调用 `find` 函数。

匹配函数 (`match`): 用于比较给定路径下的文件名是否与指定的文件名相匹配。从路径中提取文件名, 并与目标文件名进行比较。

格式化函数 (fmtname)：格式化文件名，确保文件名的长度符合要求。如果文件名长度小于 DIRSIZ，则用空格填充至 DIRSIZ 长度；如果文件名长度大于或等于 DIRSIZ，则直接返回文件名。

```
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"
#include "kernel/fs.h"

// 格式化函数: 格式化文件路径中的文件名
char* fmtname(char* path) {
    // 定义一个静态的缓冲区，大小为 DIRSIZ+1，确保能够存放文件名
    static char buf[DIRSIZ + 1];

    char* p;
    // 查找最后一个 '/' 后的第一个字符
    for (p = path + strlen(path); p >= path && *p != '/'; p--);
    p++;
    // 如果文件名长度大于或等于 DIRSIZ，直接返回文件名
    if (strlen(p) >= DIRSIZ)
        return p;
    // 否则，将文件名拷贝到缓冲区 buf，并用空格填充至 DIRSIZ 长度
    memmove(buf, p, strlen(p));
    memset(buf + strlen(p), ' ', DIRSIZ - strlen(p));
    // 返回格式化后的文件名
    return buf;
}

// 匹配函数: 判断路径中的文件名是否与给定的文件名匹配
int match(char* path, char* name) {
    char* p;
    // 查找最后一个 '/' 后的第一个字符，也就是文件名的起始位置
```

```

    for (p = path + strlen(path); p >= path && *p != '/'; p--);
    p++; // 将指针移到文件名的第一个字符
    // 比较文件名是否与给定的文件名匹配
    if (strcmp(p, name) == 0)
        return 1; // 匹配成功，返回 1
    else
        return 0; // 匹配失败，返回 0
}

// 查找函数: 递归地在目录树中查找文件
void find(char* path, char* name) {
    char buf[512], *p;
    int fd;
    struct dirent de;
    struct stat st;
    // 打开目录
    if ((fd = open(path, 0)) < 0) {
// 如果打开失败，输出错误信息并返回
        fprintf(2, "find: cannot open %s\n", path);
        return;
    }
    // 获取文件状态
    if (fstat(fd, &st) < 0) {
        // 如果获取失败，输出错误信息并关闭文件描述符
        fprintf(2, "find: cannot stat %s\n", path);
        close(fd);
        return;
    }
    // 判断当前路径对应的是文件还是目录

```

```

    if (st.type == T_FILE) {
        // 如果是文件，调用 match 函数检查文件名是否匹配
        if (match(path, name)) {
            printf("%s\n", path); // 如果匹配，输出文件路径
        }
    } else if (st.type == T_DIR) {
        // 如果是目录，准备进行递归查找

// 检查路径长度，防止缓冲区溢出
        if (strlen(path) + 1 + DIRSIZ + 1 > sizeof(buf)) {
            printf("find: path too long\n");
            close(fd);
            return;
        }

// 将当前路径复制到缓冲区，并在末尾添加 '/'
        strcpy(buf, path);
        p = buf + strlen(buf);
        *p++ = '/';

// 读取目录中的每个条目
        while (read(fd, &de, sizeof(de)) == sizeof(de)) {
            // 跳过无效目录项
            if (de.inum == 0) continue;

            // 跳过"."和".."目录
            if (de.name[0] == '.' && (de.name[1] == '\0' || (de.name[1] == '.'
&& de.name[2] == '\0'))) continue;

            // 将目录项的名称添加到路径缓冲区
            memmove(p, de.name, DIRSIZ);
            p[DIRSIZ] = '\0'; // 确保路径字符串以 '\0' 结尾

            // 获取新路径的状态信息

```

```

        if (stat(buf, &st) < 0) {
            printf("find: cannot stat %s\n", buf);
            continue;
        }
        // 递归地查找子目录或文件
        find(buf, name);
    }
}

// 关闭目录文件描述符
close(fd);
}

// 主函数: 负责接收用户输入的路径和文件名, 并启动查找过程
int main(int argc, char* argv[]) {
    // 检查命令行参数的数量, 确保用户提供了路径和文件名
    if (argc < 3) {
        printf("Usage: find [path] [filename]\n");
        exit(-1); // 参数不足, 退出程序
    }
    // 调用 find 函数, 开始从指定路径查找指定文件名
    find(argv[1], argv[2]);
    exit(0); // 程序执行结束, 正常退出
}

```

(3) 编译与运行:

使用 `make clean` 命令清理文件系统, 然后创建一些测试文件和目录。
 执行 `find` 命令, 并检查输出结果是否正确。

```

init: starting sh
$ echo>b
$ mkdir a
$ echo>a/b
$ find.b
exec find.b failed
$ find . b
./b
./a/b
$

```

(4) 评估与测试:

使用 xv6 自带的测评工具 `./grade-lab-util find` 来进行自动测试。评估结果如下，测试通过，程序的行为符合预期。

```

lleell@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-util find
make: 'kernel/kernel' is up to date.
== Test find, in current directory == find, in current directory: OK (1.2s)
== Test find, recursive == find, recursive: OK (1.0s)

```

1.5.3 实验中遇到的问题和解决办法

(1) 问题 1: 在递归调用 `find` 函数时，路径长度可能会超出缓冲区大小，导致内存溢出。

解决办法: 增加路径长度的检查逻辑，如果路径过长则提前返回错误信息。

(2) 问题 2: 当文件不存在或无法访问时，程序崩溃。

解决办法: 增加异常处理机制，如在打开文件失败时打印错误消息而不是直接崩溃。

1.5.4 实验心得

通过这次实验，我对文件系统的内部结构有了更深入的理解。特别是，我学会了如何处理文件路径，如何遍历目录树，以及如何在复杂的文件系统中查找文件。此外，我也意识到了递归搜索在实现这类工具时的重要性，并学会了如何有效地避免无限递归的问题。最后，通过调试和优化代码，我提高了自己的编程技巧和解决问题的能力。

1.6 xargs

1.6.1 实验目的

本实验的目标是实现一个简化版的 `xargs` 程序，该程序能够从标准输入中按行读取文本，并为每一行执行指定的命令，将每一行的内容作为参数传递给该命令。这个工具在处理管道输出时非常有用，因为它允许用户将其他命令的输出作为参数传递给另一个命令。

1.6.2 实验步骤

(1) 分析需求, 创建 `xargs` 文件:

在 `user` 目录下创建名为 `xargs.c` 的新文件。



将 `xargs` 程序添加到 `MAKEFILE` 文件的 `UPROGS` 部分。

```
$U/_xargs\
```

(2) 设计 `xargs` 程序:

首先初始化参数数组, 从命令行参数复制命令名称和参数到 `xargs` 数组中, 为后续添加从标准输入读取的参数做准备。

使用无限循环从标准输入读取一行文本, 直到用户输入为空或遇到 EOF。每次读取一行后, 将这一行作为参数添加到 `xargs` 数组中。

遍历输入的字符串, 处理空格和换行符, 将字符串分割成多个参数, 并将这些参数添加到 `xargs` 数组中。

使用 `fork()` 创建一个新的子进程。在子进程中使用 `exec()` 调用指定的命令, 将 `xargs` 数组作为参数传递给该命令。如果 `exec()` 失败, 则向标准错误输出错误信息, 并让子进程退出。

使用 `wait()` 系统调用来等待子进程的完成, 以避免僵尸进程的产生。

```
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"

int main(int argc, char* argv[]) {
    char buf[512]; // 缓冲区, 用于存储输入的字符串
    char* xargs[32]; // 用于存储命令行参数的数组
    int i;

    // 将命令行参数复制到 xargs 数组中
    for (i = 1; i < argc; i++) {
        xargs[i - 1] = argv[i];
    }
}
```

```
// 无限循环，直到用户输入为空
while (1) {
    int x = argc - 1;
    // 从标准输入读取一行
    if (!gets(buf, sizeof(buf)) || buf[0] == 0)
        break;

    // 将输入的字符串作为参数添加到 xargs 数组中
    xargs[x++] = buf;

    // 遍历输入的字符串，处理空格和换行符
    for (char* p = buf; *p; p++) {
        if (*p == ' ') {
            *p = 0;
            xargs[x++] = p + 1;
        } else if (*p == '\n') {
            *p = 0;
        }
    }
}

// 创建子进程
if (fork() == 0) {
    // 在子进程中执行命令
    exec(argv[1], xargs);
    // 如果 exec 失败，输出错误信息并退出
    fprintf(2, "exec %s failed\n", argv[1]);
    exit(1);
}
```

```
// 等待子进程结束

wait(0);

}
```

(3) 编译与运行：

使用 `make qemu` 编译并运行程序。

进入 xv6 后，输入 `starting sh` 命令来执行程序。

```
hart 1 starting
hart 2 starting
init: starting sh
$ sh < xargstest.sh
$ $ $ $ $ $ hello
hello
hello
$ $
```

(4) 评估与测试：

使用 xv6 自带的测评工具 `./grade-lab-util xargs` 来进行自动测试。评估结果如下，测试通过，程序的行为符合预期。

```
llee11@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-util xargs
make: 'kernel/kernel' is up to date.
== Test xargs == xargs: OK (2.4s)
```

1.6.3 实验中遇到的问题和解决办法

(1) 问题 1：当 `exec()` 成功执行后，原进程中的 `xargs` 数组没有被释放，可能导致内存泄漏。

解决方法：在实际应用中应该考虑内存释放。

1.6.4 实验心得

完成 `xargs` 程序的编写让我深刻地理解了 Unix 系统调用的工作原理，尤其是在进程管理和执行外部命令方面。实验过程中遇到了诸如处理输入中的空格和换行符、内存管理以及 `exec()` 失败的处理等问题，这些问题促使我学习了更多的编程技巧和调试技巧。

1.7 Lab1 实验成绩

在 xv6 目录下添加 `time.txt` 文本文件，写入完成该实验的小时数。在终端中执行 `./grade-lab-util`，即可对整个实验进行自动评分。评估结果如下，测试通过，程序的行为符合预期。

```
llee11@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-util
make: 'kernel/kernel' is up to date.
== Test sleep, no arguments == sleep, no arguments: OK (1.0s)
== Test sleep, returns == sleep, returns: OK (0.9s)
== Test sleep, makes syscall == sleep, makes syscall: OK (1.0s)
== Test pingpong == pingpong: OK (1.0s)
== Test primes == primes: OK (1.1s)
== Test find, in current directory == find, in current directory: OK (1.0s)
== Test find, recursive == find, recursive: OK (1.1s)
== Test xargs == xargs: OK (1.0s)
== Test time ==
time: OK
Score: 100/100
```

实验小结：

通过对这些实验的学习，我对进程的概念有了更全面的理解。进程是操作系统中的一个执行单元，它可以拥有独立的地址空间，并且能够与其他进程进行通信。在 xv6 这样的简单操作系统中，进程的生命周期、进程间通信和进程调度等概念变得更加清晰。实验不仅加深了我对理论知识的理解，而且提高了我的实践技能，尤其是在低级别的系统编程方面。通过这些实验，我认识到理论知识与实际操作之间的联系，并且意识到动手实践对于学习操作系统至关重要。

2. Lab2: system calls

切换到 syscall 分支，输入：

```
git fetch
```

```
git checkout
```

```
syscall make clean
```

```
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git checkout syscall
M      Makefile
Branch 'syscall' set up to track remote branch 'syscall' from 'origin'.
Switched to a new branch 'syscall'
lleell@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
  */*.o */*.d */*.asm */*.sym \
  user/initcode user/initcode.out kernel/kernel fs.img \
  mkfs/mkfs .gdbinit \
  user/usys.S \
  user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln
user/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_gr
ind user/_wc user/_zombie user/_sleep user/_pingpong user/_primes user/_find u
ser/_xargs \
  ph barrier
```

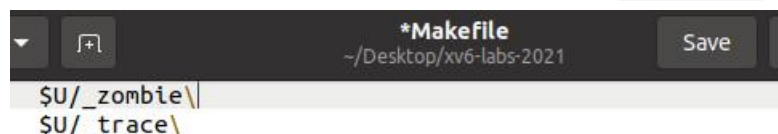
2.1 System call tracing

2.1.1 实验目的

本实验旨在通过创建一个新的系统调用 `trace` 来实现对用户程序中系统调用的追踪功能。具体目标是：实现一个系统调用 `trace`，它接受一个整数参数 `mask` 作为跟踪掩码，用于指定要追踪的系统调用。修改 `xv6` 内核，使得在每个系统调用即将返回时能够打印进程 ID、系统调用名称和返回值。使 `trace` 系统调用能够启用调用它的进程及其后续 `fork` 出的所有子进程的跟踪，不影响其他进程。

2.1.2 实验步骤

(1) 在 `Makefile` 的 `UPROGS` 环境变量中添加 `$U/_trace`



```
*Makefile
~/Desktop/xv6-labs-2021
$U/_zombie\
$U/_trace\
```

(2) 添加系统调用声明和存根：

在 `user/user.h` 中添加 `trace` 系统调用原型：`int trace(int);`

```

Open  ▾  [icon]  *user.h  ~/Desktop/xv6-labs-2021/user  Save
16 int unlink(const char*);
17 int fstat(int fd, struct stat*);
18 int link(const char*, const char*);
19 int mkdir(const char*);
20 int chdir(const char*);
21 int dup(int);
22 int getpid(void);
23 char* sbrk(int);
24 int sleep(int);
25 int uptime(void);
26 int trace(int);
27

```

在 user/usys.pl 脚本中添加 trace 对应的 entry: entry("trace")

```

Open  ▾  [icon]  *usys.pl  ~/Desktop/xv6-labs-2021/user
35 entry("getpid");
36 entry("sbrk");
37 entry("sleep");
38 entry("uptime");
39 entry("trace");

```

在 kernel/syscall.h 中添加 trace 的系统调用号: #define SYS_trace 22

```

Open  ▾  [icon]  *syscall.h  ~/Desktop/xv6-labs-2021/kern
20 #define SYS_link  19
21 #define SYS_mkdir  20
22 #define SYS_close  21
23 #define SYS_trace  22

```

(3) 编写 trace 系统调用函数:

在 kernel/sysproc.c 中实现 sys_trace 函数, 用于处理 trace 系统调用。

```

Open  ▾  [icon]  proc.h  ~/Desktop/xv6-labs-2021
105 struct file *ofile[NOFILE]; // of
106 struct inode *cwd; // cu
107 char name[16]; // Pr
108 int trace_mask;
109 }

```

```

uint64 sys_trace(void) {
    int n;

    // 获取整数类型的系统调用参数
    if(argint(0, &n) < 0){
        return -1;
    }

    struct proc *pro = myproc();
    printf("trace pid: %d\n", pro->pid);
    pro->trace_mask = n;
}

```

```
    return 0;

}
```

(4) 更新 fork 函数:

在 kernel/proc.c 中修改 fork 函数，确保子进程继承父进程的跟踪掩码。

```
np->trace_mask=p->trace_mask;
```

```
proc.c
~/Desktop/xv6-labs-2021/kernel

305
306 pid = np->pid;
307
308 np->trace_mask=p->trace_mask;
309
310 release(&np->lock);
311
```

(5) 添加系统调用名称，修改 syscall 函数:

在 syscall.c 中添加一个 syscall_names 数组，用于存储系统调用名称。

```
*syscall.c
~/Desktop/xv6-labs-2021/ke

106 extern uint64 sys_uptime(void);
107 extern uint64 sys_trace(void);
```

修改 syscall 函数，使其能够在每个系统调用返回时打印追踪信息。

// 定义一个字符数组，用于存储系统调用的名称

// 这个数组的索引与系统调用的编号对应，例如 syscall_name[1]为"fork"，表示系统调用 fork

```
static char *syscall_name[] = {

    "", "fork", "exit", "wait", "pipe", "read", "kill", "exec",

    "fstat", "chdir", "dup", "getpid", "sbrk", "sleep", "uptime",

    "open", "write", "mknod", "unlink", "link", "mkdir", "close",

    "trace"

};
```

// syscall 函数用于处理系统调用

// 每当用户程序请求一个系统调用时，都会调用此函数

```
void syscall(void)
```

```
{
```

```
    int num; // 定义一个变量 num，用于存储系统调用编号
```

```

struct proc *p = myproc(); // 获取当前进程的指针

// 从当前进程的 trapframe 结构体中获取系统调用编号
// trapframe 是一个保存进程在陷入内核模式时的 CPU 寄存器状态的结构
体

num = p->trapframe->a7;

// 检查系统调用编号是否有效（大于 0 且小于 syscalls 数组的元素数量）
// 并且 syscalls[num]指向有效的系统调用函数
if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
    // 调用对应的系统调用函数，并将返回值保存在 trapframe 的 a0 寄存器
    中

    p->trapframe->a0 = syscalls[num]();

    // 检查当前进程的 trace_mask 是否设置了对应的系统调用位
    // trace_mask 用于跟踪特定的系统调用
    if(p->trace_mask & (1 << num)) {
        // 打印追踪信息，显示进程 ID、系统调用名称和返回值
        printf("%d: syscall %s -> %d\n", p->pid, syscall_name[num],
        p->trapframe->a0);
    }
} else {
    // 如果系统调用编号无效，则打印错误信息
    printf("%d %s: unknown sys call %d\n", p->pid, p->name, num);
    // 将返回值设置为-1，表示错误
    p->trapframe->a0 = -1;
}
}

```

（6）编译与运行：

在终端中执行 `make qemu` 编译并运行 `xv6`。

在命令行中输入 `trace 32 grep hello README`，只跟踪 `read()` 系统调用（其中 32 是 `1 << SYS_read` 即 `1 << 5`）。

```
$ trace 32 grep hello README
trace pid: 67
67: syscall read -> 1023
67: syscall read -> 968
67: syscall read -> 235
67: syscall read -> 0
$
```

输入 `trace 2 usertests forkforkfork`，观察一系列系统调用信息的输出。

```
hart 1 starting
hart 2 starting
init: starting sh
$ trace 2 usertests forkforkfork
trace pid: 3
usertests starting
3: syscall fork -> 4
test forkforkfork: 3: syscall fork -> 5
5: syscall fork -> 6
6: syscall fork -> 7
6: syscall fork -> 8
7: syscall fork -> 9
6: syscall fork -> 10
7: syscall fork -> 11
8: syscall fork -> 12
6: syscall fork -> 13
7: syscall fork -> 14
6: syscall fork -> 15
7: syscall fork -> 16
6: syscall fork -> 17
7: syscall fork -> 18
6: syscall fork -> 19
7: syscall fork -> 20
6: syscall fork -> 21
9: syscall fork -> 22
7: syscall fork -> 23
8: syscall fork -> 24
```

（8）评估与测试：

使用 `xv6` 实验自带的测评工具，在终端里输入 `./grade-lab-syscall trace` 进行自动评测。评估结果如下，测试通过，程序的行为符合预期。

```
lllell@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-syscall trace
make: 'kernel/kernel' is up to date.
== Test trace 32 grep == trace 32 grep: OK (0.8s)
== Test trace all grep == trace all grep: OK (1.0s)
== Test trace nothing == trace nothing: OK (1.0s)
== Test trace children == trace children: OK (14.8s)
    (Old xv6.out.trace_children failure log removed)
lllell@ubuntu:~/Desktop/xv6-labs-2021$
```

2.1.3 实验中遇到的问题和解决办法

（1）问题 1：系统调用追踪功能影响其他进程

解决办法：

确认 `trace_mask` 是否只在调用 `trace` 系统调用的进程及其子进程中生效。检查 `fork` 函数中的实现，确保子进程正确地继承了父进程的 `trace_mask`。在 `fork` 函数中使用调试信息来确认子进程是否正确地继承了父进程的 `trace_mask`。

(2) 问题 2：系统调用追踪导致性能下降

解决办法：

测量系统调用追踪功能开启前后程序的执行时间。如果发现性能显著下降，考虑减少输出信息的数量或者仅在调试模式下启用系统调用追踪。使用性能分析工具来识别瓶颈，并优化关键路径上的代码。

2.1.4 实验心得

通过完成本实验，我更深入地了解了如何在用户级程序中调用新增的系统调用，并在实验中验证系统调用的正确性。在系统调用的实现中，用户态和内核态之间的数据传递涉及多个方面。用户程序通过系统调用传递参数，例如使用 `sys_trace` 调用将整数参数 `mask` 传递给内核，内核通过 `argint` 函数从用户空间读取这些参数。系统调用触发了从用户态到内核态的切换，该过程由系统调用机制完成。在内核中，用户传递的参数需要从用户空间复制到内核空间，以确保数据的正确性。同时，内核还需要验证参数的合法性和用户的权限，以防止恶意操作或错误参数的传递。通过仔细阅读系统调用相关文档和源代码，我们能够找到正确的数据传递和状态切换方法，确保系统调用的正确实现。

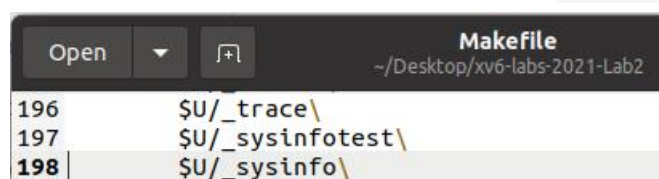
2.2 Sysinfo

2.2.1 实验目的

本次实验的目标是在 `xv6` 操作系统中添加一个新的系统调用 `sysinfo`，该系统调用允许用户空间程序查询当前系统的某些基本信息，包括空闲内存的字节数 (`freemem`) 和非闲置状态的进程数量 (`nproc`)。通过实现这一功能，学生可以深入了解操作系统内核的设计原理，特别是系统调用的机制及其与用户空间程序交互的过程。

2.2.2 实验步骤

(1) 在 `Makefile` 的 `UPROGS` 环境变量中添加 `$(U)/_sysinfotest\` 和 `$(U)/_sysinfo\`



(2) 定义并声明系统调用：

在 `kernel/syscall.h` 文件中添加宏定义 `SYS_sysinfo` 以标识新的系统调用序号。

```
Open  ▾  [+]
```

```
*syscall.h
~/Desktop/xv6-labs-2021/kerne
19 #define SYS_unlink 18
20 #define SYS_link 19
21 #define SYS_mkdir 20
22 #define SYS_close 21
23 #define SYS_trace 22
24 #define SYS_sysinfo 23
```

在 `user/user.h` 文件中预先声明 `struct sysinfo` 和 `sysinfo()` 函数原型。

```
Open  ▾  [+]
```

```
*user.h
~/Desktop/xv6-labs-2021/U
1 struct stat;
2 struct rtcdate;
3 struct sysinfo;
4 int sysinfo(struct sysinfo *);
5
```

(3) 实现系统调用：

在 `kernel/syscall.c` 文件中为 `sysinfo` 系统调用添加对应的处理函数 `sys_sysinfo`。在 `syscall_names` 中新增一个 `sysinfo`。

```
Open  ▾  [+]
```

```
*syscall.c
~/Desktop/xv6-labs-2021/kerne
106 extern uint64 sys_uptime(void);
107 extern uint64 sys_trace(void);
108 extern uint64 sys_sysinfo(void);
109
110
111
112
113
114
115
116
117
```

```
Open  ▾  [+]
```

```
syscall.c
~/Desktop/xv6-labs-2021-Lab2/kernel
Save  ≡  -  □  ×
134 };
135
136 static char *syscall_name[]={ "", "fork", "exit", "wait", "pipe",
    "read", "kill", "exec", "fstat", "chdir", "dup", "getpid", "sbrk",
    "sleep", "uptime", "open", "write", "mknod", "unlink", "link",
    "mkdir", "close", "trace", "sysinfo" };
137
```

在 `kernel/kalloc.c` 文件中实现 `free_mem()` 函数来计算空闲内存的数量。

```
// free_mem 函数用于计算系统中剩余的空闲内存数量

uint64 free_mem(void){
    struct run *r; // 定义一个指向 run 结构体的指针 r，表示当前空闲的内存块

    uint64 num = 0; // 定义一个 64 位无符号整数 num，用于存储空闲内存块的数量

    acquire(&kmem.lock); // 获取 kmem 的自旋锁，确保对空闲内存链表的访问是线程安全的
```

```

r = kmem.freelist; // 将指针 r 指向空闲内存链表的头部
while(r){ // 遍历空闲内存链表
    num++; // 每找到一个空闲内存块，num 递增
    r = r->next; // 移动到下一个空闲内存块
}

release(&kmem.lock); // 释放 kmem 的自旋锁
return num * PGSIZE; // 返回空闲内存的总大小，单位是字节
}

