

# 操作系统JOS实习第六次报告

张弛 00848231,  
zhangchitc@gmail.com

May 30, 2011

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Initialization and transmitting packets</b>	<b>2</b>
2.0.1	The Timer Environment . . . . .	2
2.0.2	The Output Environment . . . . .	3
2.0.3	The Input Environment . . . . .	3
2.1	The Network Interface Card . . . . .	3
2.1.1	PCI Interface . . . . .	4
2.1.2	E100 Reset . . . . .	11
2.1.3	E100 Structure . . . . .	14
2.1.4	DMA Rings . . . . .	14
2.2	Device Driver Organization . . . . .	15
2.3	Transmitting Packets . . . . .	15
2.3.1	C Structures . . . . .	17
2.4	Transmitting Packets: Network Server . . . . .	31
<b>3</b>	<b>Receiving packets and the web server</b>	<b>32</b>
3.1	Receiving Packets . . . . .	32
3.2	Receiving Packets: Network Server . . . . .	40
3.3	The Web Server . . . . .	46

## 1 Introduction

此次Lab是所有JOS实验中最恶心的一次，需要阅读的东西超过了以前所有的总和。所以请做好准备。

在真正下手之前，最好的话请完整的将MIT的材料完整的读一遍，对各个名词和部分有个大致的印象。其他要读的材料还有很多，具体的部分我在报告中会着重提到。

纵观全局，这次Lab的最大难点就是在于你需要**从零开始**写出一个E100网卡的驱动程序。这个驱动程序从Web Server接收IPC调用向网卡发送数据，然后从网卡接收数据发回给Web Server。这和我们以前的实验都不一样，以前都是给出了结构的框架，我们只需要针对一个具体的功能函数进行细节的填补即可，相关的数据结构、接口设置都为我们设计好了。这次需要我们从头到尾完成整个网卡驱动，困难可想而知。

除开网卡之外的部分都相对简单。因此我们这篇报告重点介绍如何完成这个网卡驱动。操纵网卡我们需要了解的方面有：

1. 一个PCI设备在JOS中的设置和相关数据结构
2. 扫描和初始化网卡
3. 网卡的关键结构
4. 网卡如何和操作系统交互数据
5. 如何对网卡发送指令

报告在后面会一步一步从课程给出的资料中抽取出这些细节。同时我在写的过程中参考了<http://code.google.com/p/os-xv6-network>项目主页提供的一份源代码，说的比较极端的话我基本是完全照搬了它的代码，不过这并不影响，学习别人优秀的代码本来就是编程中获得提高最有效的方式。重要的仍是个人的理解。

## 2 Initialization and transmitting packets

对于Network Server的架构，我们只需要大致了解模块即可，这次的实验很少需要对Server进行大规模的修改。

### 2.0.1 The Timer Environment

**Exercise 1.** Add a call to `time_tick` for every clock interrupt in `kern/trap.c`. Implement `sys_time_msec` and add it to `syscall` in `kern/syscall.c` so that user space has access to the time.

这个Exercise太简单了就不贴代码了，唯一需要注意的是在测试用户程序testtime之前由于我们还没有实现网络服务器的部分，所以需要注释掉JOS载入网络服务器的部分：

```
kern/init.c: i386_init()
1      // Should always have an idle process as first one.
2      ENV_CREATE(user_idle);
3
4      // Start fs.
5      ENV_CREATE(fs_fs);
6
7  #if !defined(TEST_NO_NS)
8      // Start ns.
9      //ENV_CREATE(net_ns);
10 #endif
```

这样运行客户程序才不会出错。

## 2.0.2 The Output Environment

## 2.0.3 The Input Environment

这两部分在材料中也提到了我们需要先实现驱动程序和系统调用部分才可以完成。所以我们先放下他们关注最重要的驱动部分。

## 2.1 The Network Interface Card

**Exercise 2.** Browse the Intel 82559 page and look at these two documents:

1. Intel 8255x 10/100 Mbps Ethernet Controller Family Open Source Software Developer Manual (local copy)
2. 82559ER Fast Ethernet PCI Controller Datasheet (local copy)

Do not worry about the details in your first pass. It is more important to read this assignment write-up first to get a high level pictures of how the Intel chip is organized and what is needed to create a device driver.

When you do read the open source developer manual in depth, glance over Section 4 to learn about the PCI interface but pay very close attention to Section 6 as it deals with the Software Interface. In fact, most everything you need is in Section 6. Use the datasheet solely as a reference if you find the developer manual vague.

A simple E100 driver needs only a fraction of the features and interfaces that the card provides. When you're reading through the developer manual, think carefully about the easiest way to interface with the card. You're of course welcome to use its more advanced, high-performance features (in fact, some of the challenge exercises ask you to do exactly this), but it's a good idea to get a basic driver working first.

The acronyms in both documents can get overwhelming. Consult the glossary at the end of this lab assignment for some help.

打开资料是不是已经晕了？这是我觉得这个Lab设计的非常不好的原因之一，一开始就甩出一大堆手册要我们看，而不是给出一个纵览性的结构和指南。先放下吧，后面我们再慢慢说具体要看这些资料的哪些部分。

### 2.1.1 PCI Interface

在使用网卡之前，我们先要通知硬件系统扫描出所有的PCI设备，并且找出其中的E100网卡，对其进行初始化。这是这部分任务的最终目的，为我们操作网卡做好准备工作。这里我们将大致介绍一下PCI设备的一些基础知识：

什么是PCI设备：

PCI是外围设备互连（Peripheral Component Interconnect）的简称，是在目前的计算机系统中得到了非常广泛应用的通用总线接口标准。

- 在一个PCI系统中，最多可以有256根PCI总线，一般的主机上只会用到其中很少的几条，比如Bus 0和Bus 1。
- 在一根PCI总线上可以连接多个物理设备，可以是一个网卡、显卡声卡等等。最多不超过32个。一般来说因为物理特性的限制，一条总线上不会有太多的设备。
- 一个PCI物理设备可以有多个功能，比如同时提供视频解析和声音解析。最多提供8个功能。
- 每个功能对应一个256 bytes的PCI Configuration Space，这个在我们接下来对网卡进行初始化的时候会特别提到这个东西的。

所以对于一个PCI设备的具体功能，我们可以使用**总线号：设备号：功能号**来对其进行定位，比如在Ubuntu下我们使用lspci命令，就可以得到这样的输出：

```
zhangchi@zhangchi-vostrol400:~$ lspci
00:00.0 Host bridge: Intel Corporation Mobile PM965/GM965/GL960 Memory Controller Hub (rev 0c)
00:02.0 VGA compatible controller: Intel Corporation Mobile GM965/GL960 Integrated Graphics Controller (rev 0c)
00:02.1 Display controller: Intel Corporation Mobile GM965/GL960 Integrated Graphics Controller (rev 0c)
00:1a.0 USB Controller: Intel Corporation 82801H (ICH8 Family) USB UHCI Controller #4 (rev 02)
00:1a.1 USB Controller: Intel Corporation 82801H (ICH8 Family) USB UHCI Controller #5 (rev 02)
00:1a.7 USB Controller: Intel Corporation 82801H (ICH8 Family) USB2 EHCI Controller #2 (rev 02)
00:1b.0 Audio device: Intel Corporation 82801H (ICH8 Family) HD Audio Controller (rev 02)
00:1c.0 PCI bridge: Intel Corporation 82801H (ICH8 Family) PCI Express Port 1 (rev 02)
00:1c.1 PCI bridge: Intel Corporation 82801H (ICH8 Family) PCI Express Port 2 (rev 02)
00:1c.3 PCI bridge: Intel Corporation 82801H (ICH8 Family) PCI Express Port 4 (rev 02)
00:1c.5 PCI bridge: Intel Corporation 82801H (ICH8 Family) PCI Express Port 6 (rev 02)
00:1d.0 USB Controller: Intel Corporation 82801H (ICH8 Family) USB UHCI Controller #1 (rev 02)
00:1d.1 USB Controller: Intel Corporation 82801H (ICH8 Family) USB UHCI Controller #2 (rev 02)
00:1d.2 USB Controller: Intel Corporation 82801H (ICH8 Family) USB UHCI Controller #3 (rev 02)
00:1d.7 USB Controller: Intel Corporation 82801H (ICH8 Family) USB2 EHCI Controller #1 (rev 02)
00:1e.0 PCI bridge: Intel Corporation 82801 Mobile PCI Bridge (rev f2)
00:1f.0 ISA bridge: Intel Corporation 82801HEM (ICH8M) LPC Interface Controller (rev 02)
00:1f.1 IDE interface: Intel Corporation 82801HEM/HEM (ICH8M/ICH8M-E) IDE Controller (rev 02)
00:1f.2 IDE interface: Intel Corporation 82801HEM/HEM (ICH8M/ICH8M-E) SATA IDE Controller (rev 02)
00:1f.3 SMBus: Intel Corporation 82801H (ICH8 Family) SMBus Controller (rev 02)
03:01.0 FireWire (IEEE 1394): Ricoh Co Ltd R5C832 IEEE 1394 Controller (rev 05)
03:01.1 SD Host controller: Ricoh Co Ltd R5C822 SD/SDIO/MMC/MS/MSPro Host Adapter (rev 22)
03:01.2 System peripheral: Ricoh Co Ltd R5C843 MMC Host Controller (rev 12)
03:01.3 System peripheral: Ricoh Co Ltd R5C592 Memory Stick Bus Host Adapter (rev 12)
03:01.4 System peripheral: Ricoh Co Ltd xD-Picture Card Controller (rev ff)
09:00.0 Ethernet controller: Broadcom Corporation NetLink BCM5906M Fast Ethernet PCI Express (rev 02)
0c:00.0 Network controller: Intel Corporation PRO/Wireless 3945ABG [Golan] Network Connection (rev 02)
```