```

在 `kernel/proc.c` 文件中实现 `nproc()` 函数来计算非闲置状态的进程数量。

```

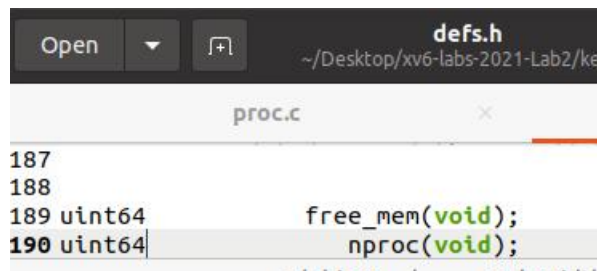
// nproc 函数用于计算当前系统中非 UNUSED 状态的进程数量
uint64 nproc(void){
    struct proc *p; // 定义一个指向 proc 结构体的指针 p，用于遍历进程表
    uint64 num = 0; // 定义一个 64 位无符号整数 num，用于存储非
UNUSED 状态的进程数量

    // 遍历进程表中的每一个进程
    for(p = proc; p < &proc[NPROC]; p++){
        acquire(&p->lock); // 获取当前进程的自旋锁，确保对进程状态的访问是
线程安全的
        if(p->state != UNUSED){ // 如果进程状态不是 UNUSED，则计数加一
            num++;
        }
        release(&p->lock); // 释放当前进程的自旋锁
    }

    return num; // 返回非 UNUSED 状态的进程数量
}

```

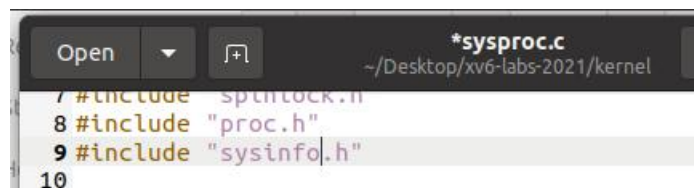
在 `kernel/defs.h` 中添加上述两个新增函数的声明。



```
187
188
189 uint64 free_mem(void);
190 uint64 nproc(void);
```

(4) 编写用户程序:

创建 `sysinfo.c` 文件作为用户程序, 调用 `sysinfo()` 函数并将结果打印出来。



```
7 #include "spinlock.h"
8 #include "proc.h"
9 #include "sysinfo.h"
10
```

// `sys_sysinfo` 函数是 `sysinfo` 系统调用的处理函数

// 它会获取系统的空闲内存和进程数量, 并将这些信息拷贝到用户空间

```
uint64 sys_sysinfo(void) {
```

```
    struct sysinfo info; // 定义一个 sysinfo 结构体实例, 用于存储系统信息
```

```
    uint64 addr; // 用于存储用户传递的地址
```

```
    struct proc *p = myproc(); // 获取当前进程的指针
```

```
    // 从系统调用参数中获取用户传递的地址
```

```
    // argaddr(0, &addr) 将第一个参数的地址存储到 addr 变量中
```

```
    // 如果获取失败, 则返回 -1
```

```
    if (argaddr(0, &addr) < 0) {
```

```
        return -1;
```

```
    }
```

```
    // 调用 free_mem() 函数获取空闲内存的大小, 并将结果存储到 info.freemem
```

```
    info.freemem = free_mem();
```

```
    // 调用 nproc() 函数获取非闲置状态的进程数量，并将结果存储到  
info.nproc
```

```
    info.nproc = nproc();
```

```
    // 将 info 结构体中的数据拷贝到用户空间
```

```
    // copyout 函数将内核空间的数据拷贝到用户空间的指定地址
```

```
    // 如果拷贝失败，则返回 -1
```

```
    if (copyout(p->pagetable, addr, (char *)&info, sizeof(info)) < 0) {
```

```
        return -1;
```

```
    }
```

```
    // 返回 0 表示成功
```

```
    return 0;
```

```
}
```

在 user 目录下添加一个 sysinfo.c 用户程序

```
#include "kernel/param.h" // 包含参数定义
```

```
#include "kernel/types.h" // 包含类型定义
```

```
#include "kernel/sysinfo.h" // 包含 sysinfo 结构体的定义
```

```
#include "user/user.h" // 包含用户程序相关的函数和宏
```

```
int main(int argc, char *argv[]) {
```

```
    // 参数错误检查
```

```
    // 如果命令行参数不为 1（程序名），则输出用法提示并退出
```

```
    if (argc != 1) {
```

```
        fprintf(2, "Usage: %s need not param\n", argv[0]);
```

```
        exit(1);
```

```
    }
```

```
    struct sysinfo info; // 定义一个 sysinfo 结构体实例，用于接收系统信息
```

```

// 调用 sysinfo 系统调用，并将结果存储到 info 结构体中
sysinfo(&info);

// 打印系统信息

// info.freemem 表示空闲内存的大小

// info.nproc 表示非闲置状态的进程数量

printf("free space: %d\nused process: %d\n", info.freemem, info.nproc);

// 退出程序

exit(0);
}

```

(5) 编译和运行：

在终端中执行 `make qemu` 编译并运行 `xv6`。

在命令行中输入 `sysinfotest` 。

```

hart 2 starting
hart 1 starting
init: starting sh
$ sysinfo
free space: 133386240
used process: 3
$ sysinfotest
sysinfotest: start
sysinfotest: OK
$ 

```

(6) 评估与测试：

使用 `xv6` 实验自带的测评工具，在终端里输入 `./grade-lab-syscall trace` 进行自动评测。评估结果如下，测试通过，程序的行为符合预期。

```

lleeell@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-syscall sysinfo
make: 'kernel/kernel' is up to date.
== Test sysinfotest == sysinfotest: OK (2.4s)

```

2.2.3 实验中遇到的问题和解决办法

(1) 问题 1：在实现 `sysinfo` 系统调用时，遇到了如何安全地访问和更新进程状态的问题。

解决办法：通过使用互斥锁 `acquire()` 和 `release()` 来确保对进程状态的原子性访问，避免了并发访问导致的数据不一致问题。

(2) 问题 2：在复制数据到用户空间时出现了错误。

解决办法：使用 `copyout()` 函数，并确保在复制前对指针进行了正确的验证，防止非法内存访问引发的崩溃或数据损坏。

(3) 问题 3：编译过程中遇到了类型不匹配的错误。

解决办法：检查并确保所有使用的类型和定义与现有的 xv6 代码库保持一致，例如确保 `uint64` 类型被正确地定义和使用。

2.2.4 实验心得

通过这次实验，我深刻理解了操作系统内核如何通过系统调用来与用户空间程序进行通信。特别是在处理内核与用户空间之间的数据传递时，我学会了如何正确地使用 `copyout()` 函数。此外，实验还加强了我在一个简单的操作系统内核环境中，对进程管理和内存管理的理解。通过调试和解决问题的过程，我也提高了自己的编程技能，尤其是对于并发编程中锁的使用有了更深入的认识。

2.3 Lab2 实验成绩

新建文件 `time.txt`，在其中写入实验所需时间(小时)，输入 `make grade`，查看 lab2 成绩。评估结果如下，测试通过，程序的行为符合预期。

```
== Test trace 32 grep ==
$ make qemu-gdb
trace 32 grep: OK (2.6s)
== Test trace all grep ==
$ make qemu-gdb
trace all grep: OK (0.6s)
== Test trace nothing ==
$ make qemu-gdb
trace nothing: OK (0.8s)
== Test trace children ==
$ make qemu-gdb
trace children: OK (13.9s)
== Test sysinfotest ==
$ make qemu-gdb
sysinfotest: OK (1.8s)
== Test time ==
time: OK
Score: 35/35
```

实验小结：

在本次实验中，通过实现系统调用追踪 (`trace`) 和系统信息查询 (`sysinfo`)，加深了我对操作系统中系统调用机制的理解。在 `trace` 实验中，我学会了如何在内核中追踪特定的系统调用，并确保追踪功能仅影响调用进程及其子进程。而在 `sysinfo` 实验中，我掌握了如何从用户空间安全地获取内核信息，如空闲内存和活跃进程数量。实验过程中遇到的问题，如并发访问和数据传递，通过使用锁机制和确保数据一致性得到了解决。这些实践增强了我的编程能力和对操作系统核心概念的理解。

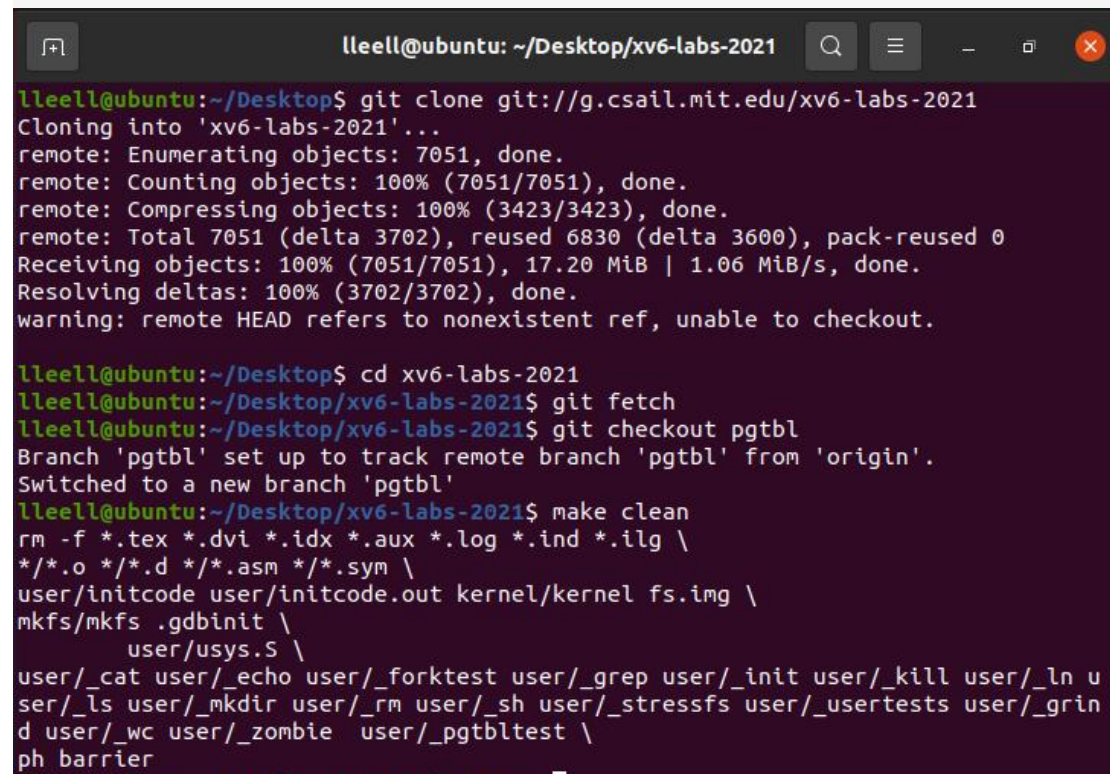
3. Lab3: page tables

切换到 `pgtbl` 分支

```
git fetch
```

```
git checkout pgtbl
```

```
make clean
```



```
lleell@ubuntu: ~/Desktop/xv6-labs-2021
lleell@ubuntu:~/Desktop$ git clone git://g.csail.mit.edu/xv6-labs-2021
Cloning into 'xv6-labs-2021'...
remote: Enumerating objects: 7051, done.
remote: Counting objects: 100% (7051/7051), done.
remote: Compressing objects: 100% (3423/3423), done.
remote: Total 7051 (delta 3702), reused 6830 (delta 3600), pack-reused 0
Receiving objects: 100% (7051/7051), 17.20 MiB | 1.06 MiB/s, done.
Resolving deltas: 100% (3702/3702), done.
warning: remote HEAD refers to nonexistent ref, unable to checkout.

lleell@ubuntu:~/Desktop$ cd xv6-labs-2021
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git checkout pgtbl
Branch 'pgtbl' set up to track remote branch 'pgtbl' from 'origin'.
Switched to a new branch 'pgtbl'
lleell@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
  */*.o */*.d */*.asm */*.sym \
  user/initcode user/initcode.out kernel/kernel fs.img \
  mkfs/mkfs .gdbinit \
  user/usys.S \
  user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln u
ser/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_grin
d user/_wc user/_zombie user/_pgtbltest \
  ph barrier
```

3.1 Speed up system calls

3.1.1 实验目的

本实验旨在学习如何通过页表机制加速系统调用，特别是针对那些仅涉及内核数据结构读取而不进行写入操作的系统调用，创建一个共享的只读页，在内核和用户程序之间建立映射，使得内核向该页写入数据时，用户程序可以直接读取这些数据，无需进行复杂的系统调用。这不仅提高了系统调用的效率，也减少了数据复制的开销。

3.1.2 实验步骤

(1) 修改内核结构:

在 kernel/proc.h 的 struct proc 中添加一个新的成员 struct usyscall *usyscallpage; 用于保存共享页面的地址。

```
Open  ~/Desktop/xv6-labs-2021/kernel
106 struct inode *cwd;           // Current
107 char name[16];               // Process name
108 struct usyscall *usyscallpage;
109 };
```

在 kernel/proc.c 的 allocproc() 函数中分配并初始化共享页。

```
Open  ~/Desktop/xv6-labs-2021/kernel  Save
130 if((p->usyscallpage = (struct usyscall *)kalloc()) == 0)
131     freeproc(p);
132     release(&p->lock);
133     return 0;
134
```

在 kernel/proc.c 的 freeproc() 函数中释放共享页。

```
Open  ~/Desktop/xv6-labs-2021/kernel
165 if(p->usyscallpage)
166     kfree((void*)p->usyscallpage);
167     p->usyscallpage = 0;
```

在 kernel/proc.c 的 proc_freepagetable() 函数中释放页表中的相应项。

```
Open  ~/Desktop/xv6-labs-2021/kernel  Save
208 uvmfree(pagetable, 0);
209 return 0;
210 }
211
212 if(mappages(pagetable, USYSCALL, PGSIZE, (uint64)(p->usyscallpage), PTE_R | PTE_U) < 0) {
213     uvmfree(pagetable, 0);
214     return 0;
215 }
216
```

(2) 页表映射:

在 kernel/proc.c 的 proc_pagetable() 函数中为新创建的进程分配一个共享的只读页。映射时设置权限为只读 (PTE_R) 且可由用户访问 (PTE_U)。

```
Open  ~/Desktop/xv6-labs-2021/kernel  Save
208 uvmfree(pagetable, 0);
209 return 0;
210 }
211
212 if(mappages(pagetable, USYSCALL, PGSIZE, (uint64)(p->usyscallpage), PTE_R | PTE_U) < 0) {
213     uvmfree(pagetable, 0);
214     return 0;
215 }
216
```

(3) 编译和运行:

在终端中执行 make qemu 编译并运行 xv6。在命令行中输入 pgtbltest。

```

hart 1 starting
hart 2 starting
init: starting sh
$ pgtbltest
ugetpid_test starting
ugetpid_test: OK
pgaccess_test starting
pgtbltest: pgaccess_test failed: incorrect access bits set, pid=3

```

(4) 评估与测试:

使用 xv6 实验自带的测评工具，在终端里输入 `./grade-lab-pgtbl ugetpid` 进行自动评测。评估结果如下，测试通过，程序的行为符合预期。

```

lleell@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-pgtbl ugetpid
make: 'kernel/kernel' is up to date.
== Test pgtbltest == (1.6s)
== Test   pgtbltest: ugetpid ==
pgtbltest: ugetpid: OK

```

3.1.3 实验中遇到的问题和解决办法

(1) 问题 1: 共享页面的分配和初始化

解决办法: 在尝试分配和初始化共享页面时，遇到了内存分配失败的情况。我首先检查了内存分配函数是否正确调用，然后确认了内存池中是否有足够的可用内存。最终发现是因为在内存不足的情况下没有适当地处理错误情况。解决方法是在分配失败时返回错误码，并在调用处进行适当的错误处理。

(2) 问题 2: 在设置页表映射权限时，我遇到了页面无法被用户程序访问的问题。经过检查发现，权限位设置不正确，虽然设置了只读权限（`PTE_R`），却没有设置用户可访问权限（`PTE_U`）。

解决办法: 根据 xv6 的设计，用户程序不应直接写入代码或数据段，操作系统负责这些操作。因此，映射这些页面时应使用只读权限（`PTE_R`）和用户权限（`PTE_U`），而不是可写权限（`PTE_W`）。修正权限位设置后，确保页面既具有只读权限，又具有用户可访问权限，问题得到了解决。

3.1.4 实验心得

通过本次实验，我对页表机制及如何优化系统调用有了更深刻的理解。我学习了如何在内核和用户程序之间创建共享的只读页，并利用页表映射来简化数据从用户空间到内核空间的复制过程。在实现页表映射的函数时，我遇到了一些问题，例如页表权限设置不当导致页面不可访问、页表遍历逻辑错误等问题，这些问题促使我深入研究 RISC-V 页表机制，并通过查阅文档和调试逐步解决了这些问题。通过这次实验，我认识到理论与实践相结合的重要性，通过实际操作加深了对操作系统内存管理机制的理解，对未来的学习和研究有着积极的影响。

3.2 Print a page table

3.2.1 实验目的

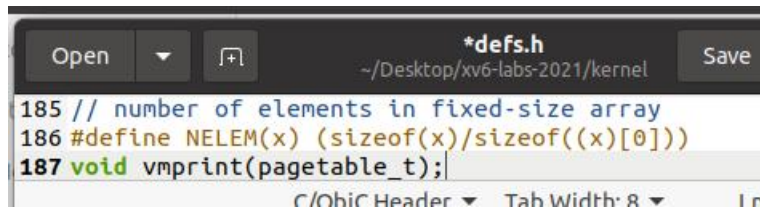
本次实验的目标是为了可视化 RISC-V 页表，并可能有助于将来的调试工作。为此，需要编写一个名为 `vmprint()` 的函数，用于打印页表内容。该函数接受一个 `pagetable_t` 类型的参数，并以指定的格式打印该页表。

3.2.2 实验步骤

(1) 定义 `vmprint()` 函数：

在 `kernel/vm.c` 文件中添加 `vmprint()` 函数。

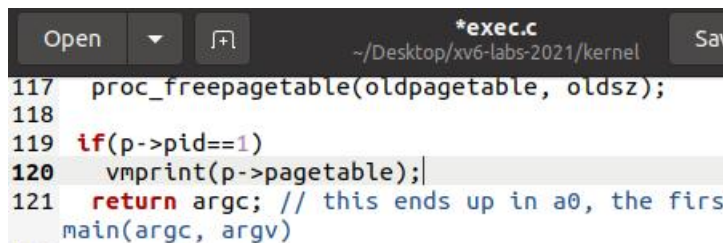
定义 `vmprint()` 函数原型在 `kernel/defs.h` 文件中，以便在 `exec.c` 中引用。函数内部遍历页表，并以指定格式打印每个页表条目的索引、十六进制表示以及对应的物理地址。



```
185 // number of elements in fixed-size array
186 #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
187 void vmprint(pagetable_t);
```

(2) 修改 `exec()` 函数：

打开 `kernel/exec.c` 文件。在 `return argc;` 语句之前插入条件语句 `if (p->pid == 1) vmprint(p->pagetable);`；这样做会在 `init` 进程被创建时打印其页表。



```
117 proc_freepagetable(oldpagetable, oldsz);
118
119 if(p->pid==1)
120     vmprint(p->pagetable);
121 return argc; // this ends up in a0, the first argument to
main(argc, argv)
```

(3) 编写辅助函数：

根据 `freewalk()` 函数的实现方式，编写一个名为 `vmprintwalk()` 的辅助函数，该函数将递归地遍历页表结构并打印相关信息。


```
Open  [icon]  *vm.c  Save  [icon]
~/Desktop/xv6-labs-2021/kernel

284 void
285 vmprint(pagetable_t pagetable)
286 {
287     // vmprint times(0,1,2)
288     static int num = 0;
289     if(num == 0)
290         printf("page table %p\n",pagetable);
291     for(int i = 0; i < 512; i++){
292         pte_t pte = pagetable[i];
293         if(pte & PTE_V){
294             for(int j=0; j<=num; ++j){
295                 printf(" ..");
296             }
297             if(num!=2){
298                 printf("%d: pte %p pa %p\n",i,pte,PTE2PA(pte));
299                 uint64 child = PTE2PA(pte);
300                 ++num;;
301                 vmprint((pagetable_t)child);
302                 --num;
303             }else{
304                 printf("%d: pte %p pa %p\n",i,pte,PTE2PA(pte));
305             }
306         }
307     }
308 }
```

(4) 编译和运行:

在终端中执行 `make qemu` 编译并运行 `xv6`。(按照一定的格式打印出每个 PTE 的索引、 PTE 的 16 进制表示和从 PTE 中提取的物理地址):

```
hart 2 starting
hart 1 starting
page table 0x0000000087f6e000
..0: pte 0x0000000021fda801 pa 0x0000000087f6a000
.. ..0: pte 0x0000000021fda401 pa 0x0000000087f69000
.. .. ..0: pte 0x0000000021fdac1f pa 0x0000000087f6b000
.. .. ..1: pte 0x0000000021fda00f pa 0x0000000087f68000
.. .. ..2: pte 0x0000000021fd9c1f pa 0x0000000087f67000
..255: pte 0x0000000021fdb401 pa 0x0000000087f6d000
.. ..511: pte 0x0000000021fdb001 pa 0x0000000087f6c000
.. .. ..509: pte 0x0000000021fdd813 pa 0x0000000087f76000
.. .. ..510: pte 0x0000000021fddc07 pa 0x0000000087f77000
.. .. ..511: pte 0x0000000020001c0b pa 0x0000000080007000
init: starting sh
```

页表条目的分析:

- 第一行:

pte (页表条目) 为 `0x0000000021fda801` , pa (物理地址) 为 `0x0000000087f6a000`。

我们可以使用 pte 和页表的起始地址来计算对应的物理地址。具体来说, 假设每个页表项的大小为 4 字节 (32 位)。

pte 拆分为两部分:

高 20 位 `0x00000000` 代表页表条目的索引。

低 12 位 `0x21fda801` 代表页表条目中的标志位。

通过将页表条目的索引左移 2 位（相当于乘以 4），然后加上页表的起始物理地址 0x0000000087f6e000，我们可以得到对应的物理地址。

- 其他行：

后续的行依次类推，表示更深层次的页表条目和对应的物理地址。

例如第二行：pte 为 0x0000000021fdac1f，pa 为 0x0000000087f6b000。类似地，我们可以使用前述的方法来计算和验证其对应的物理地址。

- 总结：

这段输出展示了一个进程的页表的层级结构以及每个页表条目的虚拟地址和物理地址的映射关系。

页表条目的 pte 值和页表的起始地址组合在一起，可以帮助操作系统将虚拟地址转换为对应的物理地址，实现内存管理和地址转换的功能。

通过这些映射，操作系统能够确保在用户态下运行的程序访问到正确的物理内存位置，从而正确执行计算任务。这就是页表和虚拟内存管理的核心功能。

（6）思考：根据实验尝试解释 vmprint 的输出。

- 第 0 页包含什么内容？

第 0 页包含指向其他页表页的指针，以构建页表的层次结构。

- 第 2 页包含什么？

第 2 页实际上是指向第 1 页的 PTE（位于第 1 页的索引 2 处）。第 1 页是第一个第 1 级页表页，它包含指向第 2 级页表页的 PTE。所以第 2 页中的 PTE 包含指向下一级页表页（第 2 级页表页）的物理地址。

- 以用户模式运行时，进程能否读/写第 1 页映射的内存？

可以参考上一个实验的经验，用户的权限通常被设置为只读（PTE_R 和 PTE_U）。因此，进程可以读取这些页面的内容，但不能写入它们，因为权限没有包括 PTE_W。

倒数第三页包含什么内容？

在 xv6 中，倒数第三页是指向用户栈的 PTE，它的权限通常是可读、可写和用户权限。用户栈用于存储进程在用户模式下调用函数时的局部变量和临时数据。

3.2.3 实验中遇到的问题和解决办法

（1）问题 1：不知道如何确定页表的起始物理地址。

解决办法：通常页表的起始物理地址是由处理器设置的，可以通过查询系统启动代码或相关配置来获取该值。在 RISC-V 架构中，可以通过读取控制寄存器（如 satp）来获取页表基址。