看到前面用冒号和点分割的数字了么，就是以总线号、设备号以及功能号进行标识的。可以看到我的机器上用到了0、3、9、12这四条PCI总线。

#### PCI配置寄存器：

每一个PCI设备都有它映射的内存地址空间和它的I/O区域。除此之外，PCI设备还有它的配置寄存器（即Configuration Space）。对于所有的PCI设备，配置地址空间一共256 bytes，其中前64 bytes是标准化的，它提供了厂商号，设备号，版本号等信息，唯一标识一个PCI设备。同时，它也提供了最多可多达6个的I/O地址区域，每个区域可以是内存也可以是I/O地址。这几个I/O地址区域是驱动程序找到设备映射到内存和I/O空间的具体位置的唯一途径。关于这64个字节的配置空间的详细情况，可以参考下图：

register	bits 31-24	bits 23-16	bits 15-8	bits 7-0
00	Device ID		Vendor ID	
04	Status		Command	
08	Class code	Subclass	Prog IF	Revision ID
0C	BIST	Header type	Latency Timer	Cache Line Size
10	Base address #0 (BAR0)			
14	Base address #1 (BAR1)			
18	Base address #2 (BAR2)			
1C	Base address #3 (BAR3)			
20	Base address #4 (BAR4)			
24	Base address #5 (BAR5)			
28	Cardbus CIS Pointer			
2C	Subsystem ID		Subsystem Vendor ID	
30	Expansion ROM base address			
34	Reserved			Capabilities Pointer
38	Reserved			
3C	Max latency	Min Grant	Interrupt PIN	Interrupt Line

#### PCI设备启动过程：

在系统引导阶段，PCI硬件设备保持未激活状态，每个设备都没有被分配内存空间和I/O端口。但每个PCI主板均配备有能够处理PCI的固件（比如BIOS），固件通过读写PCI控制器中的寄存器，提供了对设备配置地址空间的访问。系统启动以后，固件通过扫描每个PCI设备，通过读取他们的配置地址空间，为每个设备分配相应的内存和I/O端口，为后面硬件驱动程序做好准备。

### 如何使用PCI设备：

当我们想查询一个特定PCI设备的配置地址空间时，我们需要向I/O地址[0cf8, 0cfb] 写入一个4 bytes查询码指定总线号：设备号：功能号以及其配置地址空间中的查询位置。那么PCI Host Bridge将监听对于这个I/O端口的写入并在接受到写入数据后将相应的查询结果写入到[0cfc, 0cff]，我们从这个地址读出一个32位整数表示查询到的相应信息。

查询配置地址空间时，我们一般会从其6个BARS中得到特定设备的控制端口和数据端口信息，那么只要在初始化时将这些端口地址保存下来，就可以在PCI硬件驱动程序中通过向这些端口输入输出数据来达到控制PCI设备的目的了

更多更详细的内容可以参考网站：[http://xwindow.angelfire.com/page13\\_1.html](http://xwindow.angelfire.com/page13_1.html)，上面介绍的内容已经足够我们理解JOS中的相应代码了。

接下来我们看JOS中是如何对PCI设备进行编程的，这部分模块主要定义在kern/pci.c中，JOS在系统初始化时调用其中的pci\_init() 进行设备初始化，首先来看一些最基本的东西：

```

kern/pci.c

36 static void
37 pci_conf1_set_addr(uint32_t bus,
38                   uint32_t dev,
39                   uint32_t func,
40                   uint32_t offset)
41 {
42     assert(bus < 256);
43     assert(dev < 32);
44     assert(func < 8);
45     assert(offset < 256);
46     assert((offset & 0x3) == 0);
47
48     uint32_t v = (1 << 31) | // config-space
49                 (bus << 16) | (dev << 11) | (func << 8) | (offset);
50     outl(pci_conf1_addr_ioport, v);
51 }
52
53 static uint32_t
54 pci_conf_read(struct pci_func *f, uint32_t off)
55 {
56     pci_conf1_set_addr(f->bus->busno, f->dev, f->func, off);
57     return inl(pci_conf1_data_ioport);
58 }
59
60 static void
61 pci_conf_write(struct pci_func *f, uint32_t off, uint32_t v)
62 {
63     pci_conf1_set_addr(f->bus->busno, f->dev, f->func, off);
64     outl(pci_conf1_data_ioport, v);
65 }

```

这三个函数是对PCI设备最基本的读状态和写状态的函数

- `pci_conf_read()` 是读取PCI配置地址空间中特定位置的配置值
- `pci_conf_write()` 是设置PCI配置地址空间中特定位置的配置值
- 其中 `pci_conf1_set_addr()` 负责设置需要读写的具体设备。这里涉及到的两个I/O端口定义在了文件最上方：

```
kern/pci.c
12 // PCI "configuration mechanism one"
13 static uint32_t pci_conf1_addr_ioport = 0x0cf8;
14 static uint32_t pci_conf1_data_ioport = 0x0cfc;
```

正是我们前面提到的两个操作PCI设备的I/O端口。

接下来我们看看它是怎么初始化PCI设备的，看到`pic_init()`：

```
kern/pci.c
36 static int
37 pci_scan_bus(struct pci_bus *bus)
38 {
39     int totaldev = 0;
40     struct pci_func df;
41     memset(&df, 0, sizeof(df));
42     df.bus = bus;
43
44     for (df.dev = 0; df.dev < 32; df.dev++) {
45         uint32_t bhlc = pci_conf_read(&df, PCI_BHLC_REG);
46         if (PCI_HDRTYPE_TYPE(bhlc) > 1) // Unsupported or no device
47             continue;
48
49         totaldev++;
50
51         struct pci_func f = df;
52         for (f.func = 0; f.func < (PCI_HDRTYPE_MULTIFN(bhlc) ? 8 : 1); f.func++) {
53             struct pci_func af = f;
54
55             af.dev_id = pci_conf_read(&af, PCI_ID_REG);
56             if (PCI_VENDOR(af.dev_id) == 0xffff)
57                 continue;
58
59             uint32_t intr = pci_conf_read(&af, PCI_INTERRUPT_REG);
60             af.irq_line = PCI_INTERRUPT_LINE(intr);
61
62             af.dev_class = pci_conf_read(&af, PCI_CLASS_REG);
63             if (pci_show_devs)
64                 pci_print_func(&af);
65             pci_attach(&af);
66         }
67     }
68
69     return totaldev;
70 }
71
72 int
73 pci_init(void)
74 {
```

```

75     static struct pci_bus root_bus;
76     memset(&root_bus, 0, sizeof(root_bus));
77
78     return pci_scan_bus(&root_bus);
79 }

```

1. `pci_init()` 中, `root_bus` 被全部清0, 然后交给 `pci_scan_bus()` 扫描这条总线上的所有设备, 说明在 JOS 中 E100 是被放置在 0 号总线上的
2. `pci_scan_bus()` 中顺次查找 0 号总线上的 32 个设备, 如果发现其存在, 那么顺次扫描它们每个功能对应的配置地址空间, 将一些关键的控制参数读入 `pci_func` 进行保存, 其中 `pci_func` 的结构如下:

```

                                kern/pci.h
11 struct pci_func {
12     struct pci_bus *bus;           // Primary bus for bridges
13
14     uint32_t dev;
15     uint32_t func;
16
17     uint32_t dev_id;
18     uint32_t dev_class;
19
20     uint32_t reg_base[6];
21     uint32_t reg_size[6];
22     uint8_t irq_line;
23 };