（2）问题 2：输出的页表条目物理地址不正确。

解决办法：仔细检查页表条目到物理地址的计算逻辑，确保位移操作正确无误。可以使用调试工具单步执行代码，观察变量值的变化。

3.2.4 实验心得

通过本次实验，我深刻体会到了页表机制在现代操作系统内存管理中的重要性。动手实现 `vmprint()` 函数的过程让我对页表条目的结构和页表的工作原理有了直观的认识。尤其在调试过程中，逐步理解了虚拟地址到物理地址转换的细节，这对我今后深入研究内存管理技术大有裨益。实验中遇到的问题促使我查阅了大量的资料，不仅解决了难题，还拓宽了知识面。这次经历让我更加确信，实践是检验理论的最佳方式，也是提升技能不可或缺的一环。

3.3 Detecting which pages have been accessed

3.3.1 实验目的

本实验旨在实现 `pgaccess()` 系统调用，该调用能够报告哪些页面已被系统调用访问过。系统调用会通过检查 RISC-V 页表中的访问位来实现这一功能，从一个用户页表地址开始，搜索所有被访问过的页，并返回一个位图(bitmap)，以显示这些页是否被访问过。

在现代操作系统中，内存管理是关键的一环。操作系统通常使用页表来管理虚拟地址到物理地址的映射。RISC-V 架构中的页表 (Page Table) 是实现虚拟内存的核心组件之一。每个页表项 (PTE) 中包含了页面的各种信息，例如是否被访问过、是否已被写入等。

页表项 (PTE) 中的访问位：在 RISC-V 架构中，每个页表项包含多个标志位，其中一个重要的标志位是访问位 (Accessed Bit)，通常表示页面自上次清除访问位之后是否被读取或写入。这个访问位对于内存管理和垃圾回收机制尤其重要。通过监控页面的访问情况，操作系统可以更好地管理内存，例如实现页面的回收或优化内存的使用。

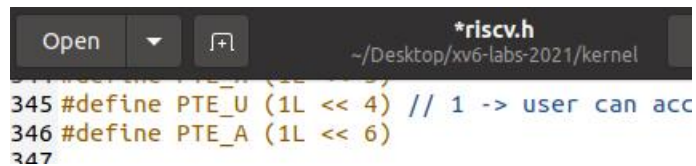
系统调用 `pgaccess()` 的应用：垃圾回收器 (Garbage Collector) 是一种自动内存管理工具，通常用于检测和回收不再使用的内存。为了高效地进行垃圾回收，垃圾回收器需要知道哪些内存页面被访问过。通过实现一个系统调用 `pgaccess()`，我们可以让用户空间的程序查询哪些页面已经被访问过。这对于垃圾回收器和其他需要跟踪内存访问情况的应用非常有用。

3.3.2 实验步骤

(1) 定义访问位

在 kernel/riscv.h 中定义了一个标志 PTE_A, 表示 RISC-V 架构中的访问位。

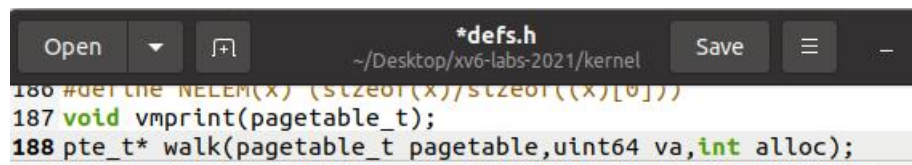
```
#define PTE_A (1L << 6)
```



```
345 #define PTE_U (1L << 4) // 1 -> user can acc
346 #define PTE_A (1L << 6)
347
```

(2) 声明 walk 函数:

在 defs.h 中声明 walk 函数, 该函数用于遍历页表并定位特定虚拟地址对应的页表项。



```
186 #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
187 void vmprint(pagetable_t);
188 pte_t* walk(pagetable_t pagetable, uint64 va, int alloc);
```

(3) 实现 sys_pgaccess 系统调用

在 kernel/sysproc.c 中实现了 sys_pgaccess() 函数, 它获取系统调用的三个参数: va: 需要检查页面的初始地址。

pagenum: 从初始地址向后检查的页面数。

abitsaddr: 存储结果位图的地址。

```
int
sys_pgaccess(void)
{
    // 定义变量

    uint64 va; // 存储传入的虚拟地址(va)参数
    int pagenum; // 存储传入的页面数量(pagenum)参数
    uint64 abitsaddr; // 存储传入的位图地址(abitsaddr)参数

    // 获取系统调用的参数值

    argaddr(0, &va); // 获取第一个参数 va: 需要检查的页面的起始虚拟地址
    argint(1, &pagenum); // 获取第二个参数 pagenum: 需要检查的页面数量
    argaddr(2, &abitsaddr); // 获取第三个参数 abitsaddr: 存储结果位图的地址

    uint64 maskbits = 0; // 用于保存已访问页面的位图, 初始为 0

    struct proc *proc = myproc(); // 获取当前进程的进程控制块 (PCB)
```



```

// 遍历每个页面，检查是否被访问过
for (int i = 0; i < pagenum; i++) {
    pte_t *pte = walk(proc->pagetable, va+i*PGSIZE, 0); // 通过 walk 函数找到虚拟地址对应的页表项 (PTE)

    // 如果页表项不存在，抛出错误
    if (pte == 0)
        panic("page not exist.");

    // 检查 PTE 中是否设置了访问位(PTE_A)
    if (PTE_FLAGS(*pte) & PTE_A) {
        maskbits = maskbits | (1L << i); // 如果该页面被访问过，则在 maskbits 中对应的位置 1
    }

    // 清除 PTE 中的访问位，将 PTE_A 位设置为 0
    *pte = ((*pte & PTE_A) ^ *pte) ^ 0;
}

// 将 maskbits 位图复制到用户空间的指定地址
if (copyout(proc->pagetable, abitsaddr, (char *)&maskbits, sizeof(maskbits)) < 0)
    panic("sys_pgaccess copyout error"); // 如果复制失败，抛出错误

return 0; // 成功执行，返回 0
}

```

(4) 编译和运行：

在终端中执行 `make qemu` 编译并运行 `xv6`。在命令行中输入 `pgtbltest`。

```

init: starting sh
$ pgtbltest
ugetpid_test starting
ugetpid_test: OK
pgaccess_test starting
pgaccess_test: OK
pgtbltest: all tests succeeded

```

3.3.3 实验中遇到的问题和解决办法

(1) 问题 1：访问位清除失败，在尝试清除访问位时，可能会导致页面无法正常访问或其他错误。

解决办法：确保清除访问位的操作不会影响到其他页表项标志位。测试清除逻辑，单独测试清除逻辑，确保它不会引起系统不稳定或其他副作用。检查硬件支持，确认目标架构（如 RISC-V）是否支持访问位的硬件清除。

3.3.4 实验心得

通过本次实验，我们不仅深入理解了操作系统如何管理内存和页面访问，还学会了如何利用页表的硬件机制来优化内存管理。特别是在处理垃圾回收机制时，利用访问位可以显著提高效率，避免不必要的内存扫描操作。

3.4 Lab3 实验成绩

新建文件 `time.txt`，在其中写入实验所需时间（小时），新建文件 `answers-pgtbl.txt` 填入之前实验的测试结果。输入 `make grade`，查看 lab3 成绩。评估结果如下，测试通过，程序的行为符合预期。

```
make[1]: Leaving directory '/home/lleell/Desktop/xv6-labs-2021'
== Test pgtbltest ==
$ make qemu-gdb
(2.7s)
== Test pgtbltest: ugetpid ==
pgtbltest: ugetpid: OK
== Test pgtbltest: pgaccess ==
pgtbltest: pgaccess: OK
== Test pte printout ==
$ make qemu-gdb
pte printout: OK (1.1s)
== Test answers-pgtbl.txt == answers-pgtbl.txt: OK
== Test usertests ==
$ make qemu-gdb
(209.9s)
== Test usertests: all tests ==
usertests: all tests: OK
== Test time ==
time: OK
Score: 46/46
```

实验小结：

通过对 RISC-V 架构中的内存管理单元(MMU)与操作系统协作机制的实验分析，我们初步理解了 RISC-V 的分页机制。然而，实验并未深入探讨 RISC-V 硬件的分页细节，因此我们参考《RISC-V 指令集手册》中 Sv-39 分页机制进行总结。Sv-39 支持 39 位虚拟地址，通过三级页表系统管理 4KB 的页面。虚拟地址的高 27 位分为三段，每段 9 位，分别对应三级页表的条目，物理地址则由 44 位物理页号和页内偏移构成。分页转换时，处理器通过 SATP 寄存器中的根页表物理页面号依次查找页表项，检查权限并更新页面状态。Sv-39 的设计体现了 RISC-V 架构的灵活性与扩展性，不仅支持大规模虚拟地址空间，还提供了精细的权限控制和状态管理机制。通过理论分析，我们进一步理解了 Sv-39 在内存管理中的关键作用，为后续更深入的研究奠定了基础。

4. Lab4: page tables

切换到 `traps` 分支:

```
git fetch
git checkout traps
make clean
```

```
llee11@ubuntu:~/Desktop$ git clone git://g.csail.mit.edu/xv6-labs-2021
Cloning into 'xv6-labs-2021'...
remote: Enumerating objects: 7051, done.
remote: Counting objects: 100% (7051/7051), done.
remote: Compressing objects: 100% (3423/3423), done.
remote: Total 7051 (delta 3701), reused 6831 (delta 3600), pack-reused 0
Receiving objects: 100% (7051/7051), 17.20 MiB | 1.72 MiB/s, done.
Resolving deltas: 100% (3701/3701), done.
warning: remote HEAD refers to nonexistent ref, unable to checkout.

llee11@ubuntu:~/Desktop$ cd xv6-labs-2021
llee11@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
llee11@ubuntu:~/Desktop/xv6-labs-2021$ git checkout traps
Branch 'traps' set up to track remote branch 'traps' from 'origin'.
Switched to a new branch 'traps'
llee11@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
*/*.o */*.d */*.asm */*.sym \
user/initcode user/initcode.out kernel/kernel fs.img \
mkfs/mkfs .gdbinit \
    user/usys.S \
user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln u
ser/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_grin
d user/_wc user/_zombie user/_call user/_btest \
ph barrier
```

4.1 RISC-V assembly

4.1.1 实验目的

本次实验的目的是了解 RISC-V 程序集, 在 `xv6 repo` 中有一个文件 `user/call.c`。 `make fs.img` 会对其进行编译, 并生成 `user/call.asm` 中程序的可读汇编版本。我们需要阅读 `call.asm` 中的函数 `g`、`f` 和 `main` 的代码, 并回答问题。

4.1.2 实验步骤

(1) 在终端输入 `make fs.img` 编译 `user/call.c`, 在 `user/call.h` 中生成可读的汇编版本, 阅读 `call.asm` 中函数 `g`、`f` 和 `main` 的代码:

```
llee11@ubuntu:~/Desktop/xv6-labs-2021$ make fs.img
gcc -DSOL_TRAPS -DLAB_TRAPS -Werror -Wall -I. -o mkfs/mkfs mkfs/mkfs.c
riscv64-linux-gnu-gcc -Wall -Werror -O -fno-omit-frame-pointer -ggdb -DSOL_TRAP
S -DLAB_TRAPS -MD -mmodel=medany -ffreestanding -fno-common -nostdlib -mno-rel
ax -I. -fno-stack-protector -fno-pie -no-pie -c -o user/ulib.o user/ulib.c
perl user/usys.pl > user/usys.S
```

```
callasm
~/Desktop/xv6-labs-2021/user

11 #include "user/user.h"
12
13 int g(int x) {
14     0: 1141          addi    sp,sp,-16
15     2: e422          sd      s0,8(sp)
16     4: 0800          addi    s0,sp,16
17     return x+3;
18 }
19     6: 250d          addiw   a0,a0,3
20     8: 6422          ld      s0,8(sp)
21     a: 0141          addi    sp,sp,16
22     c: 8082          ret
23
24 0000000000000000e <f>:
25
26 int f(int x) {
27     e: 1141          addi    sp,sp,-16
28    10: e422          sd      s0,8(sp)
29    12: 0800          addi    s0,sp,16
30     return g(x);
31 }
32    14: 250d          addiw   a0,a0,3
33    16: 6422          ld      s0,8(sp)
34    18: 0141          addi    sp,sp,16
35
```

```
callasm
~/Desktop/xv6-labs-2021/user

36
37 0000000000000001c <main>:
38
39 void main(void) {
40     1c: 1141          addi    sp,sp,-16
41     1e: e406          sd      ra,8(sp)
42     20: e022          sd      s0,0(sp)
43     22: 0800          addi    s0,sp,16
44     printf("%d %d\n", f(8)+1, 13);
45     24: 4635          li      a2,13
46     26: 45b1          li      a1,12
47     28: 00000517      auipc   a0,0x0
48     2c: 7b050513      addi    a0,a0,1968 # 7d8
49     <malloc+0xea>
50     30: 00000097      auipc   ra,0x0
51     34: 600080e7      jalr    1536(ra) # 630 <printf>
52     exit(0);
53     38: 4501          li      a0,0
54     3a: 00000097      auipc   ra,0x0
55     3e: 27e080e7      jalr    638(ra) # 2b8 <exit>
56
```

对于以下问题：

(1) Which registers contain arguments to functions? For example, which register holds 13 in main's call to printf?

寄存器 a0 到 a7（即 x10 到 x17）是用于存放函数调用的参数，由以下代码可以得到在调用 printf 的时候，寄存器 a2 保存的是 13。

```

44 printf("%d %d\n", f(8)+1, 13);
45 24: 4635 | li a2,13
46 26: 45b1 | li a1,12

```

(2) Where is the call to function f in the assembly code for main? Where is the call to g? (Hint: the compiler may inline functions.)

函数 f 调用了函数 g；函数 g 将传入的参数加上 3 后返回。编译器在生成代码时进行了内联优化，虽然在代码中看起来是通过 printf("%d %d\n", f(8)+1, 13) 调用了 f 函数，但实际上在汇编代码中，编译器直接将 f(8)+1 计算并替换为 12。这表明编译器在编译阶段已经对函数调用进行了优化，因此在 main 函数的汇编代码中，没有直接调用 f 和 g，而是编译器在运行前已经计算并替换了这些函数的结果。

```

43 22: 0800 | addi s0,sp,16
44 printf("%d %d\n", f(8)+1, 13);
45 24: 4635 | li a2,13

```

(3) At what address is the function printf located?

这个 auipc 指令是将 0x0 左移 12 位，然后加上当前的 PC 值 0x30，并存储到 ra 寄存器中。接下来的 jalr 指令会根据 ra 计算出跳转地址，并将当前 PC + 4 存入 ra，作为稍后 printf 的返回地址。所以，printf 函数的位置是 0x30+1536=0x630。查阅代码后确认其地址确实在 0x630。

```

1094 void
1095 printf(const char *fmt, ...)
1096 {
1097 630: 711d | addi sp,sp,-96
1098 632: ec06 | sd ra,24(sp)

```

(4) What value is in the register ra just after the jalr to printf in main?

执行完 jalr 跳转到 printf 之后，ra 寄存器的值为 0x38。具体过程如下：
在程序的 0x30 地址处，auipc ra, 0x0 将当前程序计数器 pc 的值存入 ra 中。接着，在 0x34 处，jalr 1536(ra) 跳转到偏移后的地址，即 printf 函数所在的 0x630 位置。根据参考资料的信息，执行这条 jalr 指令后，ra 的值会设置为当前 pc+4，也就是返回地址 0x38。

```

49 30: 00000097 | auipc ra,0x0
50 34: 600080e7 | jalr 1536(ra) # 630 <printf>

```

(5) Run the following code.

```

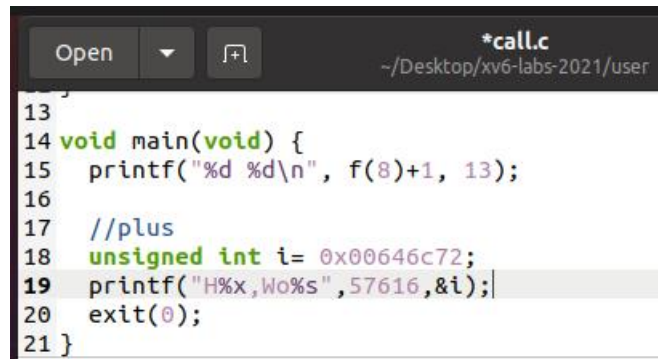
unsigned int i = 0x00646c72;
printf("H%x Wo%s", 57616, &i);

```

What is the output? Here's an ASCII table that maps bytes to characters.

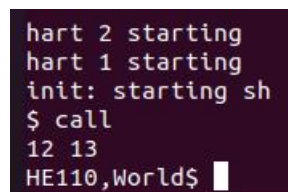
The output depends on that fact that the RISC-V is little-endian. If the RISC-V were instead big-endian what would you set `i` to in order to yield the same output? Would you need to change 57616 to a different value?

根据提示，在 `user/call.c` 中加入代码：



```
13
14 void main(void) {
15     printf("%d %d\n", f(8)+1, 13);
16
17     //plus
18     unsigned int i= 0x00646c72;
19     printf("H%x,Wo%s",57616,&i);
20     exit(0);
21 }
```

执行 `make qemu` 并输入 `call` 命令，得到了输出：打印 HE110 World。



```
hart 2 starting
hart 1 starting
init: starting sh
$ call
12 13
HE110,World$
```

在 RISC-V 体系结构中，由于其采用小端存储格式，当 `i` 的值为 `0x00646c72` 时，字节顺序在内存中存储为 `72 6c 64 00`，对应的 ASCII 字符为 `"rld"`。而数字 57616 转换为 16 进制为 `0xe110`，因此 `%x` 格式化符号打印输出 `"HE110 World"`。

在大端存储模式下，为了保证输出的字符顺序与小端一致，`i` 的值需要调整为 `0x726c6400`，这样字节顺序在内存中为 `00 64 6c 72`，对应的 ASCII 字符依然为 `"rld"`。无论是大端还是小端，57616 的十六进制值始终为 `0xe110`，因此打印结果相同。因此，无需改变 57616 的值，但需调整 `i` 的值以适应大端存储格式下的字符输出顺序。

(6) In the following code, what is going to be printed after `'y='`? (note: the answer is not a specific value.) Why does this happen?

```
printf("x=%d y=%d", 3);
```

根据提示，在 `user/call.c` 中改变代码为：



```
13
14 void main(void) {
15     //printf("%d %d\n", f(8)+1, 13);
16     printf("x=%d y=%d\n", 3);
17     //plus
18     //unsigned int i= 0x00646c72;
```

执行 `make qemu` 并输入 `call` 命令，得到了输出：x=3 y=5221

```
hart 1 starting
hart 2 starting
init: starting sh
$ call
x=3 y=5221
$
```

在这个例子中，`printf` 函数因少传递了一个参数而导致输出不确定的结果。根据函数传参规则，`y=` 后的值应该来自寄存器 `a2`。由于 `a2` 寄存器未被显式赋值，它保留了调用之前的值，这可能是随机的或受之前代码的影响。因此，`y=` 后面的值为一个无法预测的不可靠值。简单来说，`printf` 函数试图读取的参数数量超过了实际提供的参数，因此从寄存器中获取到的未定义值导致输出不可预测。

4.1.3 实验中遇到的问题和解决办法

(1) 问题 1：此次实验的问题主要集中在理解汇编代码和调试输出结果方面。

解决方法：解析 RISC-V 汇编代码中的函数调用和寄存器传递参数的过程较为复杂，尤其是在编译器进行内联优化的情况下，很难直接从汇编代码中看到函数调用的痕迹。通过查阅 RISC-V 的相关文档和调试工具，我深入理解了寄存器的使用规则以及编译器优化的机制，最终正确识别出函数调用的位置和参数传递的方式。

4.1.4 实验心得

通过本次实验，我深入学习了 RISC-V 汇编语言的基础知识，并理解了编译器优化对汇编代码的影响。在分析和调试过程中，我强化了对寄存器操作的理解，尤其是函数参数传递和调用栈的管理。此外，本次实验让我意识到，在进行系统级编程时，理解硬件架构和指令集对代码执行的影响是非常重要的，这不仅有助于优化代码性能，还能提高代码的可预测性和稳定性。

4.2 Backtrace (moderate)

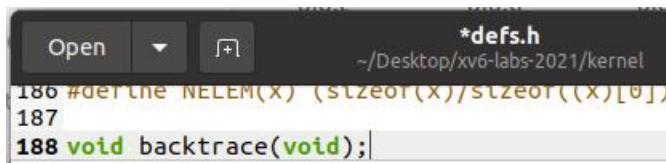
4.2.1 实验目的

本次实验旨在实现一个 Backtrace 功能，使得当内核检测到错误或异常时，能够自动输出当前执行上下文的调用栈信息。通过这种方式，开发者能够在不需要调试工具的情况下了解错误发生的具体位置及上下文，从而更高效地定位问题并解决问题。

4.2.2 实验步骤

(1) 定义和声明：

在 `defs.h` 文件中声明 `backtrace()` 函数。



```
Open  [icon] *defs.h
~/Desktop/xv6-labs-2021/kernel
186 #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
187
188 void backtrace(void);
```

(2) 定义帧指针读取函数:

在 `kernel/riscv.h` 中定义一个内联汇编函数 `r_fp()`，该函数用于获取当前函数的帧指针（frame pointer），即 `s0` 寄存器的值。



```
Open  [icon] *riscv.h Save
~/Desktop/xv6-labs-2021/kernel
368
369 static inline uint64
370 r_fp(){
371     uint64 x;
372     asm volatile("mv %0, s0" : "=r" (x));
373     return x;
374 }
```

(2) 实现 backtrace 函数:

在 `kernel/printf.c` 文件中实现了 `backtrace()` 函数。该函数使用一个循环遍历当前的栈帧，从高地址向低地址方向迭代，并输出每一个栈帧中保存的返回地址。

```
void backtrace(void)
{
    uint64 fp_address = r_fp(); // 获取当前帧指针

    while (fp_address != PGROUNDDOWN(fp_address)) {
        printf("%p\n", *(uint64*)(fp_address - 8)); // 输出返回地址
        fp_address = *(uint64*)(fp_address - 16); // 更新帧指针
    }
}
```

在函数中，我们首先获取当前帧指针，然后在一个循环中通过帧指针逐步遍历堆栈，输出每个栈帧中保存的返回地址。循环的终止条件是帧指针所在的页面起始地址（使用 `PGROUNDDOWN` 宏计算）。

(4) 在 `sys_sleep` 中调用 `backtrace`

修改 `sys_sleep` 函数，在其中调用 `backtrace()` 以便在特定条件下打印栈跟踪。


```
Open ▼ [icon] *sysproc.c
~/Desktop/xv6-labs-2021/kernel
72 release(&tickslock);
73 backtrace();
74 return 0;
```

(5) 编译与运行:

在终端中执行 `make qemu` 编译并运行 `xv6`。在命令行中输入 `bttest`，产生结果如下，符合预期。

```
hart 1 starting
hart 2 starting
init: starting sh
$ bttest
backtrace:
0x0000000080002144
0x0000000080001fa6
0x0000000080001c90
$
```

(6) 使用 `addr2line` 工具解析地址

在退出 `QEMU` 后,使用 `addr2line` 工具将 `backtrace` 输出的地址转换为相应的函数名和文件行号,以便定位错误位置。

```
addr2line -e kernel/kernel [address]
```

(7) 将 `backtrace` 函数集成到 `panic` 函数

为了在内核发生 `panic` 时能够自动输出调用堆栈信息,将 `backtrace` 函数添加到 `panic` 函数中。

```
void
panic(char *s)
{
    pr.locking = 0;
    printf("panic: ");
    printf(s);
    printf("\n");
    panicked = 1; // freeze uart output from other CPUs
    for(;;)
        ;
}
```

4.2.3 实验中遇到的问题和解决办法

(1) 问题 1: 在实现 `backtrace()` 函数时,需要确定合适的终止条件来避免无限循环。

解决办法：通过检查 `fp` 是否小于当前页面的最高地址（使用 `PGROUNDUP(fp)` 计算），这样可以确保循环在到达栈底时结束。

（2）问题 2：如何从输出的地址找到对应的源代码行号。

解决办法：使用 `riscv64-unknown-elf-addr2line` 工具，它可以将内存地址转换为源代码文件名和行号。将 `bttest` 的输出作为输入传递给该工具，可以得到清晰的调用栈信息。

4.2.4 实验心得

通过这次实验，我对 RISC-V 架构下栈的组织方式有了更深入的理解。特别是了解了如何利用寄存器 `s0` 来追踪函数调用的顺序，这对于调试内核级别的问题非常有用。此外，我还学习了如何使用 `addr2line` 这样的工具来辅助调试，这极大地提高了调试效率。