```

对于网卡驱动来说, 最重要的就是其 `reg_base` 数组, 这是我们用于向 E100 发送命令的地址端口, 在后面的初始化程序中我们需要将其记录下来。

3. 得到 `pci_func` 之后, 它被传入 `pci_attach()` 去查找是否为已存在的硬件, 如果匹配成功, 则使用预设好的程序初始化该硬件

```

                                kern/pci.c
67 static int __attribute__((warn_unused_result))
68 pci_attach_match(uint32_t key1, uint32_t key2,
69                 struct pci_driver *list, struct pci_func *pcif)
70 {
71     uint32_t i;
72
73     for (i = 0; list[i].attachfn; i++) {
74         if (list[i].key1 == key1 && list[i].key2 == key2) {
75             int r = list[i].attachfn(pcif);
76             if (r > 0)
77                 return r;
78             if (r < 0)
79                 cprintf("pci_attach_match: _attaching_"
80                        "%x.%x_ (%p) : _e\n",
81                        key1, key2, list[i].attachfn, r);
82         }
83     }
84     return 0;
85 }
86
87 static int
88 pci_attach(struct pci_func *f)

```



```

89 {
90     return
91     pci_attach_match(PCI_CLASS(f->dev_class),
92                     PCI_SUBCLASS(f->dev_class),
93                     &pci_attach_class[0], f) ||
94     pci_attach_match(PCI_VENDOR(f->dev_id),
95                     PCI_PRODUCT(f->dev_id),
96                     &pci_attach_vendor[0], f);
97 }

```

默认的情况下系统初始是没有定义的PCI设备的，看到pci\_attach\_vendor就知道：

```

                                kern/pci.c
31 // pci_attach_vendor matches the vendor ID and device ID of a PCI device
32 struct pci_driver pci_attach_vendor[] = {
33     { 0, 0, 0 },
34 };

```

待会我们初始化的时候要将我们定义的网卡加入这个数组。

到这里PCI设备的初始化就结束了，接下来我们要尝试写一个E100网卡的初始化过程。

**Exercise 3.** Implement an attach function to initialize the 82559ER. Add an entry to the pci\_attach\_vendor array in kern/pci.c to trigger your function if a matching PCI device is found. The vendor ID and device ID for the 82559ER can be found in Section 4 of the developer manual. You should also see these listed when JOS scans the PCI bus while booting.

After enabling the E100 device via pci\_func\_enable, your attach function should record the IRQ line and base I/O port assigned to the device so you'll be able to communicate with the E100.

We have provided the kern/e100.c and kern/e100.h files for you so that you do not need to mess with the make system. You may still need to include the e100.h file in other places in the kernel.

When you boot your kernel, you should see it print that the PCI function of the E100 card was enabled. Your code should now pass the pci attach test of make grade.

第一步我们查阅手册得到E100的Vendor ID为8086h，Device ID为1229h，然后将初始化程序作为驱动程序的一部分定义在kern/e100.c中，先撰写头文件kern/e100.h：

```

                                kern/e100.h
1 #ifndef JOS_KERN_E100_H
2 #define JOS_KERN_E100_H
3
4 #include <kern/pci.h>
5
6 #define E100_VENDOR            0x8086
7 #define E100_DEVICE            0x1209
8
9 int e100_attach(struct pci_func *pcif);

```

```

10
11 #endif // JOS_KERN_E100_H

```

然后是主过程，定义了一个e100记录其相应的设备信息：

```

                                kern/e100.c
1 // LAB 6: Your driver code here
2
3 #include <inc/x86.h>
4 #include <inc/stdio.h>
5
6 #include <kern/e100.h>
7
8 struct pci_func e100;
9
10 int
11 e100_attach(struct pci_func *pcif)
12 {
13     pci_func_enable(pcif);
14     e100.bus = pcif->bus;
15     e100.dev_id = pcif->dev_id;
16     e100.dev_class = pcif->dev_class;
17     int i;
18     for (i = 0; i < 6; i++) {
19         e100.reg_base[i] = pcif->reg_base[i];
20         e100.reg_size[i] = pcif->reg_size[i];
21         cprintf ("zhangchi: The %dth Bar: base = %x, size = %x\n",
22                 i, e100.reg_base[i], e100.reg_size[i]);
23     }
24     e100.irq_line = pcif->irq_line;
25
26     return 0;
27 }

```

接下来修改kern/pci.c中的pci\_attach\_vendor数组，把我们的E100初始化程序添加进去：

```

                                kern/pci.c
31 // pci_attach_vendor matches the vendor ID and device ID of a PCI device
32 struct pci_driver pci_attach_vendor[] = {
33     { E100_VENDOR, E100_DEVICE, &e100_attach },
34     { 0, 0, 0 },
35 };