4.3 Alarm

4.3.1 实验目的

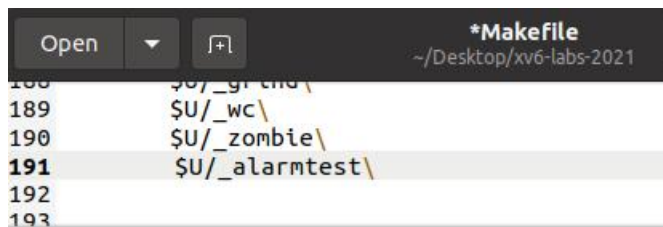
本次实验旨在为 `xv6` 内核增加一个定时提醒功能，使得用户进程可以在经过指定的 CPU 时间后执行一个预定义的函数。这个功能类似于在用户态下实现一个简化的中断处理机制，能够让进程在特定的时间间隔触发指定的操作。这种机制在处理需要周期性任务的应用程序（如定时任务、定期检查等）时尤为有用。

4.3.2 实验步骤

（1）修改内核以支持用户态警报处理函数：

- 添加必要的文件和函数声明：

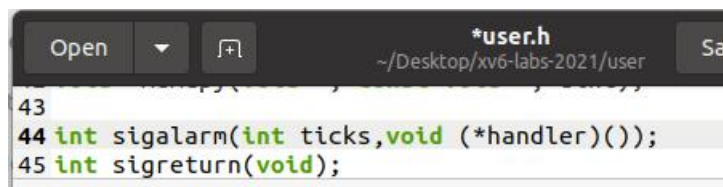
在 `Makefile` 中添加 `$U/_alarmtest`，以便将 `alarmtest.c` 编译为 `xv6` 用户程序。



```
*Makefile
~/Desktop/xv6-labs-2021

188  $U/_gr_end\
189  $U/_wc\
190  $U/_zombie\
191  $U/_alarmtest\
192
193
```

在 `user/user.h` 中添加以下函数声明



```
*user.h
~/Desktop/xv6-labs-2021/user

43
44 int sigalarm(int ticks, void (*handler)());
45 int sigreturn(void);
```

修改 `user/usys.pl` 文件，增加新的系统调用入口

```

Open  ▾  ↵  usys.pl
~/Desktop/xv6-labs-2021/u
37 entry("sleep");
38 entry("uptime");
39 entry("sigalarm");
40 entry("sigreturn");

```

- 系统调用和内核功能的实现：

在 `syscall.h` 中定义新的系统调用编号

```

Open  ▾  ↵  *syscall.h
~/Desktop/xv6-labs-2021/kernel
21 #define SYS_mkdir  20
22 #define SYS_close  21
23 #define SYS_sigalarm 22
24 #define SYS_sigreturn 23

```

在 `syscall.c` 中注册新的系统调用处理函数

```

Open  ▾  ↵  *syscall.c
~/Desktop/xv6-labs-2021/kernel
107 extern uint64 sys_sigalarm(void);
108 extern uint64 sys_sigreturn(void);
109
132 [SYS_sigalarm]    sys_sigalarm,
133 [SYS_sigreturn]   sys_sigreturn,
134 }:

```

修改 `proc.h` 文件，添加新的字段以支持警报功能

```

void (*alarm_handler)();

int alarm_interval;

int ticks;

int alarm_handling;

struct trapframe *saved_trapframe;

```

在 `sysproc.c` 中实现 `sys_sigalarm` 和 `sys_sigreturn` 的核心逻辑。
`sys_sigalarm` 用于设置警报的时间间隔和处理函数；`sys_sigreturn` 用于恢复被中断的执行状态。

```

// 定义一个函数 `sys_sigreturn`，返回类型为 `uint64`
uint64
sys_sigreturn(void)
{
    // 获取当前进程的指针，通过调用 `myproc()` 函数
    struct proc* proc = myproc();

```

```

// 重新存储陷阱帧(trapframe)，以便返回到中断发生之前的状态
*proc->trapframe = proc->savd_trapframe;

// 标记该进程已经处理过信号返回，设置为 1 表示 true
proc->have_return = 1;

// 返回当前进程的 `trapframe` 中的 `a0` 寄存器的值
// 这个值通常用于保存系统调用的返回值
return proc->trapframe->a0;
}

// 定义一个函数 `sys_sigalarm`，返回类型为 `uint64`
uint64
sys_sigalarm(void)
{
    int ticks;           // 定义一个整数变量 `ticks`，用于存储闹钟的间隔
    uint64 handler_va;   // 定义一个无符号 64 位整数变量 `handler_va`，用于
                        // 存储信号处理程序的虚拟地址

    // 获取系统调用的第一个参数（闹钟的时间间隔），存储到 `ticks` 中
    argint(0, &ticks);

    // 获取系统调用的第二个参数（信号处理程序的虚拟地址），存储到
    // `handler_va` 中
    argaddr(1, &handler_va);

    // 获取当前进程的指针，通过调用 `myproc()` 函数
    struct proc* proc = myproc();

    // 设置当前进程的闹钟时间间隔为 `ticks`

```

```

proc->alarm_interval = ticks;

// 设置当前进程的信号处理程序的虚拟地址为 `handler_va`
proc->handler_va = handler_va;

// 标记该进程已经设置了闹钟, `have_return` 为 1 表示 true
proc->have_return = 1;

// 返回 0, 表示成功设置了闹钟
return 0;
}

```

- 修改中断处理逻辑:

在 `kernel/trap.c` 的 `usertrap()` 函数中, 增加处理逻辑, 使每次时钟中断都会检查是否需要触发警报处理函数。

```

// 检查 `which_dev` 是否等于 2
if (which_dev == 2) {
    // 获取当前进程的指针, 通过调用 `myproc()` 函数
    struct proc *proc = myproc();

    // 检查当前进程的闹钟间隔 `alarm_interval` 是否非零且 `have_return` 是否为真
    if (proc->alarm_interval && proc->have_return) {
        // 将 `passed_ticks` 计数器加 1
        if (++proc->passed_ticks == proc->alarm_interval) {
            // 如果已经过的 ticks 达到闹钟间隔, 则:

            // 保存当前的陷阱帧到 `saved_trapframe` 中
            // `saved_trapframe` 用于在信号处理完成后恢复进程的状态
            proc->saved_trapframe = *proc->trapframe;

```

```

// 将 `trapframe` 的 `epc` 寄存器设置为信号处理程序的虚
拟地址

// `epc` 寄存器存储了下一条将要执行的指令的地址
proc->trapframe->epc = proc->handler_va;

// 将 `passed_ticks` 重置为 0
proc->passed_ticks = 0;

// 将 `have_return` 设置为 0，表示信号处理程序已经被设置
proc->have_return = 0;
    }
}
}

```

（2）确保执行恢复与定时器重置

- 实现 `sys_sigreturn` 函数：

`sys_sigreturn` 的作用是恢复在警报处理程序触发时保存的用户态寄存器状态，使用户进程可以继续从中断点执行。

```

uint64
sys_sigreturn(void)
{
    struct proc* proc = myproc();

    *proc->trapframe = proc->savd_trapframe;

    proc->have_return = 1; // 标记处理函数已返回

    return proc->trapframe->a0;
}

```

（3）编译与运行：

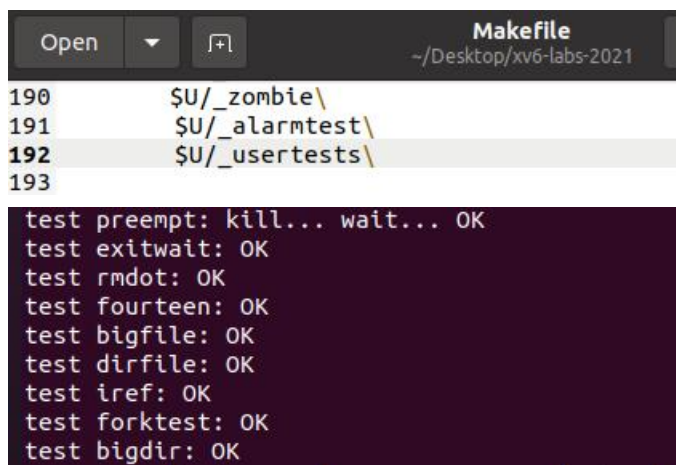
在终端中执行 `make qemu` 编译并运行 `xv6`。在命令行中输入 `alarmtest`，产生结果如下，符合预期。

```

hart 1 starting
hart 2 starting
init: starting sh
$ alarmtest
test0 start
.....alarm!
test0 passed
test1 start
.alarm!
..alarm!
.alarm!
..alarm!
....alarm!
..alarm!
....alarm!
...alarm!
.....alarm!
...alarm!
.test1 passed
test2 start
.....alarm!
test2 passed

```

在 Makefile 中添加 `$U/_usertests`, 运行用户测试程序以确保内核的其他部分没有被意外修改或破坏。



The screenshot shows a text editor window titled "Makefile" with the path "~/Desktop/xv6-labs-2021". The editor contains a list of user tests being added to the Makefile:

```

190     $U/_zombie\
191     $U/_alarmtest\
192     $U/_usertests\
193

```

Below the editor window, a terminal window shows the output of the tests:

```

test preempt: kill... wait... OK
test exitwait: OK
test rmdot: OK
test fourteen: OK
test bigfile: OK
test dirfile: OK
test iref: OK
test forktest: OK
test bigdir: OK

```

4.3.3 实验中遇到的问题和解决办法

(1) 问题 1: 初次实现时, 由于对用户态和内核态之间的切换机制理解不深, 导致漏掉了对用户层面的函数声明和入口设置。

解决方法: 通过查阅 xv6 文档和代码, 逐步理解了系统调用的完整流程, 并修正了这些问题。

(2) 问题 2: 在实现过程中, 没有正确地保存和恢复中断时的寄存器状态, 导致用户进程无法从中断恢复。

解决方案: 通过仔细阅读 `trapframe` 结构的定义, 理解了在中断时如何保存和恢复寄存器状态, 最终解决了恢复的问题。

(3) 问题 3: 出现超时, 导致实验成绩未达标

解决方法：在 grade-lab-traps 中设置时间 300 为 800。同时在 param.h 中设置 1000 为 10000。

4.3.4 实验心得

本次实验让我深入理解了操作系统的定时中断和用户态中断处理机制。通过在 xv6 内核中添加定时提醒功能，我掌握了如何在内核态和用户态之间进行信息传递，以及如何在中断处理程序中实现特定的功能。此外，本次实验也让我体会到，在修改内核代码时，确保系统的稳定性和正确性是至关重要的，任何一个细微的错误都有可能导致系统崩溃或不稳定。因此，深入理解系统调用和中断处理的原理，做好充分的测试和验证，是确保内核功能扩展成功的关键。

4.4 Lab4 实验成绩

```
== Test answers-traps.txt == answers-traps.txt: OK
== Test backtrace test ==
$ make qemu-gdb
backtrace test: OK (3.0s)
== Test running alarmtest ==
$ make qemu-gdb
(3.9s)
== Test  alarmtest: test0 ==
alarmtest: test0: OK
== Test  alarmtest: test1 ==
alarmtest: test1: OK
== Test  alarmtest: test2 ==
alarmtest: test2: OK
== Test usertests ==
$ make qemu-gdb
usertests: OK (495.3s)
(Old xv6.out.usertests failure log removed)
== Test time ==
time: OK
Score: 85/85
```

实验小结：

在完成章节 4 的实验后，我对 RISC-V 架构下的陷阱处理有了深刻的认识。从学习 RISC-V 汇编语言的底层指令，到实现回溯功能以追踪调用路径，再到设置定时提醒功能处理系统调用，我经历了复杂但充实的过程。RISC-V 汇编的挑战让我熟悉了指令集和异常处理，而回溯功能加深了我对系统状态恢复的理解。定时提醒功能的实现则强化了内核和用户态交互的知识。这些实验不仅提高了我的编程技能，也加深了对操作系统设计的理解，为今后的系统编程打下了坚实的基础。

5. Lab5: Copy-on-Write Fork for xv6

切换到 `cow` 分支:

```
git fetch
```

```
git checkout cow
```

```
make clean
```

```
lleell@ubuntu:~/Desktop$ cd xv6-labs-2021
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git checkout cow
Branch 'cow' set up to track remote branch 'cow' from 'origin'.
Switched to a new branch 'cow'
lleell@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
*/*.o */*.d */*.asm */*.sym \
user/initcode user/initcode.out kernel/kernel fs.img \
mkfs/mkfs .gdbinit \
    user/usys.S \
user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln u
ser/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_grin
d user/_wc user/_zombie user/_cowtest \
ph barrier
```

5.1 Implement copy-on write

5.1.1 实验目的

传统的 `fork()` 系统调用会将父进程的整个用户空间内存复制到子进程，导致资源浪费和效率低下。通过 COW `fork`，我们只在需要写入时才进行物理页面的分配和复制，这样可以显著减少内存消耗和提高系统性能。本实验旨在通过实践掌握 COW 机制的工作原理，了解如何在操作系统中实现内存共享与写时复制技术。

5.1.2 实验步骤

(1) 设置 COW 标记位:

在 `kernel/riscv.h` 中定义一个新的页表项标志位 `PTE_RSW` 用于标记 COW 页面。`#define PTE_RSW (1L << 8) // RSW` 这个标记位将用于区分 COW 页面，便于在页表项中记录该信息。

```

Open  ▾  [+l]  *riscv.h  ~/Desktop/xv6-labs-2021/kernel  Save  ≡
340
341 #define PTE_V (1L << 0) // valid
342 #define PTE_R (1L << 1)
343 #define PTE_W (1L << 2)
344 #define PTE_X (1L << 3)
345 #define PTE_U (1L << 4) // 1 -> user can access
346 #define PTE_COW (1L << 8) |
347

```

(2) 修改 `uvmcopy` 函数：

在 `kernel/vm.c` 文件中，修改 `uvmcopy()` 函数，使其不再复制父进程的物理内存页面，而是仅复制父进程的页表，将父进程的物理页面映射到子进程中。

```

Open  ▾  [+l]  *defs.h  ~/Desktop/xv6-labs-2021/kernel  Save
62 // kalloc
63 void*      kalloc(void);
64 void      kfree(void *);
65 void      kinit(void);
66 void      adjustref(uint64, int);
67

```

清除父子进程页表项的 `PTE_W`（写标志）并设置 `PTE_RSW` 标志，标记这些页面为只读的 `COW` 页面。

```

// uvmcopy 函数将父进程的页表复制到子进程中
// 同时将物理页面标记为只读，使用 COW（写时复制）机制
void uvmcopy(pagetable_t old, pagetable_t new, uint64 sz) {
    pte_t* pte;
    uint64 pa, i;
    uint flags;

    // 遍历每一页
    for (i = 0; i < sz; i += PGSIZE) {
        // 获取 old 页表中虚拟地址 i 的页表项
        if ((pte = walk(old, i, 0)) == 0)
            panic("uvmcopy: pte should exist"); // 页表项不存在则出错
        if ((*pte & PTE_V) == 0)
            panic("uvmcopy: page not present"); // 页面不存在则出错

        // 获取物理地址
        pa = PTE2PA(*pte);

```

```

        // 清除写标志，将页面标记为只读
        *pte &= ~PTE_W;
        flags = PTE_FLAGS(*pte);
        // 在新页表中映射物理页面
        if (mappages(new, i, PGSIZE, pa, flags) != 0) {
            goto err;
        }
        // 调整引用计数，增加 1
        adjustref(pa, 1);
    }
    return;

err:
    // 如果出现错误，解除映射并返回-1
    uvmunmap(new, 0, i / PGSIZE, 1);
    return -1;
}

// adjustref 函数调整物理页面的引用计数
void adjustref(uint64 pa, int num) {
    if (pa >= PHYSTOP) {
        panic("addr: pa too big"); // 物理地址过大则出错
    }
    acquire(&kmem.lock);
    cowcount[PA2INDEX(pa)] += num; // 增加或减少引用计数
    release(&kmem.lock);
}

```

(3) 页面错误处理:

修改 `kernel/trap.c` 文件中的 `usertrap()` 函数，增加对页面错误的处理逻辑。

```
// 用户态陷阱处理函数

void usertrap(void) {
    // 处理页面错误（异常码 15）
    if (r_scause() == 15) {
        if (cowalloc(p->pagetable, r_stval()) < 0) {
            p->killed = 1; // 分配页面失败，标记进程为被杀死
        }
    }
    // 其他处理逻辑...
}
```

当页面错误是由于写入 COW 页面引起时，调用 `cowfault()` 函数进行处理。在该函数中，为发生错误的进程分配一个新的物理页面，将原页面的内容复制到新页面，然后更新页表项以指向新页面，并设置为可写。



（4）管理引用计数：

在 `kernel/kalloc.c` 中，为每个物理页面维护一个引用计数，记录每个物理页面被多少个页表项引用。

修改 `kalloc()` 函数，在分配新页面时将引用计数初始化为 1。

```
// kalloc 函数分配一个新的物理页面

// 初始化新页面的引用计数为 1

void* kalloc(void) {
    struct run *r;

    // 申请内存

    if ((r = kmem.freelist) != 0) {
        kmem.freelist = r->next;

        memset((char *)r, 5, PGSIZE); // 填充垃圾数据
```

零

```
int idx = PA2INDEX(r);

if (cowcount[idx] != 0) {

    panic("kalloc: cowcount[idx] != 0"); // 页面的引用计数不应非

}

cowcount[idx] = 1; // 新分配的页面的引用计数为 1

}

return (void *)r;

}

// kfree 函数释放一个物理页面

void kfree(void *pa) {

    struct run *r;

    if (((uint64)pa % PGSIZE) != 0 || (char *)pa < end || (uint64)pa >=
PHYSTOP)

        panic("kfree");

    acquire(&kmem.lock);

    int remain = --cowcount[PA2INDEX(pa)]; // 减少引用计数

    release(&kmem.lock);

    if (remain > 0) {

        // 如果引用计数大于 0，则不释放页面

        return;

    }

    // 释放页面

    r = (struct run *)pa;

    r->next = kmem.freelist;
```

```

        kmem.freelist = r;
    }

    // freerange 函数初始化物理页面范围
    void freerange(void *pa_start, void *pa_end) {
        char *p;
        p = (char *)PGROUNDUP((uint64)pa_start);
        for (; p + PGSIZE <= (char *)pa_end; p += PGSIZE) {
            cowcount[PA2INDEX(p)] = 1; // 初始化每个页面的引用计数为 1
            kfree(p);
        }
    }
}

```

(5) 处理 COW 页面在 `copyout()` 函数中的操作:

修改 `kernel/vm.c` 中的 `copyout()` 函数, 处理写入 COW 页面的情况。如果目标页面为 COW 页面, 函数会触发页面错误处理逻辑, 分配新的物理页面并复制原数据。

```

// copyout 函数将用户数据拷贝到目标虚拟地址
// 处理写入 COW 页面时的情况
int copyout(pagetable_t pagetable, uint64 dstva, char *src, uint64 len) {
    uint64 n, va0, pa0;

    while (len > 0) {
        va0 = PGROUNDDOWN(dstva); // 获取虚拟地址所在页面的起始
地址
        if (va0 >= MAXVA) {
            printf("copyout: va exceeds MAXVA\n");
            return -1;
        }
        pte_t *pte = walk(pagetable, va0, 0); // 获取页表项
    }
}

```

```

        if (pte == 0 || (*pte & PTE_U) == 0 || (*pte & PTE_V) == 0) {
            printf("copyout: invalid pte\n");
            return -1;
        }
        if ((*pte & PTE_W) == 0) {
            // 如果页面为 COW 共享页，则需要进行写时复制
            if (cowalloc(pagetable, va0) < 0) {
                return -1;
            }
        }
        pa0 = walkaddr(pagetable, va0); // 获取物理地址
        if (pa0 == 0)
            return -1;
        n = PGSIZE - (dstva - va0);
        if (n > len)
            n = len;
        memmove((void *)(pa0 + (dstva - va0)), src, n); // 拷贝数据

        len -= n;
        src += n;
        dstva = va0 + PGSIZE; // 移动到下一页
    }
    return 0;
}

```

(7) 定义 `cow_alloc()` 函数分配新的物理页：

```

// cowalloc 函数处理写时复制的页面分配
int cowalloc(pagetable_t pagetable, uint64 va) {
    if (va >= MAXVA) {
        printf("cowalloc: exceeds MAXVA\n");
        return -1;
    }
}

```

```

    }

    pte_t* pte = walk(pagetable, va, 0);
    if (pte == 0) {
        panic("cowalloc: pte not exists"); // 页表项不存在则出错
    }
    if ((*pte & PTE_V) == 0 || (*pte & PTE_U) == 0) {
        panic("cowalloc: pte permission err"); // 页表项权限错误则出错
    }
    uint64 pa_new = (uint64)kalloc(); // 分配新的物理页面
    if (pa_new == 0) {
        printf("cowalloc: kalloc fails\n");
        return -1;
    }
    uint64 pa_old = PTE2PA(*pte);
    memmove((void *)pa_new, (const void *)pa_old, PGSIZE); // 复制旧页面
    内容到新页面
    kfree((void *)pa_old); // 释放旧页面
    *pte = PA2PTE(pa_new) | PTE_FLAGS(*pte) | PTE_W; // 更新页表项，
    设置为可写
    return 0;
}

```

(7) 编译与运行:

在终端中执行 `make qemu` 编译并运行 `xv6`。在命令行中输入 `cowtest`，产生结果如下，符合预期。

```

hart 1 starting
hart 2 starting
init: starting sh
$ cowtest
simple: ok
simple: ok
three: ok
three: ok
three: ok
file: ok
ALL COW TESTS PASSED

```


在命令行中输入 `usertests`，产生结果如下，符合预期。

```
test throot: OK
test fourteen: OK
test bigfile: OK
test dirfile: OK
test iref: OK
test forktest: OK
test bigdir: OK
ALL TESTS PASSED
```

5.1.3 实验中遇到的问题和解决办法

(1) 问题 1：在初次实现时，未能正确管理物理页面的引用计数，导致某些页面在没有实际引用时未被释放，从而引发内存泄漏问题。

解决办法：通过在 `kalloc.c` 中为每个物理页面增加引用计数的管理，确保每次页面释放前都检查其引用计数，只有当计数为 0 时才真正释放内存。

(2) 问题 2：在页面错误处理过程中，有时因未正确检测 COW 页面，导致误判或系统崩溃。

解决办法：增加了对 PTE 标志的严格检查，确保仅在符合 COW 条件时进行页面分配和复制。

(3) 问题 3：在多个进程同时操作共享页面时，可能发生引用计数竞争条件，导致数据不一致或系统崩溃。

解决办法：在 `kalloc.c` 中引入锁机制，确保引用计数的更新是原子的，从而避免竞争条件。

5.1.4 实验心得

通过本次实验，我深入理解了 COW 机制的实现原理及其在操作系统中的应用。COW 通过延迟物理内存的分配和复制，优化了进程创建的性能，节省了系统资源。实现过程中遇到的内存管理问题和并发控制问题，使我更加明白操作系统开发中细节处理的重要性。同时，这次实验也让我体会到了现代操作系统中内存管理的复杂性，以及合理设计数据结构和同步机制对系统稳定性的关键作用。通过本次实验，我学到了许多实际编程技巧和调试经验，为今后的系统开发奠定了良好的基础。

5.2 Lab5 实验成绩

在 `xv6` 目录下添加 `time.txt` 文本文件，写入完成该实验的小时数。在终端中执行 `make grade`，即可对整个实验进行自动评分。评估结果如下，测试通过，程序的行为符合预期。

```

make[1]: Leaving directory '/home/cteeit/desktop/xv6-lab
== Test running cowtest ==
$ make qemu-gdb
(7.4s)
== Test    simple ==
    simple: OK
== Test    three ==
    three: OK
== Test    file ==
    file: OK
== Test usertests ==
$ make qemu-gdb
(419.3s)
    (Old xv6.out.usertests failure log removed)
== Test    usertests: copyin ==
    usertests: copyin: OK
== Test    usertests: copyout ==
    usertests: copyout: OK
== Test    usertests: all tests ==
    usertests: all tests: OK
== Test time ==
    time: OK
Score: 110/110

```

实验小结：

在本实验中，我们为 xv6 实现了一个基于写时复制（Copy-on-Write, COW）的 fork 机制。传统的 fork 系统调用会复制父进程的整个地址空间，导致内存资源浪费和效率低下。在实验中，我们改进了 xv6 的 fork 实现，使其在初始时不复制父进程的内存页，而是让父子进程共享这些页，并将页标记为只读。当任一进程试图写入这些共享页时，COW 机制会触发页面错误处理器，真正复制内存页，从而避免不必要的内存拷贝，提高系统性能。