```

然后make qemu启动JOS，应该能看到E100网卡被顺利激活了：

```

enabled interrupts: 1 2
        Setup timer interrupts via 8259A
enabled interrupts: 0 1 2
        unmasked timer interrupt
PCI: 00:00.0: 8086:1237: class: 6.0 (Bridge device) irq: 0
PCI: 00:01.0: 8086:7000: class: 6.1 (Bridge device) irq: 0
PCI: 00:01.1: 8086:7010: class: 1.1 (Storage controller) irq: 0
PCI: 00:01.3: 8086:7113: class: 6.80 (Bridge device) irq: 9
PCI: 00:02.0: 1013:00b8: class: 3.0 (Display controller) irq: 0
PCI: 00:03.0: 8086:1209: class: 2.0 (Network controller) irq: 11
PCI function 00:03.0 (8086:1209) enabled
zhangchi: The 0th Bar: base = f2020000, size = 1000
zhangchi: The 1th Bar: base = c040, size = 40
zhangchi: The 2th Bar: base = f2040000, size = 20000

```

```

zhangchi: The 3th Bar: base = 0, size = 0
zhangchi: The 4th Bar: base = 0, size = 0
zhangchi: The 5th Bar: base = 0, size = 0
FS is running
FS can do I/O

```

看到那个以c开头的第二个地址了么，这个明显不是内存地址，它应该就是我们以后进行操作的I/O端口c040，空间大小为40h = 64 bytes

### 2.1.2 E100 Reset

**Exercise 4.** Add code to your attach function to reset the 82559ER. If you set the `-debug-e100` flag, QEMU should tell you if the reset was successfully. It will print something like this after JOS starts scanning the PCI bus:

```
EE100  nic_reset          0xacea498
```

There will also be a few `nic_reset`'s before JOS starts; those are the BIOS itself resetting the device.

这一段的MIT提供的资料是相对来说比较详尽的，还记得我们前面打印出来的c040地址么？这个地方就是我们要写入CSR的端口。CSR(Control/Status Registers)是我们对于E100网卡的控制字，如前面所说，它是一个64 bytes的地址空间，其中我们最需要关注的是它的前12个bytes，称为SCB(System Control Block)，我们对网卡的主要控制主要是对于SCB相应参数的进行设置。其布局如下，更详细可以参考提供的手册里6.3.1 Control / Status Registers (CSR)：

Upper Word		Lower Word		Offset
31	16	15	0	
SCB Command Word		SCB Status Word		0h
SCB General Pointer				4h
PORT				8h
EEPROM Control Register		Reserved		Ch
MDI Control Register				10h
RX DMA Byte Count				14h
PMDR	Flow Control Register		Reserved	18h
Reserved		General Status	General Control	1Ch
Reserved				20h-2Ch
Function Event Register				30h
Function Event Mask Register				34h
Function Present State Register				38h
Force Event Register				3Ch

对于重设网卡，是采用的PORT Interface的形式进行控制，详见Manual的6.3.3 PORT Interface。这里E100允许我们只对SCB中的PORT设置特定值以后就执行相应的功能，比如说有：

Function	Pointer Field (Bits 31:4)	Opcode (Bits 3:0)
Software Reset	Don't care	0000
Self-test	Self-test results pointer (16 byte alignment)	0001
Selective Reset	Don't care	0010
Dump	Dump area pointer (16 byte alignment)	0011
Dump Wake-up	Dump area pointer (16 byte alignment)	0111

从这个表格看出，我们只需要对PORT字段写入全0就可以达到重启的目的了。但是在进行编码之前，我们需要将一些常用的常数放入头文件，方便后续的程序编写。

```

                                kern/e1000.h
1  #ifndef JOS_KERN_E100_H
2  #define JOS_KERN_E100_H
3
4  #include <kern/pci.h>
5
6  #define E100_VENDOR          0x8086
7  #define E100_DEVICE          0x1209
8
9  #define E100_MEMORY          0
10 #define E100_IO              1
11 #define E100_FLASH           2
12
13 #define CSR_SCB               0x0
14 #define CSR_STATUS            0x0
15 #define CSR_US                0x0
16 #define CSR_STATAK            0x1
17 #define CSR_COMMAND           0x2
18 #define CSR_UC                0x2
19 #define CSR_INT               0x3
20 #define CSR_GP                0x4
21 #define CSR_PORT              0x8
22
23
24 #define PORT_SW_RESET         0x0
25 #define PORT_SELF_TEST        0x1
26 #define PORT_SEL_RESET        0x2
27
28 int e100_attach(struct pci_func *pcif);
29
30 #endif // JOS_KERN_E100_H