通过这个实验，我们深入理解了 COW 机制的实现细节，包括页面表管理、页面错误处理和内存管理中的细微差别。实验结果表明，COW fork 能显著减少内存消耗并提升进程创建的效率，这对于多进程操作系统的优化具有重要意义。通过本实验，我们不仅掌握了 xv6 的内核结构，也增强了对操作系统内存管理机制的理解。

6. Lab6: Multithreading

这个 lab 主要由三部分组成：

- 实现一个用户级线程的创建和切换

- 使用 UNIX pthread 线程库实现一个线程安全的 Hash 表

- 利用 UNIX 的锁和条件变量实现一个 barrier

切换到 thread 分支：

```
git fetch
```

```
git checkout thread
```

```
make clean
```

```
lleell@ubuntu:~/Desktop$ cd xv6-labs-2021
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git checkout thread
Branch 'thread' set up to track remote branch 'thread' from 'origin'.
Switched to a new branch 'thread'
lleell@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
*/*.o */*.d */*.asm */*.sym \
user/initcode user/initcode.out kernel/kernel fs.img \
mkfs/mkfs .gdbinit \
    user/usys.S \
user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln u
ser/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_grin
d user/_wc user/_zombie user/_uthread \
ph barrier
```

6.1 Uthread: switching between threads

6.1.1 实验目的

本实验的主要目的是设计和实现一个用户级线程系统的上下文切换机制。目标是在用户态下管理线程，而不依赖内核进行调度。我们需要实现线程的创建和上下文切换功能，这涉及到保存和恢复线程的寄存器状态，确保不同线程之间的切换能够正确进行并通过测试。

6.1.2 实验步骤

(1) 理解现有代码：

阅读已给出的 `uthread.c` 和 `uthread_switch.S` 文件，理解其中已有的线程操作代码和基础测试代码，为后续实现线程切换功能打下基础。

设计线程数据结构，在 `uthread.c` 中创建一个数据结构 `uthread`，用来表示线程的状态。该结构体包含了保存线程上下文的寄存器信息以及线程栈：

```
Open uthread.c ~/Desktop/xv6-labs-2021-Lab6/user Save
12
13 // Saved registers for kernel context switches.
14 struct uthread_context {
15     uint64 ra;
16     uint64 sp;
17
18     // callee-saved
19     uint64 s0;
20     uint64 s1;
21     uint64 s2;
22     uint64 s3;
23     uint64 s4;
24     uint64 s5;
25     uint64 s6;
26     uint64 s7;
27     uint64 s8;
28     uint64 s9;
29     uint64 s10;
30     uint64 s11;
31 };
32
```

(2) 实现线程创建:

在 `uthread.c` 中完成 `thread_create()` 函数。在此函数中, 设置线程的上下文: 将寄存器 `ra` 设置为要执行的函数的地址, 将 `sp` 寄存器设置为线程栈的底部, 确保线程在切换时, 能够正确执行并使用独立的栈空间。

```
void
thread_create(void (*func)()) // thread_create 函数用于创建新线程
{
    struct thread *t;

    // 遍历线程池 (all_thread 数组), 找到一个空闲的线程结构 (状态为 FREE 的线程)
    for (t = all_thread; t < all_thread + MAX_THREAD; t++) {
        if (t->state == FREE) break; // 找到空闲线程, 退出循环
    }

    // 设置找到的线程的状态为 RUNNABLE, 表示该线程可以被调度执行
    t->state = RUNNABLE;

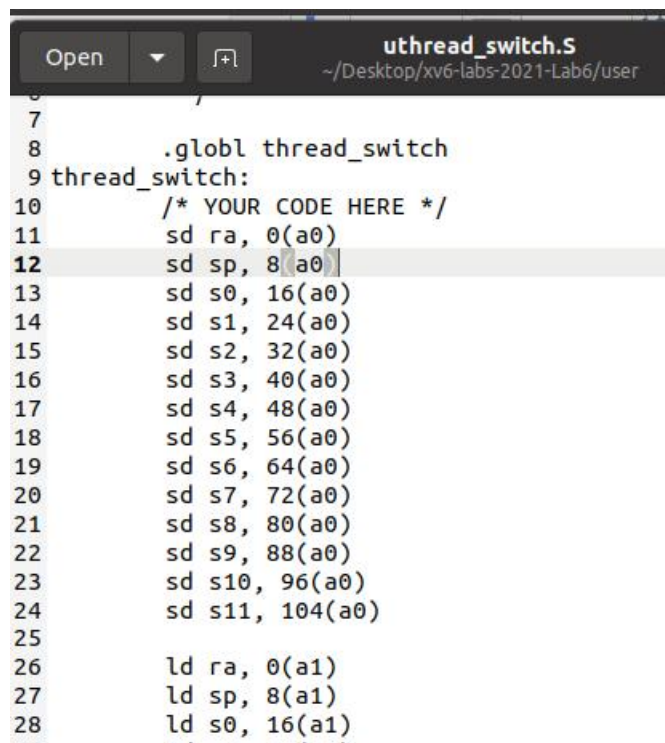
    // 将该线程的返回地址寄存器 (ra) 设置为参数 func, 即要执行的函数的入口地址
    // 当线程被调度运行时, 会从此函数地址开始执行
    t->context.ra = (uint64)func;
```

```
// 设置该线程的栈指针（sp）寄存器，指向线程栈的顶部（栈是向下增长的，所以 sp 指向栈顶） // 这样线程在执行时有自己的独立栈空间，防止与其他线程的栈混淆
```

```
t->context.sp = (uint64)(t->stack + STACK_SIZE);  
  
}
```

（3）实现线程切换：

在 `uthread_switch.S` 中实现 `thread_switch` 函数。该函数负责保存当前线程的上下文（即寄存器状态），然后切换到下一个线程并恢复其上下文。可以参考 `kernel/switch.S` 文件的写法，按照 `struct context` 中寄存器在内存中的位置进行保存和恢复。



```
7  
8     .globl thread_switch  
9 thread_switch:  
10    /* YOUR CODE HERE */  
11    sd ra, 0(a0)  
12    sd sp, 8(a0)  
13    sd s0, 16(a0)  
14    sd s1, 24(a0)  
15    sd s2, 32(a0)  
16    sd s3, 40(a0)  
17    sd s4, 48(a0)  
18    sd s5, 56(a0)  
19    sd s6, 64(a0)  
20    sd s7, 72(a0)  
21    sd s8, 80(a0)  
22    sd s9, 88(a0)  
23    sd s10, 96(a0)  
24    sd s11, 104(a0)  
25  
26    ld ra, 0(a1)  
27    ld sp, 8(a1)  
28    ld s0, 16(a1)
```

（4）实现调度器：

在 `uthread.c` 中编写 `thread_schedule()` 函数，负责调用 `thread_switch` 来切换线程。函数需要正确传递当前线程和下一个线程的上下文指针给 `thread_switch`，实现线程之间的切换。

```
// thread_schedule 函数实现调度器的功能，负责调用 thread_switch 进行线程切换
```

```
// 检查是否存在可以运行的线程
```

```
if(next_thread == 0) {  
  
    printf("thread_schedule: no runnable threads\n");  
  
    exit(-1); // 没有可运行的线程，退出程序
```

```

    }

    // 如果存在可运行的线程，且当前线程不是即将运行的线程

    if (current_thread != next_thread) {          /* 是否需要切换线程？ */

        next_thread->state = RUNNING;// 将即将运行的线程状态设置为
RUNNING

        t = current_thread;// 保存当前线程指针到 t

        current_thread = next_thread;// 更新当前线程指针为即将运行的线程


        /* YOUR CODE HERE

        * 调用 thread_switch 从当前线程 t 切换到下一个线程 next_thread
        * 需要传递两个参数：当前线程的上下文地址和下一个线程的上下文
地址
        */

        thread_switch((uint64)&t->context, (uint64)&next_thread->context);

    }

```

(5) 编译与运行：

在终端中执行 `make qemu` 编译并运行 xv6。在命令行中输入 `uthread`，产生结果如下，符合预期。

```

thread_c: exit after 100
thread_a: exit after 100
thread_b: exit after 100
thread_schedule: no runnable threads

```

(6) 评估与测试：

使用 xv6 自带的测评工具 `./grade-lab-thread uthread` 来进行自动测试。评估结果如下，测试通过，程序的行为符合预期。

```

llee11@ubuntu:~/Desktop/xv6-labs-2021$ ./grade-lab-thread uthread
make: 'kernel/kernel' is up to date.
== Test uthread == uthread: OK (1.8s)

```

6.1.3 实验中遇到的问题和解决办法

(1) 问题 1：在创建线程时，正确分配和管理线程的栈空间是关键。初始实现时，因没有正确设置线程的栈指针 `sp`，导致在执行线程时发生错误。

解决方案：为解决这个问题，我查看了 RISC-V 架构文档，明确了 `sp` 应该指向栈底（较高地址），确保栈的生长方向正确。

(2) 问题 2: 参数传递错误, 在实现 `thread_switch` 时, 起初传递的上下文指针有误, 导致无法正确切换线程。

解决方案: 通过使用 `gdb` 调试工具, 我设置了断点, 逐步检查寄存器状态, 发现问题所在。最终, 通过正确传递参数指针, 解决了这个问题。

6.1.4 实验心得

通过本实验, 我深入理解了用户级线程系统的基本工作原理, 特别是线程上下文的保存与恢复机制。实现过程中, 使用汇编代码编写上下文切换功能让我更好地理解了寄存器的使用和栈的管理。同时, 掌握了如何在 **RISC-V** 架构下进行寄存器操作和指针管理, 有助于理解操作系统底层的线程管理机制。实验中的问题虽然复杂, 但通过仔细调试和查阅文档, 我成功地解决了这些问题, 最终实现了预期功能。通过本实验, 我加深了对多线程系统的理解, 为将来更复杂的操作系统设计打下了坚实基础。

6.2 Using threads

6.2.1 实验目的

本实验的目的是通过使用 **POSIX** 线程库 (`pthread`) 实现多线程编程, 并在多线程环境下操作哈希表。本实验中, 我们将学习如何创建和管理线程, 以及如何使用锁来保护共享资源, 从而实现线程安全的哈希表, 确保在多线程环境中的数据一致性和程序性能。

6.2.2 实验步骤

(1) 了解 `pthread` 库的基本概念:

在实验开始之前, 我们首先需要阅读有关 `pthread` 库的文档, 了解基本概念、函数和用法。这将帮助我们正确地使用线程和锁来实现并行编程。

(2) 构建和运行实验代码:

使用提供的代码文件构建实验程序。在命令行中运行 `make ph` 来编译 `notxv6/ph.c`。

运行 `./ph 1`, 验证在单线程模式下写入哈希表的数据被完整地读出, 没有遗漏。

```
lleell@ubuntu:~/Desktop/xv6-labs-2021$ make ph
make: 'ph' is up to date.
lleell@ubuntu:~/Desktop/xv6-labs-2021$ ./ph 1
100000 puts, 7.246 seconds, 13801 puts/second
0: 0 keys missing
100000 gets, 7.297 seconds, 13705 gets/second
```

运行./ph 2，测试在多线程模式下的性能和正确性。观察输出，注意在多线程情况下可能会出现缺少的键（missing keys）问题，即一些数据没有被正确写入到哈希表中。

```
lleell@ubuntu:~/Desktop/xv6-labs-2021$ ./ph 2
100000 puts, 5.204 seconds, 19217 puts/second
0: 16728 keys missing
1: 16728 keys missing
200000 gets, 16.733 seconds, 11952 gets/second
```

（3）分析多线程问题：

分析多线程环境中缺少的键问题，找出可能导致问题的竞争条件。假设两个线程同时向哈希表中添加条目，如果两个键的哈希值相同，它们会尝试插入相同的位置，导致其中一个键被覆盖，造成数据丢失。

```
29 static void
30 insert(int key, int value, struct entry **p, struct entry *n)
31 {
32     struct entry *e = malloc(sizeof(struct entry));
33     e->key = key;
34     e->value = value;
35     e->next = n;
36     *p = e;
37 }
38
```

上面这段代码是多线程操作哈希表时出错的根源，假设某个线程调用了insert 但没有返回，此时另一个线程调用 insert，它们的第四个参数 n(bucket 的链表头) 如果值相同，就会发生漏插入键值对的现象。

（4）添加锁来保护共享资源：

因此，我们需要定义一个互斥锁 pthread_mutex_t 来保护哈希表的访问。

```
Open  ▼  [+]
```

ph.c
~/Desktop/xv6-labs-2021-Lab6/notesxv6 Save

```
19
20 pthread_mutex_t lock[NBUCKET]; //每个桶分配一个锁
```

在 main 函数中使用 pthread_mutex_init 初始化锁。

```
Open  ▼  [+]
```

ph.c
~/Desktop/xv6-labs-2021-Lab6/notesxv6 Save

```
25 for (int i = 0; i < NBUCKET; ++i)
26     pthread_mutex_init(&lock[i], NULL);
27
28 //
```

在调用 insert 函数之前，使用 pthread_mutex_lock 上锁，确保 insert 操作的原子性；操作完成后，使用 pthread_mutex_unlock 解锁。


```
ph.c
~/Desktop/xv6-labs-2021-Lab6/notxv6
Save

50 }
51
52 pthread_mutex_lock(&lock[i]);
53 if(e){
54     // update the existing key.
55     e->value = value;
56 } else {
57     // the key is new.
58     insert(key, value, &table[i], table[i]);
59 }
60 pthread_mutex_unlock(&lock[i]);
61
62 }
63
```

使用完毕后，使用 `pthread_mutex_destroy` 销毁锁，防止占用系统资源。

```
ph.c
~/Desktop/xv6-labs-2021-Lab6/notxv6
Save

139
140 for (int i = 0; i < NBUCKET; ++i)
141     pthread_mutex_destroy(&lock[i]);
142
```

(5) 编译与运行：

在终端中执行 `make ph` 编译并运行 `xv6`。在命令行中重新输入 `./ph 1` 和 `./ph 2`，产生结果如下，keys missing 在单线程和多线程下都为 0，符合预期。

```
lleell@ubuntu:~/Desktop/xv6-labs-2021$ make ph
gcc -o ph -g -O2 -DSOL_THREAD -DLAB_THREAD notxv6/ph.c -pthread
lleell@ubuntu:~/Desktop/xv6-labs-2021$ ./ph 1
100000 puts, 7.217 seconds, 13856 puts/second
0: 0 keys missing
100000 gets, 6.943 seconds, 14403 gets/second
lleell@ubuntu:~/Desktop/xv6-labs-2021$ ./ph 2
100000 puts, 3.387 seconds, 29525 puts/second
0: 0 keys missing
1: 0 keys missing
200000 gets, 6.599 seconds, 30307 gets/second
lleell@ubuntu:~/Desktop/xv6-labs-2021$
```

(6) 优化性能：

通过为每个哈希桶设置独立的锁，减少锁的粒度，最大化并行性能。这样可以在不保护重叠区域时实现更高效的并行操作。

6.2.3 实验中遇到的问题和解决办法

(1) 问题 1：在初次实现锁机制时，由于锁的使用不当导致了死锁，即多个线程都在等待锁的释放而无法继续执行，而产生死锁。

解决方法：确保在每次获取锁之后，代码都能正常释放锁。同时，使用 `pthread_mutex_destroy` 销毁锁，避免未销毁的锁占用系统资源。

(2) 问题 2：使用过多的锁会导致性能下降，因为锁会限制多线程的并行执行。

解决方法：是将锁的粒度缩小，为每个哈希桶设置独立的锁，避免不必要的锁竞争，提升并行性和性能。

6.2.4 实验心得

通过本次实验，我深刻地体会到多线程编程的复杂性和重要性。在多线程环境中，线程安全是一个非常关键的问题。使用锁可以确保在访问共享资源时的线程安全性，防止竞争条件和数据不一致的问题。我学会了如何正确使用 `pthread` 库中的锁机制，并且在使用锁时需要遵循严格的获取和释放顺序，以避免死锁。同时，我还了解到，合理的锁管理对于程序性能的影响非常大。通过优化锁的使用策略，如缩小锁的粒度和使用独立的锁来减少锁竞争，可以显著提升程序的并行性能。

6.3 Barrier

6.3.1 实验目的

本实验的主要目标是通过实现一个线程屏障（Barrier）来加深对多线程编程中同步和互斥机制的理解。在多线程环境下，线程屏障可以确保所有参与的线程在到达特定的同步点后都等待，直到所有线程都到达此点才能继续执行。这种机制在需要多个线程在某个时间点协调或同步时非常有用。本实验旨在利用 `pthread` 库的条件变量和互斥锁来实现一个线程屏障，解决多线程编程中常见的竞争条件和同步问题。

6.3.2 实验步骤

（1）理解 Barrier 的基本原理：

线程屏障用于在某个预定点同步多个线程的执行。每个线程在到达屏障时都会进入等待状态，直到所有其他线程也到达屏障为止。只有当所有线程都到达屏障，所有线程才会被唤醒，继续执行接下来的任务。

```
pthread_cond_wait(&cond, &mutex);  
  
// 根据条件休眠，释放锁互斥，唤醒后获取  
  
pthread_cond_broadcast(&cond);  
  
// 唤醒在 cond 上每个处于睡眠状态的线程
```

（2）学习 barrier.c 文件的结构：

首先阅读 `barrier.c` 文件，理解其基本结构和已有函数，如 `barrier_init()`，该函数用于初始化屏障结构体 `struct barrier` 的状态。熟悉结构体的成员变量，这些变量用于存储当前的线程计数、同步轮次和用于同步的互斥锁与条件变量。

(2) 分析并实现 `barrier()` 函数：

在 `barrier.c` 文件中的 `barrier()` 函数内添加逻辑，使其在每个线程到达屏障后进行正确的同步操作：

使用 `pthread_mutex_lock()` 加锁以保护共享资源。

每个线程进入 `barrier()` 函数时，将计数器 `bstate.nthread` 加 1。

如果未达到所需的线程数 `nthread`，调用 `pthread_cond_wait()` 等待其他线程到达。

当最后一个线程到达时，重置 `bstate.nthread`，增加轮次 `bstate.round`，并通过 `pthread_cond_broadcast()` 唤醒所有在等待的线程。

最后，通过 `pthread_mutex_unlock()` 释放锁。

```
static void
barrier()
{
    // YOUR CODE HERE

    //

    // Block until all threads have called barrier() and
    // then increment bstate.round.

    //

    // 上锁

    pthread_mutex_lock(&bstate.barrier_mutex);

    // 来了一个线程，累计 ++
    ++ bstate.nthread;

    // 查看条件变量
    if (bstate.nthread != nthread) {
        // 放弃锁并且睡眠，等待条件变量来唤醒

        pthread_cond_wait(&bstate.barrier_cond, &bstate.barrier_mutex);
    } else {
        // 增加 barrier 里的计数 (只计一次)

        ++ bstate.round;
```

```

        // 清零

        bstate.nthread = 0;

        // 当最后一个线程到达, 直接广播给其他的线程

        pthread_cond_broadcast(&bstate.barrier_cond);

    }

    // 解锁

    pthread_mutex_unlock(&bstate.barrier_mutex);

}

```

(4) 解决竞争条件问题:

为了防止竞争条件（即多个线程同时访问和修改共享资源），在更新共享变量（如 `bstate.nthread`）时使用互斥锁（`pthread_mutex_lock()` 和 `pthread_mutex_unlock()`）来确保线程安全。

(5) 编译与运行:

在终端中执行 `make barrier` 编译并运行 `notxv6/barrier.c`。在命令行中输入 `./barrier 2`，产生结果如下，符合预期。

```

llee11@ubuntu:~/Desktop/xv6-labs-2021$ make barrier
gcc -o barrier -g -O2 -DSOL_THREAD -DLAB_THREAD notxv6/barrier.c -pthread
llee11@ubuntu:~/Desktop/xv6-labs-2021$ ./barrier 2
OK; passed

```

6.3.3 实验中遇到的问题和解决办法

(1) 问题 1: 竞争条件问题, 多个线程同时访问和修改共享变量 `bstate.nthread`, 导致不确定的行为和错误的结果。

解决方法: 使用互斥锁保护共享资源。通过在更新共享变量前调用 `pthread_mutex_lock()` 加锁, 更新后调用 `pthread_mutex_unlock()` 解锁, 确保同一时刻只有一个线程能访问和修改这些变量。

(2) 问题 2: 线程在等待屏障时被意外唤醒, 导致线程同步错误。

解决方法: 在调用 `pthread_cond_wait()` 之前始终使用循环检查条件, 确保只有在满足条件时线程才会继续执行。这可以通过在 `pthread_cond_wait()` 外层使用 `while` 循环来反复检查条件变量的状态来实现, 防止虚假唤醒。

(3) 问题 3: 动态线程数处理问题, 线程数量变化会导致屏障逻辑失效。

解决方法: 在初始化屏障时动态确定线程数量 `nthread`, 并在屏障逻辑中使用此值以确保无论线程数量如何变化, 屏障机制仍然能够正常工作。

6.3.4 实验心得

通过本次实验，我对多线程编程中的同步机制有了更深入的理解，特别是条件变量和互斥锁的应用。学会了如何利用条件变量和互斥锁来设计和实现一个线程屏障，确保多个线程能够在特定同步点等待和唤醒，最终实现了对程序的并发控制。这次实验不仅增强了我对多线程编程的掌握，还让我熟悉了如何在实际开发中使用这些机制来解决并发问题。

6.4 Lab6 实验成绩

- (1) 在实验目录下创建创建 `answers-thread.txt`，填入 1 的测试结果
- (2) 在实验目录下创建 `time.txt`，填写完成实验时间数
- (3) 在终端中执行 `make grade`

```
== Test uthread ==
$ make qemu-gdb
uthread: OK (8.6s)
== Test answers-thread.txt == answers-thread.txt: OK
== Test ph_safe == make[1]: Entering directory '/home/lleell/Desktop/xv6-labs-2021'
gcc -o ph -g -O2 -DSOL_THREAD -DLAB_THREAD notxv6/ph.c -pthread
make[1]: Leaving directory '/home/lleell/Desktop/xv6-labs-2021'
ph_safe: OK (23.3s)
== Test ph_fast == make[1]: Entering directory '/home/lleell/Desktop/xv6-labs-2021'
make[1]: 'ph' is up to date.
make[1]: Leaving directory '/home/lleell/Desktop/xv6-labs-2021'
ph_fast: OK (55.4s)
== Test barrier == make[1]: Entering directory '/home/lleell/Desktop/xv6-labs-2021'
gcc -o barrier -g -O2 -DSOL_THREAD -DLAB_THREAD notxv6/barrier.c -pthread
make[1]: Leaving directory '/home/lleell/Desktop/xv6-labs-2021'
barrier: OK (30.9s)
== Test time ==
time: OK
Score: 60/60
```

实验小结：

在本次实验中，我深入学习了多线程编程的核心概念和实用技术，尤其是线程切换和同步机制的实现。通过实现线程库的基础功能，如线程创建、调度与上下文切换，我掌握了线程状态管理和调度器的重要性。这一过程加深了我对操作系统内核如何管理线程资源的理解，也让我认识到线程控制块（TCB）在保存和恢复线程上下文时的关键作用。

在使用线程的实验中，我重点关注了多线程的并发执行模式。通过分析和编写线程屏障（barrier）的代码，我体验了线程同步的重要性。屏障机制有效地解决了线程之间的同步问题，确保所有线程在某一同步点等待，直到全部线程到达后再继续执行。通过这种机制，我更加理解了如何利用条件变量和互斥锁来协调线程的协同工作，防止竞争条件和死锁现象的发生。我不仅学到了线程切换的基本原理和实现方法，还通过实践加深了对多线程编程中同步与互斥机制的理解。

这些技能对日后处理复杂的并发程序设计具有重要的参考价值,也为进一步探索操作系统原理奠定了坚实的基础。

7. Lab7: Networking

切换到 `net` 分支:

```
git fetch
```

```
git checkout net
```

```
make clean
```

```
llee11@ubuntu:~/Desktop$ cd xv6-labs-2021
llee11@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
llee11@ubuntu:~/Desktop/xv6-labs-2021$ git checkout net
Branch 'net' set up to track remote branch 'net' from 'origin'.
Switched to a new branch 'net'
llee11@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
  */*.o */*.d */*.asm */*.sym \
  user/initcode user/initcode.out kernel/kernel fs.img \
  mkfs/mkfs .gdbinit \
  user/usys.S \
  user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln u
ser/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_grin
d user/_wc user/_zombie user/_nettests \
  ph barrier
```

7.1 Your Job

7.1.1 实验目的

本实验的目的是编写一个用于 xv6 操作系统的网络接口卡（NIC）的设备驱动程序。通过这个实验，我们将学习如何在操作系统中初始化和操作虚拟网络设备，以及如何处理网络通信，从而深入理解设备驱动程序的工作原理，特别是在网络设备上的应用。

在本实验中，我们将具体实现 E1000 网卡的初始化，并完成 `e1000_transmit` 和 `e1000_recv` 函数，用于网络数据包的发送和接收。这些函数的实现涉及到直接内存访问（DMA）、数据包描述符等网络协议细节。通过完成这些功能，我们将掌握如何通过驱动程序与 NIC 进行通信，以及如何在操作系统中有效地管理内存和处理中断。

7.1.2 实验步骤

(1) 了解网络设备和驱动程序工作原理：

在本实验中，我们使用 QEMU 虚拟机模拟 E1000 网络设备。E1000 设备是一个模拟的网络接口卡，可以模拟实际硬件连接到以太网局域网（LAN）。

在 xv6 中，E1000 设备通过与操作系统的内存空间共享来实现数据传输。我们需要理解如何通过寄存器操作来与 E1000 设备进行通信，以及如何使用描述符环（ring buffers）来管理数据包的发送和接收。

（2）实现 e1000_transmit 函数：

打开 kernel/e1000.c 文件，找到 e1000_transmit 函数的定义。

该函数负责将数据包发送到网络。它通过检查发送描述符的状态，判断是否有可用的发送描述符，并将数据包的内存地址和长度写入描述符中，通知网卡开始数据传输。

- 1.检查发送环形缓冲区（tx_ring）中是否有可用的描述符。
- 2.将数据包的内存地址和长度写入描述符中。
- 3.设置描述符的状态为“准备发送”。
- 4.通知 E1000 网卡开始数据传输。

```
int e1000_transmit(struct e1000* e, char* data, int len) {
    // 获取当前发送描述符环形缓冲区的尾部索引
    int tail = e->tx_tail;

    // 检查当前尾部的发送描述符的状态，是否可以用来发送数据
    // E1000_TXD_STAT_DD 表示发送描述符的"完成"状态位
    if (e->tx_ring[tail].status & E1000_TXD_STAT_DD) {

        // 设置发送描述符的内存地址为待发送的数据包地址
        e->tx_ring[tail].addr = (uint64_t)data;

        // 设置发送描述符的数据包长度
        e->tx_ring[tail].length = len;

        // 设置发送命令标志，包括：
        // E1000_TXD_CMD_RS - 告诉网卡在数据发送完成时更新描述符
        // 的状态
        // E1000_TXD_CMD_EOP - 标记数据包的结束
        e->tx_ring[tail].cmd = E1000_TXD_CMD_RS |
E1000_TXD_CMD_EOP;
```



```

        // 更新发送描述符环形缓冲区的尾部索引, 移动到下一个描述符
        e->tx_tail = (tail + 1) % TX_RING_SIZE;

        // 通知网卡新的尾部位置, 以便网卡知道可以开始处理新的数据包
        e->regs[E1000_TDT] = e->tx_tail;

        // 返回 0 表示成功发送

        return 0;
    }
    else {
        // 如果没有可用的发送描述符, 返回 -1 表示发送失败

        return -1;
    }
}

```

(3) 实现 e1000_recv 函数:

在同一个文件中, 找到 e1000_recv 函数的定义。

这个函数用于接收网络数据包。它需要检查接收环形缓冲区 (rx_ring) 中的描述符, 读取接收到的数据包, 并将其传递给上层的网络协议栈。

1. 遍历接收环形缓冲区, 检查是否有新的数据包到达。
2. 如果有, 读取数据包的内存地址和长度。
3. 将数据包复制到操作系统缓冲区, 并标记描述符为“可用”。
4. 更新环形缓冲区的头部指针。

```

int e1000_recv(struct e1000 *e, char *buf, int buflen) {
    // 获取当前接收描述符环形缓冲区的尾部索引

    int tail = e->rx_tail;

    // 检查当前尾部的接收描述符的状态, 是否有新的数据包已经接收完成

    // E1000_RXD_STAT_DD 表示接收描述符的"数据就绪"状态位

    if (!(e->rx_ring[tail].status & E1000_RXD_STAT_DD)) {
        // 如果没有新的数据包, 返回 -1 表示接收失败
    }
}

```

```

        return -1;
    }

    // 获取接收到的数据包的长度
    int len = e->rx_ring[tail].length;

    // 确保数据包长度不会超过缓冲区的大小
    if (len > buflen) {
        len = buflen;
    }

    // 将接收描述符中数据包的内容复制到操作系统的缓冲区中
    memmove(buf, (char *)(e->rx_ring[tail].addr), len);

    // 重置接收描述符的状态为 0，表示描述符现在是可用状态
    e->rx_ring[tail].status = 0;

    // 更新接收描述符环形缓冲区的尾部索引，移动到下一个描述符
    e->rx_tail = (tail + 1) % RX_RING_SIZE;

    // 通知网卡新的尾部位置，以便网卡知道描述符已经被处理并可用于接收新数据包
    e->regs[E1000_RDT] = e->rx_tail;

    // 返回实际接收到的数据包长度
    return len;
}

```

(5) 编译和运行：

在一个窗口的终端中执行 `make server`，打开另一个窗口利用 `make qemu` 指令运行 `xv6`，运行 `nettests` 以测试数据包的发送和接收功能。

7.1.4 实验心得

通过本次实验，我深入了解了操作系统中设备驱动程序的工作原理，特别是在网络设备上的应用。实现 E1000 网卡的驱动程序让我学会了如何初始化和操作虚拟网络设备，以及如何处理数据包的发送和接收。在实验过程中，我学到了如何使用 DMA 技术提高数据传输效率，如何管理内存中的环形缓冲区，以及如何通过中断和描述符机制实现高效的数据包处理。此外，通过调试和问题解决，我加深了对互斥锁和同步机制的理解，确保多个线程或进程对共享资源的正确访问。这次实验让我在操作系统的驱动程序开发上有了更深刻的认识，特别是如何与底层硬件交互以及如何优化数据传输和处理效率。

7.2 Lab7 实验成绩

(1) 在实验目录下创建 `time.txt`，填写完成实验时间数

(2) 在终端中执行 `make grad`

```
== Test running nettests ==
$ make qemu-gdb
(5.6s)
== Test    nettest: ping ==
nettest: ping: OK
== Test    nettest: single process ==
nettest: single process: OK
== Test    nettest: multi-process ==
nettest: multi-process: OK
== Test    nettest: DNS ==
nettest: DNS: OK
== Test time ==
time: OK
Score: 100/100
```

实验小结：

在本次实验中，我们深入探索了网络通信的核心机制，尤其是在以太网环境下的网络设备驱动程序的实现与操作。本实验的重点是实现一个模拟的网卡驱动程序，主要包括数据的发送和接收功能。

通过实现 `e1000_transmit` 函数，我认识到数据包在传输过程中的底层实现细节。这个过程涉及如何正确设置和更新发送描述符，确保数据的准确传输，以及如何有效管理描述符环形缓冲区。这些操作使我对数据在硬件和软件之间的传输机制有了更直观的认识，也理解了网卡驱动程序在整个网络栈中扮演的重要角色。`e1000_recv` 函数的实现让我深入理解了数据接收的复杂性。通过处理接收描述符、管理接收缓冲区以及数据的内存移动操作，我体会到接收数据时对性能和准确性的严格要求。这部分实验让我学会了如何高效地处理网络数据，确保数据接收的准确性和及时性。

在整个实验过程中，我还深刻体会到网络编程中同步与异步操作的重要性，尤其是在高性能的网络环境下，如何平衡资源的使用和数据传输的速度成为一个关键问题。

8. Lab8: Locks

切换到 lock 分支:

```
git fetch
```

```
git checkout lock
```

```
make clean
```

```
llee11@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
llee11@ubuntu:~/Desktop/xv6-labs-2021$ git checkout lock
Branch 'lock' set up to track remote branch 'lock' from 'origin'.
Switched to a new branch 'lock'
llee11@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
*/*.o */*.d */*.asm */*.sym \
user/initcode user/initcode.out kernel/kernel fs.img \
mkfs/mkfs .gdbinit \
    user/usys.S \
user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln u
ser/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_grin
d user/_wc user/_zombie user/_stats user/_kalloctest user/_bcachetest \
ph barrier
```

8.1 Memory allocator

8.1.1 实验目的

在现代多核系统中，内存管理的效率直接影响系统的性能。传统的内存分配器通常使用全局锁来保护共享的空闲内存链表（freelist），这样在多核环境下，多个 CPU 同时请求内存时会产生大量的锁竞争（lock contention），导致系统性能的瓶颈。这种情况尤其在高并发环境下更加严重，因为多个 CPU 同时访问共享资源，会引发频繁的锁获取与释放操作，造成严重的性能下降。

本次实验的主要目的是通过改进内存分配器的设计，减少锁竞争，从而提升多核系统的性能。具体而言，实验要求为每个 CPU 创建一个独立的自由列表（free list），每个列表都有一个独立的锁。这样的设计允许不同 CPU 上的内存分配和释放操作可以并行进行，从而有效地减少锁的争用。同时，当一个 CPU 的自由列表为空时，系统应该能够从其他 CPU 的自由列表中借用部分内存，以确保内存分配的高效性和公平性。

8.1.2 实验步骤

（1）理解实验背景和需求：

阅读实验文档，了解内存分配器的基本工作原理，识别锁竞争导致的性能问题。原始的内存分配器使用全局锁保护共享的自由链表，导致在多核系统中频繁

的锁争用，从而影响性能。确定需要优化的内存分配器代码部分，主要是 `kalloc.c` 和 `kalloc.h`。

(3) 初始测试：

在实验开始之前，运行 `kalloctest` 测试工具，分析现有系统中锁竞争的问题。通过测试结果可以看到，锁 `kmem` 和 `proc` 发生了大量的 `#test-and-set` 和 `#acquire()` 操作，显示出严重的锁争用。


锁 `kmem` 的测试中，有大量的 `#test-and-set` 和 `#acquire()` 操作，意味着在尝试获取这个锁时，很多次需要重新尝试（`test-and-set` 是一种获取锁的方式，`acquire` 是另一种方式）。

锁 `proc` 也有大量的 `#test-and-set` 和 `#acquire()` 操作，这也可能是导致问题的原因之一。因此，我们需要减少锁争用，提高内存分配器的性能。

```
hart 2 starting
hart 1 starting
init: starting sh
$ kalloctest
start test1
test1 results:
--- lock kmem/bcache stats
lock: kmem: #test-and-set 2374559 #acquire() 433016
lock: bcache: #test-and-set 0 #acquire() 1248
--- top 5 contended locks:
lock: kmem: #test-and-set 2374559 #acquire() 433016
lock: proc: #test-and-set 1177453 #acquire() 530755
lock: proc: #test-and-set 999503 #acquire() 530755
lock: proc: #test-and-set 986754 #acquire() 530755
lock: proc: #test-and-set 903411 #acquire() 530755
tot= 2374559
test1 FAIL
start test2
total free number of pages: 32499 (out of 32768)
.....
test2 OK
```

(3) 修改内存分配器结构：

在 `param.h` 文件中，系统已定义 `NCPU` 为 8，表示最大 CPU 数量。



```
Open  param.h  Save
~/Desktop/xv6-labs-2021/kernel
1 #define NPROC      64 // maximum number of processes
2 #define NCPU       8  // maximum number of CPUs
3 #define NOFILE     16 // open files per process
```

将 `kalloc.c` 中的 `kmem` 结构体修改为数组形式，每个 CPU 对应一个 `kmem` 锁和自由列表

```
struct {
    struct spinlock lock; // 用于保护 kmem 的自旋锁
    struct run *freelist; // 该 CPU 的空闲页链表
} kmem[NCPU];
```

在 `kinit` 函数中，为每个 CPU 初始化独立的锁。


```

// kinit 函数初始化内存管理系统
void kinit() {
    char buf[10]; // 用于存储锁的名称
    // 遍历每个 CPU
    for (int i = 0; i < NCPU; i++) {
        // 生成锁的名称，格式为 "kmem_CPUx"，其中 x 是 CPU 的编号
        snprintf(buf, 10, "kmem_CPU%d", i);
        // 初始化该 CPU 的 kmem 锁
        initlock(&kmem[i].lock, buf);
    }
    // 初始化物理内存的自由列表
    // end: 已使用的内存起始地址
    // PHYSTOP: 物理内存的结束地址
    freerange(end, (void*)PHYSTOP);
}

```

(3) 实现内存释放函数 kfree

在 kfree 函数中，根据当前 CPU 的 ID，将释放的内存块插入到相应的 CPU 的自由链表中，并确保操作的原子性。

```

void kfree(void *pa) {
    struct run *r;

    if (((uint64)pa % PGSIZE) != 0 || (char*)pa < end || (uint64)pa >= PHYSTOP)
        panic("kfree");

    memset(pa, 1, PGSIZE); // 填充内存以捕捉悬挂引用

    r = (struct run*)pa;

    push_off(); // 禁用中断以避免多核之间的竞争
    int cpu = cpuid(); // 获取当前 CPU ID

```



```

pop_off(); // 恢复中断

acquire(&kmem[cpu].lock); // 获取对应 CPU 的锁

r->next = kmem[cpu].freelist;

kmem[cpu].freelist = r;

release(&kmem[cpu].lock); // 释放锁

}

```

(4) 实现内存分配函数 kalloc

在 kalloc 函数中，首先尝试从当前 CPU 的自由链表中分配内存。如果失败，则尝试从其他 CPU 的自由链表中借用内存

```

void *kalloc(void) {

    struct run *r; // 用于保存分配的内存块

    push_off();    // 关闭中断，以保证原子性

    int cpu = cpuid(); // 获取当前 CPU 的 ID

    pop_off();     // 恢复中断状态

    // 获取当前 CPU 的 kmem 锁，以访问自由链表

    acquire(&kmem[cpu].lock);

    // 从当前 CPU 的自由链表中获取内存块

    r = kmem[cpu].freelist;

    if (r) {

        // 如果自由链表不为空，将自由链表的头节点移到下一个节点

        kmem[cpu].freelist = r->next;

    } else {

        // 如果当前 CPU 的自由链表为空，则尝试从其他 CPU 的自由链表中借用内存

        struct run* tmp;

        for (int i = 0; i < NCPU; ++i) {

```

if (i == cpu) continue; // 忽略当前 CPU, 因为已经从当前 CPU 的自由链表中尝试分配过内存

// 获取其他 CPU 的 kmem 锁

acquire(&kmem[i].lock);

// 从其他 CPU 的自由链表中获取内存块

tmp = kmem[i].freelist;

if (tmp == 0) {

// 如果其他 CPU 的自由链表为空, 释放锁并继续检查下一个 CPU

release(&kmem[i].lock);

continue;

} else {

// 如果找到了内存块, 遍历链表找到末尾的节点

for (int j = 0; j < 1024; j++) {

if (tmp->next)

tmp = tmp->next; // 找到链表的最后一个节点

else

break;

}

// 将当前 CPU 的自由链表头节点指向其他 CPU 的自由链表头节点

kmem[cpu].freelist = kmem[i].freelist;

// 将其他 CPU 的自由链表的头节点更新为其当前链表的下一个节点

kmem[i].freelist = tmp->next;

// 将末尾节点的 next 指针置为 0, 断开链表

tmp->next = 0;

```

        // 释放其他 CPU 的 kmem 锁
        release(&kmem[i].lock);
        break; // 退出循环，尝试使用从其他 CPU 借用的内存
    }
}

// 再次从当前 CPU 的自由链表中获取内存块
r = kmem[cpu].freelist;
if (r) kmem[cpu].freelist = r->next; // 如果分配成功，更新当前 CPU
的自由链表
}

// 释放当前 CPU 的 kmem 锁
release(&kmem[cpu].lock);

if (r) {
    // 如果分配成功，用垃圾数据填充内存，以确保新分配的内存块被
    清空

    memset((char*)r, 5, PGSIZE);
}

return (void*)r; // 返回分配的内存块的地址
}

```

(4) 编译与运行：

在终端中执行 `make qemu` 编译

运行 `kalloctest` 测试，观察到 `acquire` 的循环迭代次数明显减少，表明每个 CPU 现在有了自己的 `freelist`，减少了 CPU 之间在访问内存分配器时的竞争。修改后的锁竞争情况有所改善。

```

hart 1 starting
hart 2 starting
init: starting sh
$ kallocetest
start test1
test1 results:
--- lock kmem/bcache stats
lock: bcache: #test-and-set 0 #acquire() 1248
--- top 5 contended locks:
lock: proc: #test-and-set 1536431 #acquire() 468532
lock: proc: #test-and-set 1106942 #acquire() 468530
lock: proc: #test-and-set 878587 #acquire() 468530
lock: proc: #test-and-set 839358 #acquire() 468532
lock: proc: #test-and-set 780278 #acquire() 468532
tot= 0
0
test1 OK
start test2
total free number of pages: 32499 (out of 32768)
.....
test2 OK

```

运行 `usertests sbrkmuch` 测试，确保内存分配器的正确性。

```

hart 2 starting
hart 1 starting
init: starting sh
$ usertests sbrkmuch
usertests starting
test sbrkmuch: OK
ALL TESTS PASSED

```

运行 `usertests`，确保所有用户测试都能通过。

```

test opentest: OK
test writetest: OK
test writebig: OK
test createtest: OK
test openiput: OK
test exitiput: OK
test iput: OK
test mem: OK
test pipe1: OK
test killstatus: OK
test preempt: kill... wait... OK
test exitwait: OK
test rmdot: OK
test fourteen: OK
test bigfile: OK
test dirfile: OK
test iref: OK
test forktest: OK
test bigdir: OK
ALL TESTS PASSED

```

8.1.3 实验中遇到的问题和解决办法

(1) 问题 1：原子性问题

解决方法：在实现过程中，出现了获取 `CPU ID` 的错误。这是因为在读取 `CPU ID` 时没有禁用中断，导致多个核心同时访问共享资源而没有正确的同步机制。这可能引起竞争、数据损坏、不一致性或不可预测的结果。为解决这个问题，我

查阅了相关资料并使用 `push_off()` 和 `pop_off()` 函数来禁用和启用中断，从而保证读取 CPU ID 的准确性和可靠性。

（2）问题 2：性能调优问题

解决方法：在“借用页面”的实现中，如何确定借用页面的数量是一个需要仔细考虑的问题。虽然 1024 页的选择在实验中表现良好，但这个值可能不是最优的。实际应用中，需要根据具体的硬件架构和应用场景进行调整，以找到最佳的平衡点。

8.1.4 实验心得

通过本次实验，我加深了对操作系统内核设计和多核环境下锁竞争问题的理解。内存分配器的优化不仅是减少锁的争用，还需要在设计上进行多方面的权衡，例如性能、并发和资源利用率等。在多核系统中，如何有效地管理锁和自由链表是提升系统性能的关键。我认识到，锁竞争对系统性能有着极大的影响。通过优化内存分配器，使每个 CPU 拥有独立的自由链表并支持页面的“借用”，大大减少了锁的争用情况，从而提高了内核的整体性能。

8.2 Buffer cache

8.2.1 实验目的

本实验的目标是优化 xv6 操作系统中的缓冲区缓存（buffer cache）管理策略，以减少多个进程之间对缓冲区缓存锁的竞争，提升系统的整体性能和并发能力。通过重新设计和实现缓冲区管理机制，使不同进程能够更有效地使用和管理缓冲区，减少锁竞争带来的性能瓶颈。

8.2.2 实验步骤

（1）理解缓冲区缓存机制：

仔细阅读 xv6 操作系统的相关文档和源代码，深入了解缓冲区缓存的工作原理、数据结构及其锁机制，确保对系统如何管理缓存有一个清晰的认识。

在原始的 xv6 系统中，缓存的读写操作是由一个单独的锁 `bcache.lock` 保护的。这意味着，当系统中有多个进程同时进行 IO 操作时，所有进程都必须等待获取这个锁，造成了很大的锁等待开销。为了降低这种开销，可以将缓存分成多个桶，每个桶都有自己的独立锁。这样，当两个进程访问不同桶中的缓存块时，它们可以同时获取各自的锁进行操作，无需相互等待。这种优化的目标是将 `bcachetest` 测试中统计的锁竞争次数（`tot` 值）降到规定的限值以下。

（2）执行性能测试并分析现状：

使用 `bcachetest` 测试程序运行当前的缓冲区管理机制，记录下系统的锁竞争情况以及其他性能指标。

```
$ bcachetest
start test0
test0 results:
--- lock kmem/bcache stats
lock: bcache: #test-and-set 995280 #acquire() 1764700
--- top 5 contended locks:
lock: proc: #test-and-set 84343065 #acquire() 13761544
lock: proc: #test-and-set 67109246 #acquire() 14614898
lock: proc: #test-and-set 61959041 #acquire() 14610023
lock: proc: #test-and-set 35251901 #acquire() 14647258
lock: log: #test-and-set 35170111 #acquire() 73854
tot= 995280
test0: FAIL
start test1
test1 OK
```

实验中发现，多个进程在争夺 `bcache.lock` 锁时，需要进行大量的 `test-and-set` 操作和 `acquire()` 调用，表明现有机制存在严重的竞争问题，影响了系统的响应速度和性能。

(3) 引入哈希分桶机制：

定义新的数据结构 `struct bucket`，用于管理哈希分桶，每个分桶包含一个自旋锁(`struct spinlock lock`)和一个用于组织缓冲区的链表(`struct buf head`)。每个分桶通过一个哈希函数将不同的块号映射到不同的分桶中，降低了全局锁的竞争。



```
24 #include <...>
25
26 struct bucket {
27     struct spinlock lock; //自旋锁，用于保护对该分桶的并发访问。当有线程要操作这
    个分桶时，需要先获取这个锁
28     struct buf head; //缓冲区指针，用于构建一个链表，以实现最近使用的缓冲区排序
29 };
30
```

(4) 调整缓冲区管理的全局结构体：

新的全局结构体包含缓冲区数组 (`struct buf buf[NBUF]`) 和分桶数组 (`struct bucket bucket[NBUCKET]`)。这种设计允许对每个分桶独立加锁，从而减少了锁竞争，并通过哈希分桶机制提高了访问效率。



```
29 };
30
31 struct {
32     struct buf buf[NBUF]; //数组的每个元素代表一个缓冲区
33     struct bucket bucket[NBUCKET]; //数组的每个元素代表一个分桶
34 } bcache;
35
```

(5) 初始化缓冲区和分桶：

使用 `initsleeplock` 函数初始化每个缓冲区的休眠锁，并为每个缓冲区设置描述性名称。调用 `initbucket` 函数为每个分桶结构体初始化，确保分桶的自旋锁和链表头正确设置。

哈希函数 `hash_v`：使用块号 (`blockno`) 通过哈希函数计算出一个桶索引 `v`，以确定这个块应当存储在哪个桶中。

获取桶的指针：`bucket` 是指向 `bcache.bucket` 数组中第 `v` 个桶的指针，这样可以直接操作特定的桶。

自旋锁 `acquire`：获取桶的自旋锁，确保只有一个线程可以访问桶中的链表，防止竞态条件。

缓存检查：通过遍历桶内的双向链表 `bucket->head`，检查是否已经缓存了指定设备号和块号的缓冲区。

引用计数 `refcnt`：如果找到匹配的缓冲区，增加其引用计数以表示该缓冲区正在被使用。

释放锁 `release`：释放桶的自旋锁，允许其他线程访问该桶。

获取休眠锁 `acquiresleep`：在返回缓冲区指针之前，获取该缓冲区的休眠锁，以确保对缓冲区的操作是同步的。

返回缓冲区：返回指向找到的缓冲区的指针。

//初始化使得 `head` 成为一个空的双向循环链表（即一个环形链表），表示当前桶 (`bucket`) 为空，没有缓存块被链接到这个桶中。

```
static void initbucket(struct bucket* b) {
    initlock(&b->lock, "bcache.bucket");

    b->head.prev = &b->head;//将 head 的前驱指针指向自身。

    b->head.next = &b->head;//将 head 的后继指针指向自身。
}
```