```

里面定义了三组常数分别以E100，CSR和PORT开头：

1. 在手册中的4.1 PCI Configuration Space，对于E100而言，PCI配置中提供的6个地址中的前三个分别为
  - (a) CSR Memory Mapped Base Address Register
  - (b) CSR I/O Mapped Base Address Register
  - (c) Flash Memory Mapped Base Address Register

因为我们只使用I/O端口对CSR进行控制，不使用内存地址的原因资料中也提到了，有可能因为编译器的原因使得地址端口失效，所以最

稳固的方法还是使用I/O的方式。这三个地址在初始化时已经被载入到e100.reg\_base[0-2]中了。在使用他们的基址的时候，我们为他们定义了相应的数组索引位置

2. 定义了一系列SCB字段在CSR中的位移，以便于后面我们使用in和out指令对他们进行读写操作
3. 预定义了三条PORT Interface指令

作这项工作中我大量参考了<http://code.google.com/p/os-xv6-network/source/browse/trunk/dev/e100.h>提供的参数，节省了我大量的时间，对作者表示感谢。

然后就可以真正开始对网卡进行重启了：

```

kern/e100.c
1 // LAB 6: Your driver code here
2
3 #include <inc/x86.h>
4 #include <inc/stdio.h>
5
6 #include <kern/e100.h>
7
8 struct pci_func e100;
9
10 static void e100_sw_reset(struct pci_func e100);
11
12 int
13 e100_attach(struct pci_func *pcif)
14 {
15     pci_func_enable(pcif);
16     e100.bus = pcif->bus;
17     e100.dev_id = pcif->dev_id;
18     e100.dev_class = pcif->dev_class;
19     int i;
20     for (i = 0; i < 6; i++) {
21         e100.reg_base[i] = pcif->reg_base[i];
22         e100.reg_size[i] = pcif->reg_size[i];
23     }
24     e100.irq_line = pcif->irq_line;
25
26     e100_sw_reset(e100);
27 }
28
29 static void
30 e100_sw_reset(struct pci_func e100) {
31     outl(e100.reg_base[E100_IO] + CSR_PORT, PORT_SW_RESET);
32
33     // delay about 10us
34     int i = 0;
35     for (i = 0; i < 8; i++) {
36         inb(0x84);
37     }
38 }

```

注意不要忘了按照MIT材料的提示重启后delay一段时间再返回。使用make qemu QEMUEXTRA="-debug-e100"启动JOS应该可以看到网卡的重启消息：

```

EE100  nic_init
EE100  pci_reset          0x9566008
EE100  nic_init          macaddr: 52 54 00 12 34 56
EE100  nic_reset          0x9566008
EE100  nic_selective_reset checksum=0xbe34
EE100  nic_init          model=i82559er,macaddr=52:54:00:12:34:56
EE100  nic_reset          0x9566008
EE100  nic_selective_reset checksum=0xbe34
EE100  pci_mmio_map      region 0, addr=0xf2020000, size=0x00001000, type=8
EE100  pci_map           region 1, addr=0x0000c040, size=0x00000040, type=1
EE100  pci_mmio_map      region 2, addr=0xf2040000, size=0x00020000, type=0
6828 decimal is 15254 octal!
Hooray! Passed all test cases for stdlib!!
Physical memory: 66556K available, base = 640K, extended = 65532K
check_page_alloc() succeeded!
page_check() succeeded!
check_boot_pgdir() succeeded!
enabled interrupts: 1 2
    Setup timer interrupts via 8259A
enabled interrupts: 0 1 2
    unmasked timer interrupt
PCI: 00:00.0: 8086:1237: class: 6.0 (Bridge device) irq: 0
PCI: 00:01.0: 8086:7000: class: 6.1 (Bridge device) irq: 0
PCI: 00:01.1: 8086:7010: class: 1.1 (Storage controller) irq: 0
PCI: 00:01.3: 8086:7113: class: 6.80 (Bridge device) irq: 9
PCI: 00:02.0: 1013:00b8: class: 3.0 (Display controller) irq: 0
PCI: 00:03.0: 8086:1209: class: 2.0 (Network controller) irq: 11
EE100  pci_mmio_map      region 0, addr=0xf2020000, size=0x00001000, type=8
EE100  pci_map           region 1, addr=0x0000c040, size=0x00000040, type=1
EE100  pci_mmio_map      region 2, addr=0xf2040000, size=0x00020000, type=0
PCI function 00:03.0 (8086:1209) enabled
EE100  eeprol100_write4   addr=Port+0 val=0x00000000
EE100  nic_reset          0x9566008
EE100  nic_selective_reset checksum=0xbe34
FS is running