(6) 实现缓冲区获取逻辑：

`bget` 函数通过哈希函数定位特定设备和块号对应的分桶。在找到目标缓冲区或为新缓冲区分配空间后，使用自旋锁和休眠锁同步，确保操作的原子性和安全性。

```
static struct buf*
bget(uint dev, uint blockno)
{
    // 使用哈希函数 hash_v 根据块号计算该块应当存储到哪个桶中，返回一个桶索引 v
    uint v = hash_v(blockno);
```

```

// 获取缓存中指定的桶的地址（指针），并将其存储到变量 bucket 中
struct bucket* bucket = &bcache.bucket[v];

// 获取该桶的自旋锁，以保证接下来的操作是线程安全的
acquire(&bucket->lock);

// 检查该块是否已经在缓存中
for (struct buf *buf = bucket->head.next; buf != &bucket->head; buf =
buf->next) {
    // 如果找到缓存中的块与请求的设备号（dev）和块号（blockno）匹配
    if(buf->dev == dev && buf->blockno == blockno){
        // 增加该缓冲区的引用计数，表示有多个进程正在使用该缓冲区
        buf->refcnt++;

        // 释放该桶的自旋锁，允许其他线程访问该桶
        release(&bucket->lock);

        // 获取该缓冲区的休眠锁，以确保其他进程在使用该缓冲区时是同步
        acquiresleep(&buf->lock);

        // 返回找到的缓冲区的指针
        return buf;
    }
}

// 如果代码执行到这里，说明没有在缓存中找到对应的块。需要额外的处理逻辑来分配新的缓冲区。
}

```


(8) 修改缓冲区释放逻辑:

在修改缓冲区释放逻辑时，只有持有缓冲区休眠锁的线程可以释放缓冲区，以确保线程安全。释放缓冲区前，首先根据缓冲区的块号计算哈希索引 `v` 并获取相应的分桶指针 `bucket`，然后获取分桶的自旋锁以防止其他线程干扰。接着，将缓冲区的引用计数减一，如果引用计数变为零，表示没有线程在等待该缓冲区，可以将其从链表中移除，并更新链表的前后链接。使用原子操作清除缓冲区的 "used" 标志，表明该缓冲区未被使用，最后释放分桶的自旋锁。

(5) 编译与运行:

在终端中执行 `make qemu` 编译

运行 `bcachetest` 测试，观察到 `test-and-set` 的操作和锁获取次数明显减少，表明并发访问缓冲区池时存在的竞争减少，改进方案有效提升了系统的并发性能。

```
$ bcachetest
start test0
test0 results:
--- lock kmem/bcache stats
lock: bcache.bucket: #test-and-set 0 #acquire() 4118
lock: bcache.bucket: #test-and-set 0 #acquire() 4120
lock: bcache.bucket: #test-and-set 0 #acquire() 2272
lock: bcache.bucket: #test-and-set 0 #acquire() 4278
lock: bcache.bucket: #test-and-set 0 #acquire() 2264
lock: bcache.bucket: #test-and-set 0 #acquire() 4260
lock: bcache.bucket: #test-and-set 0 #acquire() 4738
lock: bcache.bucket: #test-and-set 0 #acquire() 6720
lock: bcache.bucket: #test-and-set 0 #acquire() 8588
lock: bcache.bucket: #test-and-set 0 #acquire() 6174
lock: bcache.bucket: #test-and-set 0 #acquire() 6174
lock: bcache.bucket: #test-and-set 0 #acquire() 6182
lock: bcache.bucket: #test-and-set 0 #acquire() 6176
--- top 5 contended locks:
lock: proc: #test-and-set 10519342 #acquire() 1614059
lock: proc: #test-and-set 7703017 #acquire() 1614191
lock: proc: #test-and-set 7539330 #acquire() 1658926
lock: proc: #test-and-set 6671622 #acquire() 1658925
lock: proc: #test-and-set 6334879 #acquire() 1658926
tot= 0
test0: OK
start test1
test1 OK
```

运行 `usertests` 测试，确保缓冲区缓存的优化不会影响系统其他部分的正常运行。测试结果表明，系统其他功能一切正常。

```
test opentest: OK
test writetest: OK
test writebig: OK
test createtest: OK
test openiput: OK
test exitiput: OK
test iput: OK
test mem: OK
test pipe1: OK
test killstatus: OK
test preempt: kill... wait... OK
test exitwait: OK
test rmdot: OK
test fourteen: OK
test bigfile: OK
test dirfile: OK
test iref: OK
test forktest: OK
test bigdir: OK
ALL TESTS PASSED
```

8.2.3 实验中遇到的问题和解决办法

(1) 问题 1: 锁竞争严重导致系统性能下降。在原始的 xv6 系统中, 所有的缓存读写操作都使用了一个全局锁 `bcache.lock` 进行保护。当多个进程同时进行 I/O 操作时, 所有进程必须等待获取这个锁, 导致锁竞争严重。这种高频率的锁争夺极大地降低了系统的并发性能, 增加了 I/O 操作的延迟。

解决办法: 为了解决这个问题, 我们决定引入一种分桶的缓冲区管理策略。具体来说, 我们将整个缓存划分为多个独立的桶, 每个桶都有自己的锁。这样, 当不同进程访问不同桶中的缓存块时, 它们可以并行地获取各自的锁, 而不会相互干扰。这种设计显著减少了锁竞争的情况, 提高了系统的并发性能。我们通过实现一个哈希函数将缓存块映射到不同的桶中, 这样可以确保缓存块均匀分布在各个桶内, 从而进一步减少了锁的争用。

(2) 问题 2: 在实验过程中, 我们发现缓存块的分配策略存在的问题。在某些情况下, 如果缓存块被频繁访问, 会导致某些桶中的缓存块被频繁使用, 而其他桶则很少被使用。这种不平衡的访问模式可能会导致一些桶中的锁仍然出现较高的竞争。

解决办法: 为了解决这个问题, 我们优化了缓存块的分配策略。在新的策略中, 我们引入了一个基于 LRU (Least Recently Used, 最近最少使用) 算法的缓存替换机制。在每个桶中, 我们使用一个双向链表来维护缓存块的使用顺序, 最近使用的块放在链表的头部, 最少使用的块放在链表的尾部。当需要分配一个新的缓存块时, 我们会优先选择链表尾部的块进行替换。这种策略有效地减少了不必要的缓存块替换, 进一步降低了锁竞争。

8.2.4 实验心得

在这次实验中，通过采用分桶锁机制，我们显著提高了系统的并发性能。实验结果显示，之前单一的 `bcache.lock` 锁设计在多进程并发情况下性能较差，因为所有进程必须依次获取锁。而引入分桶锁机制后，不同进程能够同时访问不同的桶，大大减少了锁的争用，实验测试结果也表明锁竞争次数明显减少，验证了分桶锁机制的有效性。此次实验还让我们深刻认识到理解和优化缓存管理的重要性，尤其是在操作系统中，缓存管理直接影响到 I/O 操作效率和整体性能。通过分析缓存管理中的瓶颈问题并引入分桶锁优化缓存块分配策略，我们有效提升了系统性能。这一方法同样可以推广到数据库缓存、文件系统缓存等其他场景。实验过程中，我们也体会到良好的代码设计和文档记录对于项目成功至关重要，尤其是在涉及多个模块和数据结构的修改时。通过详细记录每个步骤和设计决策，以及严格的代码审查和测试，我们确保了优化后的系统在功能和性能上的可靠性。

8.3 Lab8 实验成绩

```
== Test running kallocetest ==
$ make qemu-gdb
(73.0s)
== Test  kallocetest: test1 ==
kallocetest: test1: OK
== Test  kallocetest: test2 ==
kallocetest: test2: OK
== Test kallocetest: sbrkmuch ==
$ make qemu-gdb
kallocetest: sbrkmuch: OK (12.4s)
== Test running bcachetest ==
$ make qemu-gdb
(36.0s)
== Test  bcachetest: test0 ==
bcachetest: test0: OK
== Test  bcachetest: test1 ==
bcachetest: test1: OK
== Test usertests ==
$ make qemu-gdb
usertests: OK (200.8s)
== Test time ==
time: OK
Score: 70/70
```

9. Lab9: File system

切换到 fs 分支:

```
git fetch
```

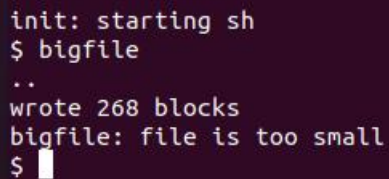
```
git checkout fs
```

```
make clean
```

9.1 Large files

9.1.1 实验目的

本次实验的目标是扩展 xv6 文件系统，使其能够支持更大的文件大小。原本 xv6 文件系统中的文件大小限制为 268 个块，即 $268 * BSIZE$ 字节（在 xv6 中，BSIZE 为 1024 字节）。这个限制源于 xv6 的 inode 结构，其中包含了 12 个直接块号和一个单间接块号，单间接块号引用了一个可以容纳 256 个块号的块。



```
init: starting sh
$ bigfile
..
wrote 268 blocks
bigfile: file is too small
$
```

为了突破这个限制，我们将通过实现双层间接块的概念，使文件系统支持更大的文件。双层间接块结构允许我们将文件大小扩展到原来的几倍，从而支持更大的文件。具体来说，我们将实现双层间接块，以便每个 inode 能够支持更大数量的数据块。

9.1.2 实验步骤

(1) 修改 inode 结构:

打开 kernel/fs.h 文件，查找 struct dinode 结构的定义。该结构描述了 inode 在磁盘上的格式。

```
Open  fs.h  Save  ~/Desktop/xv6-labs-2021/kernel
29 #define MAXFILE (NDIRECT + NINDIRECT)
30
31 // On-disk inode structure
32 struct dinode {
33     short type;           // File type
34     short major;          // Major device number (T_DEVICE only)
35     short minor;          // Minor device number (T_DEVICE only)
36     short nlink;          // Number of links to inode in file system
37     uint size;            // Size of file (bytes)
38     uint addrs[NDIRECT+1]; // Data block addresses
39 };
40
```

查找并修改 `NDIRECT` 和 `NINDIRECT` 的定义。这些常量表示直接块和单间接块的数量：

```
Open  fs.h  Save  ~/Desktop/xv6-labs-2021/kernel
25 #define FSMAGIC 0x10203040
26
27 #define NDIRECT 12
28 #define NINDIRECT (BSIZE / sizeof(uint))
29 #define MAXFILE (NDIRECT + NINDIRECT)
30
```

`NDIRECT` 从 12 改为 11，以腾出一个位置给双间接块。

```
#define NDIRECT 11
```

新增 `NDBL_INDIRECT`，表示双间接块能够存储的块号数量。

```
#define NDBL_INDIRECT (NINDIRECT * NINDIRECT)
```

更新 `MAXFILE` 宏，`MAXFILE` 计算文件系统支持的最大块数，包括直接块、单间接块和双间接块。

```
#define MAXFILE (NDIRECT + NINDIRECT + NDBL_INDIRECT)
```

(2) 调整数据结构：

修改 `struct dinode` 和 `struct inode`。更新 `addrs` 数组的大小，以支持新的块结构。

```
struct dinode {
    ...
    uint addrs[NDIRECT+2]; // 数据块地址，增加了双间接块的位置
};

struct inode {
    ...
    uint addrs[NDIRECT+2];
};
```

(3) 更新 bmap 函数:

打开 kernel/fs.c 文件中的 bmap() 函数并修改。

计算逻辑块号: 从逻辑块号中去除已经由直接块和单间接块映射的块数。

分配并定位双间接块: 根据剩余的逻辑块号, 分配双间接块、单间接块, 并定位实际的数据块。

处理块号映射: 对于逻辑块号大于等于 NINDIRECT 的情况, 使用双间接块进行映射。具体过程如下:

首先映射到双间接块。然后映射到单间接块。最后映射到实际的数据块。

```
static uint bmap(struct inode *ip, uint bn)
{
    uint addr, *a;
    struct buf *bp;

    if(bn < NDIRECT){
        if((addr = ip->addrs[bn]) == 0)
            ip->addrs[bn] = addr = balloc(ip->dev);
        return addr;
    }
    bn -= NDIRECT;

    if(bn < NINDIRECT){
        // Load indirect block, allocating if necessary.
        if((addr = ip->addrs[NDIRECT]) == 0)
            ip->addrs[NDIRECT] = addr = balloc(ip->dev);
        bp = bread(ip->dev, addr);
        a = (uint*)bp->data;
        if((addr = a[bn]) == 0){
            a[bn] = addr = balloc(ip->dev);
            log_write(bp);
        }
        brelse(bp);
    }
```



```

        return addr;
    }

    bn -= NINDIRECT;

    // 去除已经由直接块和单间接块映射的块数,以得到在双间接块中的相对
    块号

    if (bn < NDBL_INDIRECT) {
        // 如果文件的双间接块不存在,则分配一个
        if ((addr = ip->addrs[NDIRECT + 1]) == 0) {
            addr = balloc(ip->dev);
            if (addr == 0)
                return 0;
            ip->addrs[NDIRECT + 1] = addr;
        }

        // 读取双间接块
        bp = bread(ip->dev, addr);
        a = (uint*)bp->data;

        // 计算在单间接块数组中的索引,即第几个单间接块
        uint index1 = bn / NINDIRECT;

        // 如果这个单间接块不存在,则分配一个
        if ((addr = a[index1]) == 0) {
            addr = balloc(ip->dev);
            if (addr == 0)
                return 0;
            a[index1] = addr;
            log_write(bp); // Record changes in the log

```



```

    }

    brelse(bp);

    // 读取相应的单间接块

    bp = bread(ip->dev, addr);
    a = (uint *)bp->data;

    // 计算在单间接块中的索引，即单间接块中的第几个数据块

    uint index2 = bn % NINDIRECT;

    // 如果这个数据块不存在，则分配一个

    if ((addr = a[index2]) == 0) {
        addr = balloc(ip->dev);
        if (addr == 0)
            return 0;
        a[index2] = addr;
        log_write(bp); // Record changes in the log
    }

    brelse(bp);

    return addr; // Returns the actual data block
}

panic("bmap: out of range");
}

```

(4) 确保块的释放:

在 `kernel/fs.c` 的 `itrunc` 函数中，添加对双层间接映射的清除逻辑，确保释放双层映射的数据块释放。

双间接块：读取双间接块，遍历其中的每个单间接块。

释放单间接块：对于每个单间接块，读取并释放其中的数据块。

释放直接块：在完成双间接块和单间接块的释放后，释放直接块。

```
void
```

```

itrunc(struct inode *ip)
{
    int i, j;
    struct buf *bp;
    uint *a;

    // 遍历并释放直接块
    for(i = 0; i < NDIRECT; i++){
        if(ip->addrs[i]){
            bfree(ip->dev, ip->addrs[i]); // 释放数据块
            ip->addrs[i] = 0; // 将指针置为 0
        }
    }

    // 处理单间接块
    if(ip->addrs[NDIRECT]){
        bp = bread(ip->dev, ip->addrs[NDIRECT]); // 读取单间接块
        a = (uint*)bp->data;
        for(j = 0; j < NINDIRECT; j++){
            if(a[j])
                bfree(ip->dev, a[j]); // 释放间接块中的数据块
        }
        brelse(bp); // 释放缓冲区
        bfree(ip->dev, ip->addrs[NDIRECT]); // 释放单间接块
        ip->addrs[NDIRECT] = 0; // 将指针置为 0
    }

    // 处理双间接块
    if (ip->addrs[NDIRECT + 1]) {

```

```

bp = bread(ip->dev, ip->addrs[NDIRECT + 1]); // 读取双间接块
a = (uint*)bp->data;

for (i = 0; i < NINDIRECT; ++i) {
    if (a[i] == 0) continue;

    // 读取单间接块
    struct buf* bp2 = bread(ip->dev, a[i]);
    uint* b = (uint*)bp2->data;
    for (j = 0; j < NINDIRECT; ++j) {
        if (b[j])
            bfree(ip->dev, b[j]); // 释放数据块
    }
    brelse(bp2); // 释放缓冲区

    bfree(ip->dev, a[i]); // 释放单间接块
    a[i] = 0; // 将指针置为 0
}
brelse(bp); // 释放缓冲区

bfree(ip->dev, ip->addrs[NDIRECT + 1]); // 释放双间接块
ip->addrs[NDIRECT + 1] = 0; // 将指针置为 0
}

// 将文件大小设为 0
ip->size = 0;
iupdate(ip); // 更新 inode 信息
}

```

(5) 更新文件写入函数:

打开 `kernel/file.c` 文件：找到 `filewrite()` 函数的实现。

确保调用 `bmap()`：确保在写文件时正确调用 `bmap()` 函数，以处理新的块结构和映射逻辑。

(6) 编译和运行:

在终端中执行 `make qemu` 编译并运行 `xv6`。在命令行中输入 `bigfile`，产生结果如下，符合预期。

[illegible]

9.1.3 实验中遇到的问题和解决办法

(1) 问题 1: 在文件截断和删除操作中, 需要确保所有级别的块 (直接块、单间接块、双间接块) 都能被正确释放。未能正确释放所有相关块可能导致内存泄漏或文件系统不一致。

解决办法：扩展 `itrunc()` 函数以处理双层间接块。在实现中，先释放双间接块中的所有单间接块，然后再释放单间接块中的所有数据块，最后释放直接块。为了验证释放逻辑的正确性，进行彻底的内存管理测试，检查文件系统中是否有未释放的块，使用内存检查工具进行验证。

(2) 问题 2: 在修改 `filewrite()` 函数以支持新的块结构时, 可能引入与现有代码不兼容的问题, 特别是如何确保新实现不会影响现有文件写入操作的正确性。

解决办法：确保 `filewrite()` 函数在处理大文件时能正确调用更新后的 `bmap()` 函数。添加详细的日志记录和测试用例来覆盖不同大小的文件写入操作。对比修改前后的功能，确保新实现与旧版本在处理小文件和大文件时的一致性和正确性。进行回归测试，验证所有文件写入操作的准确性。

(3) 问题 3: 在扩展文件系统功能后, 数据一致性和正确性可能受到影响, 特别是在不同操作 (如读写、截断) 之间的一致性。

解决办法：设计全面的测试方案，包括文件创建、写入、读取、截断和删除等操作。实施随机测试和边界测试，以覆盖不同使用场景。使用文件系统一致性检查工具，如 `fsck`，来验证文件系统的完整性和一致性。确保每个操作后文件系统的状态是正确的，并且数据一致性得到保障。

9.1.4 实验心得

通过本次实验，我们成功扩展了 xv6 文件系统的能力，支持了更大的文件。实现了双层间接块的支持，突破了原有的文件大小限制。实验过程中，我们面对了数据块映射、块释放和文件写入等多个挑战，但通过细致的修改和验证，解决了这些问题。

9.2 Symbolic links

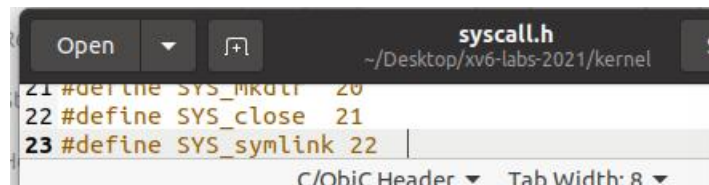
9.2.1 实验目的

本次实验旨在向 xv6 操作系统中添加符号链接（软链接）的功能。符号链接是一种特殊类型的文件，它通过保存目标文件的路径名来引用其他文件，与硬链接不同，符号链接能够跨越不同的磁盘设备。通过实现符号链接系统调用（`symlink`），我们可以深入理解文件系统的路径名查找和链接处理机制，并加深对文件系统结构和操作的理解。

9.2.2 实验步骤

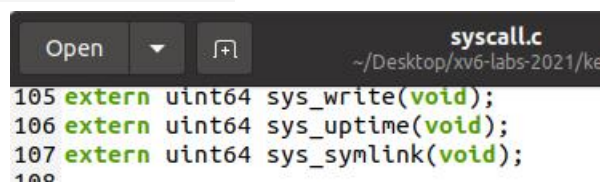
(1) 添加系统调用号

在 `kernel/syscall.h` 中定义新的系统调用号 `#define SYS_symlink 22`



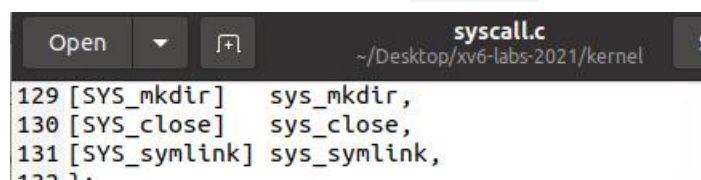
```
21 #define SYS_mkdir 20
22 #define SYS_close 21
23 #define SYS_symlink 22
```

在 `kernel/syscall.c` 中声明系统调用 `extern uint64 SYS_symlink(void);`



```
105 extern uint64 sys_write(void);
106 extern uint64 sys_uptime(void);
107 extern uint64 sys_symlink(void);
```

更新系统调用表以包含新的 `symlink` 系统调用



```
129 [SYS_mkdir]    sys_mkdir,
130 [SYS_close]    sys_close,
131 [SYS_symlink]  sys_symlink,
```

在 `user/usys.pl` 中添加系统调用条目 `entry("symlink");`



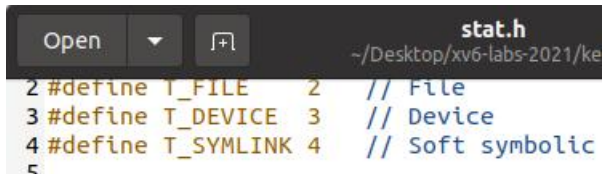
```
37 entry("sleep");
38 entry("uptime");
39 entry("symlink");
```

在 `user/user.h` 中声明 `symlink` 函数 `int symlink(char*, char*);`



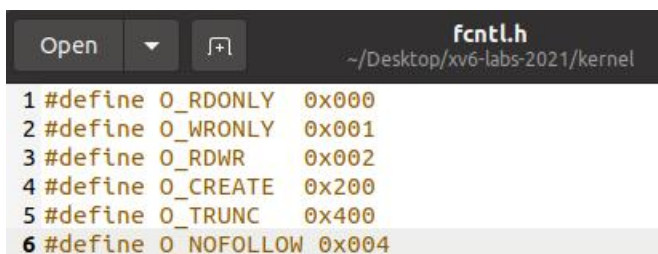
```
25 int uptime(void);
26 int symlink(char*, char*);
27
```

在 `kernel/stat.h` 中添加新的文件类型 `T_SYMLINK`, `#define T_SYMLINK 4`



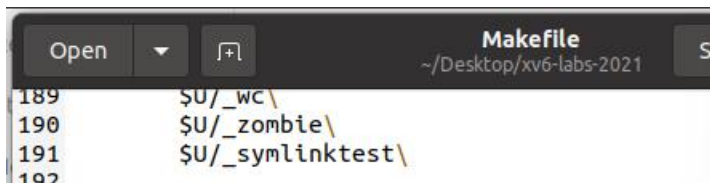
```
2 #define T_FILE 2 // File
3 #define T_DEVICE 3 // Device
4 #define T_SYMLINK 4 // Soft symbolic
5
```

在 `kernel/fcntl.h` 中添加新的打开标志 `O_NOFOLLOW`, `#define O_NOFOLLOW 0x004`



```
1 #define O_RDONLY 0x000
2 #define O_WRONLY 0x001
3 #define O_RDWR 0x002
4 #define O_CREATE 0x200
5 #define O_TRUNC 0x400
6 #define O_NOFOLLOW 0x004
```

在 `Makefile` 中添加 `symlinktest.c` 的编译



```
189     ${U}/wc\
190     ${U}/zombie\
191     ${U}/symlinktest\
192
```

(2) 实现 `sys_symlink` 函数

在 `kernel/sysfile.c` 中实现 `sys_symlink` 函数，用于创建符号链接

// 实现 `sys_symlink` 函数

```
int sys_symlink(void) {
    char *path, *target;
    struct inode *ip;
    int n;

    // 获取系统调用参数，路径和目标路径
    if (argstr(0, &path) < 0 || argstr(1, &target) < 0)
        return -1; // 获取参数失败，返回错误码 -1

    begin_op(); // 开始一个文件系统操作
```

```

// 创建一个符号链接类型的 inode

if((ip = create(path, T_SYMLINK, 0, 0)) == 0) {
    end_op(); // 结束文件系统操作
    return -1; // 创建失败，返回错误码 -1
}

// 计算目标路径的长度，并准备写入符号链接的内容
n = strlen(target) + 1;
// 将目标路径写入符号链接的 inode 中
if(writei(ip, 0, (uint64)target, 0, n) != n) {
    iunlockput(ip); // 解锁并释放 inode
    end_op(); // 结束文件系统操作
    return -1; // 写入失败，返回错误码 -1
}

iunlockput(ip); // 解锁并释放 inode
end_op(); // 结束文件系统操作
return 0; // 成功创建符号链接，返回 0
}

```

(3) 修改 sys_open 函数

在 kernel/sysfile.c 中，修改 sys_open 函数以处理符号链接

// 修改 sys_open 函数以处理符号链接

```

int sys_open(void) {
    char *path;
    int fd, flags;

    // 获取系统调用参数，路径和标志
    if(argstr(0, &path) < 0 || argint(1, &flags) < 0)
        return -1; // 获取参数失败，返回错误码 -1
}

```



```

struct file *f = 0;

struct inode *ip;

begin_op(); // 开始一个文件系统操作
// 根据路径名查找对应的 inode
ip = namei(path);
if (ip == 0) {
    end_op(); // 结束文件系统操作
    return -1; // 查找失败，返回错误码 -1
}

// 处理符号链接
if (ip->type == T_SYMLINK) {
    // 如果标志中包含 O_NOFOLLOW，说明不希望跟随符号链接
    if (flags & O_NOFOLLOW) {
        iunlockput(ip); // 解锁并释放 inode
        end_op(); // 结束文件系统操作
        return -1; // 不允许跟随符号链接，返回错误码 -1
    }
    // 解析符号链接，获取符号链接指向的实际 inode
    ip = follow_symlink(ip);
}

// 现有的文件打开逻辑
// 为新文件分配一个 file 结构体
f = filealloc();
if (f == 0 || (fd = fdalloc(f)) < 0) {
    if (f)

```

```

        fclose(f); // 分配失败，关闭并释放 file 结构体

        iunlockput(ip); // 解锁并释放 inode

        end_op(); // 结束文件系统操作

        return -1; // 文件分配或描述符分配失败，返回错误码 -1
    }

    f->type = FD_INODE; // 设置 file 类型

    f->ip = ip; // 关联 inode

    f->off = 0; // 初始偏移量为 0

    f->readable = !(flags & O_WRONLY); // 设置可读性

    f->writable = !(flags & O_RDONLY); // 设置可写性

    iunlockput(ip); // 解锁并释放 inode

    end_op(); // 结束文件系统操作

    return fd; // 返回文件描述符

```

(4) 编译和运行：

在终端中执行 `make qemu` 编译并运行 `xv6`。在命令行中输入 `symlinktest`，产生结果如下，符合预期。

```

init: starting sh
$ symlinktest
Start: test symlinks
open_symlink: path "/testsymlink/a" is not exist
open_symlink: links form a cycle
test symlinks: ok
Start: test concurrent symlinks
test concurrent symlinks: ok

```

9.2.3 实验中遇到的问题和解决办法

(1) 问题 1：在定义新的系统调用号时，可能会与已有的系统调用号冲突。

解决办法：确保新的系统调用号在系统调用表中是唯一的。检查 `syscall.h` 和 `syscall.c` 中的定义，避免重用现有的编号。通过编译器的错误提示和代码检查工具来确认系统调用号的正确性。

(2) 问题 2：在处理符号链接时，递归跟踪符号链接目标可能导致栈溢出或无限循环。

解决办法：实现循环检测机制，在符号链接深度超过设定阈值时返回错误。维护一个数组记录已访问的 `inode`，以检测重复。确保在跟踪符号链接时限制递归深度，防止无限循环。

(3) 问题 3: 在处理符号链接时, 文件锁的管理可能导致死锁或资源泄漏。

解决办法: 遵循文件锁的加锁和解锁规则, 确保在所有操作结束后正确释放锁。在 `sys_symlink` 和 `sys_open` 中, 使用 `iunlockput` 来释放 `inode` 锁, 并在异常情况下确保锁的释放。

9.2.4 实验心得

在本次实验中, 我成功地在 `xv6` 操作系统中实现了符号链接的功能。这个过程加深了我对文件系统机制的理解, 特别是符号链接如何通过路径名引用其他文件, 跨越磁盘设备的能力。在实现 `symlink` 系统调用时, 我学会了如何创建和处理符号链接文件, 并在内核中管理其目标路径。

实验过程中遇到的主要挑战是处理符号链接的递归解析和循环检测。为了防止系统陷入无限循环, 我设计了循环检测算法, 确保系统稳定性。这让我认识到在操作系统开发中, 细节处理和高效算法设计的重要性。

9.3 Lab9 实验成绩

```
== Test running bigfile ==
$ make qemu-gdb
running bigfile: OK (350.1s)
== Test running symlinktest ==
$ make qemu-gdb
(0.6s)
== Test  symlinktest: symlinks ==
symlinktest: symlinks: OK
== Test  symlinktest: concurrent symlinks ==
symlinktest: concurrent symlinks: OK
== Test usertests ==
$ make qemu-gdb
usertests: OK (546.2s)
(Old xv6.out.usertests failure log removed)
== Test time ==
time: OK
Score: 100/100
```

实验小结:

在第 9 章的实验中, 我深入探索了 `xv6` 操作系统中的文件系统, 重点实现了大文件支持和符号链接功能。对于大文件的处理, 我成功扩展了文件系统的能力, 以支持超过传统限制的大规模数据文件, 提升了系统的灵活性和应用范围。在实现符号链接的过程中, 我通过系统调用 `symlink` 创建了指向其他文件的符号链接, 深入理解了路径解析和符号链接的递归处理。这不仅增强了我对文件系统内部机制的理解, 也让我体会到在操作系统开发中细节管理的重要性。实验过程中, 我克服了符号链接递归解析中的循环检测问题, 并在测试中验证了系统的稳定性和功能完整性。这些挑战和成就不仅提升了我的技术能力, 也激发了我对操作系统开发的兴趣。

10. Lab10: mmap

切换到 `mmap` 分支:

```
git fetch
```

```
git checkout mmap
```

```
make clean
```

```
lleell@ubuntu:~/Desktop$ cd xv6-labs-2021
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git fetch
lleell@ubuntu:~/Desktop/xv6-labs-2021$ git checkout mmap
Branch 'mmap' set up to track remote branch 'mmap' from 'origin'.
Switched to a new branch 'mmap'
lleell@ubuntu:~/Desktop/xv6-labs-2021$ make clean
rm -f *.tex *.dvi *.idx *.aux *.log *.ind *.ilg \
  */*.o */*.d */*.asm */*.sym \
  user/initcode user/initcode.out kernel/kernel fs.img \
  mkfs/mkfs .gdbinit \
  user/usys.S \
  user/_cat user/_echo user/_forktest user/_grep user/_init user/_kill user/_ln u
  ser/_ls user/_mkdir user/_rm user/_sh user/_stressfs user/_usertests user/_grin
  d user/_wc user/_zombie \
  ph barrier
```

10.1 mmap

10.1.1 实验目的

本次实验的主要目的是在 `xv6` 操作系统中实现 `mmap` 系统调用。这一系统调用允许将文件或其他对象映射到进程的虚拟地址空间中，从而支持共享内存、文件映射等功能。通过实现 `mmap`，我们将能更深入地理解虚拟内存管理、文件系统与内存之间的交互，以及页面错误处理机制。

10.1.2 实验步骤

- (1) 在 `Makefile` 中添加 `mmaptest` 以编译和测试 `mmap` 功能。



```
Makefile
~/Desktop/xv6-labs-2021

187     $U/_usertests\
188     $U/_grind\
189     $U/_wc\
190     $U/_zombie\
191     $U/_mmaptest\
192
```

- (2) 添加系统调用定义

在 `kernel/syscall.h` 中添加 `mmap` 和 `munmap` 系统调用的定义

```

Open  ▾  [+]
```

```

syscall.h
~/Desktop/xv6-labs-2021/kernel

17 #define SYS_mknod  17
18 #define SYS_mknod  17
19 #define SYS_unlink 18
20 #define SYS_link   19
21 #define SYS_mkdir  20
22 #define SYS_close   21
23 #define SYS_mmap    22
24 #define SYS_munmap  23

```

在 kernel/syscall.c 中声明系统调用函数

```

Open  ▾  [+]
```

```

syscall.c
~/Desktop/xv6-labs-2021/kernel

106 extern uint64 sys_uptime(void);
107 extern uint64 sys_mmap(void);
108 extern uint64 sys_munmap(void);
109

```

```

Open  ▾  [+]
```

```

syscall.c
~/Desktop/xv6-labs-2021/kernel

129 [SYS_link]      sys_link,
130 [SYS_mkdir]     sys_mkdir,
131 [SYS_close]     sys_close,
132 [SYS_mmap]      sys_mmap,
133 [SYS_munmap]    sys_munmap,
134 };

```

在 user/usys.pl 中添加用户态的系统调用接口：entry("mmap"); 和 entry("munmap");

```

Open  ▾  [+]
```

```

usys.pl
~/Desktop/xv6-labs-2021/user

37 entry("sleep");
38 entry("uptime");
39 entry("mmap");
40 entry("munmap");

```

在 user/user.h 中声明系统调用的用户接口：void *mmap(void *, int, int, int, int, int); 和 int munmap(void *, int);

```

Open  ▾  [+]
```

```

user.h
~/Desktop/xv6-labs-2021/user

41 int memcmp(const void *, const void *, uint);
42 void *memcpy(void *, const void *, uint);
43 void *mmap(void *, int, int, int, int, int);
44 int munmap(void *, int);

```

(3) 定义 vm_area 结构体

在 kernel/proc.h 中定义 vm_area 结构体，记录映射的内存区域信息。

```

struct vm_area {
    uint64 addr;
    int len;
    int prot;
    int flags;
    int offset;

```

```
    struct file* f;
};
```

在 `struct proc` 中添加 `vma` 数组，用于记录每个进程的虚拟内存区域。

```
struct vm_area vma[NVMA];
```

(4) 实现 `mmap` 系统调用

在 `kernel/sysfile.c` 中实现 `sys_mmap()`，处理用户传递的参数，如 `len` 和 `offset`。

进行参数检查，确保 `flags` 参数只能为 `MAP_SHARED` 或 `MAP_PRIVATE`。

分配 `vm_area` 结构体，并记录映射信息。

使用 `Lazy allocation`，实际内存页面在陷阱处理页面错误时分配。

返回分配的地址，失败则返回 `-1`。

// 实现系统调用 `sys_mmap`，用于将文件映射到进程的地址空间

```
uint64 sys_mmap(void) {
    uint64 addr;          // 映射的起始地址

    int len, prot, flags, offset; // 映射长度、保护标志、映射标志和偏移量

    struct file *f;       // 文件指针

    struct vm_area *vma = 0; // 虚拟内存区域

    struct proc *p = myproc(); // 获取当前进程

    int i;

    // 获取系统调用参数

    if (argaddr(0, &addr) < 0 || argint(1, &len) < 0
        || argint(2, &prot) < 0 || argint(3, &flags) < 0
        || argfd(4, 0, &f) < 0 || argint(5, &offset) < 0) {
        return -1; // 参数获取失败，返回错误码 -1
    }

    // 检查映射标志是否有效

    if (flags != MAP_SHARED && flags != MAP_PRIVATE) {
        return -1; // 无效的映射标志，返回错误码 -1
    }
}
```

```

// 如果映射标志为 MAP_SHARED, 则文件必须是可写的
if (flags == MAP_SHARED && f->writable == 0 && (prot &
PROT_WRITE)) {
    return -1; // 文件不可写且要求写入保护, 返回错误码 -1
}

// 偏移量必须是页面大小的整数倍
if (len < 0 || offset < 0 || offset % PGSIZE) {
    return -1; // 长度、偏移量无效, 返回错误码 -1
}

// 为映射的内存分配一个虚拟内存区域
for (i = 0; i < NVMA; ++i) {
    if (!p->vma[i].addr) {
        vma = &p->vma[i]; // 找到一个空闲的虚拟内存区域
        break;
    }
}

if (!vma) {
    return -1; // 没有找到空闲的虚拟内存区域, 返回错误码 -1
}

// 假设 addr 总是为 0, 内核选择一个页面对齐的地址进行映射
addr = MMAPMINADDR; // 初始化地址为最小映射地址
for (i = 0; i < NVMA; ++i) {
    if (p->vma[i].addr) {
        // 获取当前映射内存的最大地址
        addr = max(addr, p->vma[i].addr + p->vma[i].len);
    }
}

```



```

    }
}

addr = PGROUNDUP(addr); // 将地址对齐到页面边界
if (addr + len > TRAPFRAME) {
    return -1; // 映射区域超出允许范围, 返回错误码 -1
}

// 设置虚拟内存区域的属性
vma->addr = addr;
vma->len = len;
vma->prot = prot;
vma->flags = flags;
vma->offset = offset;
vma->f = f;
filedup(f); // 增加文件的引用计数

return addr; // 返回映射的起始地址
}

```

(5) 实现 munmap 系统调用

在 `kernel/sysfile.c` 中实现 `sys_munmap()`, 取消映射指定的页面, 并处理文件的修改和写回操作。

提取参数 `addr` 和 `length`。检查参数, 确保 `length` 非负, `addr` 是 `PGSIZE` 的整数倍。根据 `addr` 和 `length` 找到对应的 `vm_area` 结构体。若 `length` 为 0, 直接返回成功。

判断 `vm_area` 的标志位, 若有 `MAP_SHARED`, 则需要将修改的页面写回文件。使用脏页标志位 `PTE_D` 判断哪些页面需要写回。使用类似于 `filewrite()` 的方法, 分批次将修改写回文件。使用 `uvmunmap()` 取消用户页表中的映射。

更新 `vm_area` 结构体。

```

uint64 sys_munmap(void) {
    uint64 addr, va;           // 地址变量
    int len;                   // 长度变量

```

```

struct proc *p = myproc(); // 获取当前进程的信息

struct vm_area *vma = 0; // 虚拟内存区域结构指针

uint maxsz, n, n1; // 最大尺寸和循环计数器

int i; // 循环变量


// 从用户空间获取地址和长度参数
if (argaddr(0, &addr) < 0 || argint(1, &len) < 0) {
    return -1;
}

// 检查地址是否对齐以及长度是否合法
if (addr % PGSIZE || len < 0) {
    return -1;
}


// 查找对应的虚拟内存区域（VMA）
for (i = 0; i < NVMA; ++i) { // NVMA 是虚拟内存区域数组的大小
    if (p->vma[i].addr && addr >= p->vma[i].addr
        && addr + len <= p->vma[i].addr + p->vma[i].len) {
        vma = &p->vma[i]; // 找到匹配的 VMA
        break;
    }
}

if (!vma) { // 如果没有找到匹配的 VMA
    return -1;
}


if (len == 0) { // 如果长度为零，则不需要做任何事情
    return 0;
}

```

```

// 如果是共享映射，则需要将脏页写回文件
if ((vma->flags & MAP_SHARED)) {
    // 计算可以写入磁盘的最大尺寸
    maxsz = ((MAXOPBLOCKS - 1 - 1 - 2) / 2) * BSIZE;
    for (va = addr; va < addr + len; va += PGSIZE) { // 遍历每一页
        if (uvmgetdirty(p->pagetable, va) == 0) { // 如果页面不是脏页，则跳
过
            continue;
        }
        // 只有脏页才会写回到映射的文件中
        n = min(PGSIZE, addr + len - va); // 当前页需要写回的数据量
        for (i = 0; i < n; i += n1) {
            n1 = min(maxsz, n - i); // 计算本次写回的最大尺寸
            begin_op(); // 开始一个磁盘操作
            ilock(vma->f->ip); // 加锁以保护文件
            if (writei(vma->f->ip, 1, va + i, va - vma->addr + vma->offset + i,
n1) != n1) {
                iunlock(vma->f->ip); // 解锁文件
                end_op(); // 结束磁盘操作
                return -1; // 出错返回 -1
            }
            iunlock(vma->f->ip); // 解锁文件
            end_op(); // 结束磁盘操作
        }
    }
}

// 从页表中移除映射
uvmunmap(p->pagetable, addr, (len - 1) / PGSIZE + 1, 1);

```

```

// 更新虚拟内存区域信息

if (addr == vma->addr && len == vma->len) { // 如果整个区域都被解除映射

    vma->addr = 0;

    vma->len = 0;

    vma->offset = 0;

    vma->flags = 0;

    vma->prot = 0;

    fclose(vma->f); // 关闭文件

    vma->f = 0; // 文件指针清零

} else if (addr == vma->addr) { // 如果是从区域的开始位置解除映射

    vma->addr += len;

    vma->offset += len;

    vma->len -= len;

} else if (addr + len == vma->addr + vma->len) { // 如果是从区域的结束位置解除映射

    vma->len -= len;

} else {

    panic("unexpected munmap"); // 如果解除映射的位置不符合预期

}

return 0; // 成功返回 0

}

```

(6) 处理页面错误

在 `kernel/trap.c` 中修改 `usertrap()` 函数，处理因访问 `mmap` 区域发生的页面错误。实现惰性分配，分配物理页面并从文件读取数据。

根据 `r_scause()` 值，检查是否发生页错误，可能的值为 12、13 和 15。

对于存储错误（Store Page Fault），设置脏页标志位 `PTE_D`。

执行惰性分配，使用 `kalloc()` 分配物理页，并使用 `memset()` 清空物理页。

使用 `readi()` 从文件中读取数据到物理页，大小为 `PGSIZE`，对文件 `inode` 进行加锁。

设置 PTE 权限标志位为对应的读、写、执行权限。

使用 `mappages()` 将物理页映射到用户进程的虚拟地址。

```
void usertrap(void) {
    int which_dev = 0;

    if ((r_sstatus() & SSTATUS_SPP) != 0)
        panic("usertrap: not from user mode");

    // 发送中断和异常到 `kerneltrap()`, 因为我们现在在内核模式。
    w_stvec((uint64)kernelvec);

    struct proc *p = myproc();

    // 保存用户程序计数器。
    p->trapframe->epc = r_sepc();

    if (r_scause() == 8) { // 系统调用
        if (p->killed)
            exit(-1);

        // `sepc` 指向 `ecall` 指令, 但我们需要返回到下一条指令。
        p->trapframe->epc += 4;

        // 中断会改变 `sstatus` 等寄存器,
        // 因此在处理这些寄存器之前不要启用中断。
        intr_on();

        syscall();
    } else if (r_scause() == 12 || r_scause() == 13 || r_scause() == 15) { // 分页
```

错误

```

char *pa;

uint64 va = PGROUNDDOWN(r_stval()); // 故障虚拟地址

struct vm_area *vma = 0;

int flags = PTE_U;

int i;


// 查找虚拟内存区域 (VMA)
for (i = 0; i < NVMA; ++i) {
    // 类似于 Linux 的 `mmap`, 它可以修改映射页末尾的剩余字节
    if (p->vma[i].addr && va >= p->vma[i].addr && va < p->vma[i].addr +
p->vma[i].len) {
        vma = &p->vma[i];
        break;
    }
}

if (!vma) {
    goto err;
}


// 设置映射页的 PTE 写标志和脏页标志

if (r_scause() == 15 && (vma->prot & PROT_WRITE) &&
walkaddr(p->pagetable, va)) {
    if (uvmsetdirtywrite(p->pagetable, va)) {
        goto err;
    }
} else {
    if ((pa = kalloc()) == 0) {
        goto err;
    }

    memset(pa, 0, PGSIZE); // 清零物理页

```

```

        ilock(vma->f->ip); // 锁定文件
        if (readi(vma->f->ip, 0, (uint64)pa, va - vma->addr + vma->offset,
PGSIZE) < 0) {
            iunlock(vma->f->ip);
            goto err;
        }
        iunlock(vma->f->ip); // 解锁文件

        if ((vma->prot & PROT_READ)) {
            flags |= PTE_R; // 添加读权限标志
        }

        // 只有在存储分页故障且映射页可写的情况下才设置 PTE 写标志
        和脏页标志
        if (r_scause() == 15 && (vma->prot & PROT_WRITE)) {
            flags |= PTE_W | PTE_D; // 添加写和脏页标志
        }
        if ((vma->prot & PROT_EXEC)) {
            flags |= PTE_X; // 添加执行权限标志
        }
        if (mappages(p->pagetable, va, PGSIZE, (uint64)pa, flags) != 0) {
            kfree(pa);
            goto err;
        }
    }
} else if ((which_dev = devintr()) != 0) { // 设备中断
    // ok
} else {
err:
    printf("usertrap(): 意外的 scause %p 进程 ID=%d\n", r_scause(), p->pid);

```



```

        printf("          sepc=%p stval=%p\n", r_sepc(), r_stval());
        p->killed = 1;
    }

    if (p->killed)
        exit(-1);

    // 如果这是定时器中断，放弃 CPU。
    if (which_dev == 2)
        yield();

    usertrapret();
}

```

(7) 设置脏页标志位:

在 `kernel/riscv.h` 中定义脏页标志位 `PTE_D`。

实现 `uvmgetdirty()` 和 `uvmsetdirtywrite()` 函数。

(8) 更新系统调用

在 `kernel/proc.c` 中修改 `exit()` 和 `fork()` 函数，确保进程退出时取消所有映射，并在进程创建时正确复制映射区域。

(9) 编译和运行:

在终端中执行 `make qemu` 编译并运行 `xv6`。在命令行中输入 `mmaptest`，产生结果如下，符合预期。

```
hart 2 starting
hart 1 starting
init: starting sh
$ mmaptest
mmap_test starting
test mmap f
test mmap f: OK
test mmap private
test mmap private: OK
test mmap read-only
test mmap read-only: OK
test mmap read/write
test mmap read/write: OK
test mmap dirty
test mmap dirty: OK
test not-mapped unmap
test not-mapped unmap: OK
test mmap two files
test mmap two files: OK
mmap_test: ALL OK
fork_test starting
fork_test OK
mmaptest: all tests succeeded
```

10.1.3 实验中遇到的问题和解决办法

(1) 问题 1: 如何实现 `uvmsetdirtywrite()` 函数, 以便正确设置脏页标志位。

解决方法: 在页发生写入操作时, 硬件会将此标志位置位, 表示页面已被修改。需要考虑 `PTE_W` 标志位, 它表示页面是否可写。如果页面是只读的, 这个标志位将不会被设置。使用按位或操作将 `PTE_D` 和 `PTE_W` 标志位设置在该页表项中, 这将标识此页为脏页并且可写。

(2) 问题 2: 在实现 `munmap` 时, 需要处理文件是否被映射为只读但使用了 `MAP_SHARED` 标志位的情况。在这种情况下, 文件的内容不能修改, 但仍然需要执行操作。

解决办法: 确保在处理 `MAP_SHARED` 映射时, 检查映射权限 `prot` 是否包含写权限 `PROT_WRITE`。对于只读文件的写操作, 应该返回错误。

10.1.4 实验心得

面对 `mmap` 和 `munmap` 的实现细节时, 我感到非常困惑。尤其是如何在内核中正确处理用户空间传来的参数, 以及如何在虚拟内存系统中正确地映射文件。随着不断地调试和测试, 我逐渐理解了虚拟内存管理的基本概念, 比如如何通过页表实现虚拟地址到物理地址的转换, 以及如何利用懒加载技术来优化内存的使用。

10.2 Lab10 实验成绩

```
== Test running mmaptest ==
$ make qemu-gdb
(5.0s)
== Test  mmaptest: mmap f ==
mmaptest: mmap f: OK
== Test  mmaptest: mmap private ==
mmaptest: mmap private: OK
== Test  mmaptest: mmap read-only ==
mmaptest: mmap read-only: OK
== Test  mmaptest: mmap read/write ==
mmaptest: mmap read/write: OK
== Test  mmaptest: mmap dirty ==
mmaptest: mmap dirty: OK
== Test  mmaptest: not-mapped unmap ==
mmaptest: not-mapped unmap: OK
== Test  mmaptest: two files ==
mmaptest: two files: OK
== Test  mmaptest: fork_test ==
mmaptest: fork_test: OK
== Test usertests ==
$ make qemu-gdb
usertests: OK (477.3s)
(Old xv6.out.usertests failure log removed)
== Test time ==
time: OK
Score: 140/140
```

实验小结:

通过本次实验，我深入学习了操作系统虚拟内存管理和文件系统的基础知识。我成功实现了 `mmap` 和 `munmap` 系统调用，掌握了如何将文件映射到进程的虚拟地址空间，并能够正确处理共享和私有映射的不同场景。实验中，我理解了懒加载机制的重要性，即在页面首次访问时才进行物理内存的分配和内容填充。此外，我还学会了如何处理页面错误，以及如何通过合理的标志位设置来识别和处理脏页。这次实验不仅增强了我对操作系统底层机制的理解，还提升了我分析问题和解决问题的能力。