```

可以看到最后两行nic\_reset表明网卡已成功进行了软启动。

### 2.1.3 E100 Structure

其实这里Intel的手册挺让人费解的，我读了以后发现CU其实就是负责发送数据的模块，RU就是接收数据的模块，按道理两个是相对的，那么CU应该称为Transmit Unit才对，Intel却命名为Control Unit，让人感觉这两个模块是分立的，并且CU有控制RU一样。

对于CU和RU的讨论我们要等到将DMA Ring看完以后才能完整的描述，我们先来看DMA Rings

### 2.1.4 DMA Rings

根据资料的描述，DMA Rings就是系统为E100在内存中开辟的一片区域，用于网卡使用DMA缓存当前的收发数据的。一个DMA Rings在申请好以后，就可以将其所在的物理地址通知给E100的DMA控制器，那么E100就可以在不占用CPU的情况下自己根据内存中的数据开始进行收发操作了。**实际情况中内存里应该有两条DMA Rings**，一个专门用于放置待发送的数据（Control Block list, CBL），一个专门放置接受到的数据（Receive Frame Area, RFA）。

因为采取了这样的操作模式，所以提示我们E100是根据DMA Rings中的内容进行工作的。接下来我们会详细讲述利用CU发送数据时数据包格式是如何用Control Block(CB)进行描述的。E100读取CB中的设置，然后进行相应的发送操作，当工作完成后，**通过引发中断或者改变CB相应的状态位来提示系统工作已完成**。我们可以修改CB中相应场位的设置，来改变E100的工作模式。

上面是E100自动进行工作的一种方式，如果我们需要人为的干预E100的运行，**还可以向CSR中的SCB寄存器写入相应的控制指令**，比如终止运行等，来改变其运行状态。

所以，在完成这个部分的工作前，我们需要了解两方面重要的内容：

1. CB (Control Block)的控制设置
2. SCB (System Control Block)的控制设置

在了解DMA Rings 的相关结构以后，我们会结合DMA Rings来进一步阐述这些控制设置在DMA Rings上运行的效果。

## 2.2 Device Driver Organization

## 2.3 Transmitting Packets

前面提到DMA Rings主要分为两种用途：发送和接受。

- 发送数据的DMA Rings是由若干个CB (Control Block)组成的，这些CB通过指针连接成一个环状的结构CBL (Control Block List)。
- 接受数据的DMA Rings由若干个RFD (Receive Frame Descriptor)组成，也是连成环状，称为RFA (Receive Frame Area)

这个阶段我们主要关注发送，发送是由CU模块来完成的。

**CB是一个通用的概念**，即使是只在发送时用到，但是因为CU执行的不止一种命令，所以根据不同命令的需要对CB进行更加细致的规定。只有CB的前三个成员Control, Status和Link在各个应用场景中都是相同的，接下来的Command Specific Data才是根据不同的命令发生变化的。一个CB的结构在Intel的开源手册中的6.4.1.1 General Action Command Format中可以找到具体格式，如下图所示：

Figure 14. General Action Command Format

Offset	Command Word Bits 31:16					Status Word Bits 15:0				
00h	EL	S	I	0000000000	CMD	C	X	OK	XXXXXXXXXXXX	
04h	Link Offset									
08h	Optional Address and Data Fields									

其中有三个位置是特别值得我们关注的：

- S: 如果该位被设置为1, 那么当E100的**CU**执行完此CB的命令之后, 将停止运行进入挂起状态(Suspend), 只有使用SCB对E100下达恢复运行(Resume)指令后, CU模块才会重新运行, 重新运行开始时执行的地址是该CB的Link Address指向的下一CB。注意, 这里要注意**发送和接受是互相独立不干扰的**, 所以CU停止工作的时候有可能RU还在正常运行, 注意明确他们两个结构的概念。
- Status Word Bits中的C位: 在操作CU执行任意CB上的命令时, **程序员首先应该负责手动清除该位上的值**, 将其置为0。那么当CU模块执行完任务后, 这个位置被设置成1。就可以检测某些指令的执行情况了。比如发送数据时由于数据包比较大, 从CU获取到CB开始执行到发送完毕需要耗费一定的时间,
- CMD: CB支持的不同的操作类型, 具体操作类型可以在手册的6.4.2 Specific Action Commands中找到, 大致有以下几种操作:
  - NOP (000b)
  - Individual Address Setup (001b)
  - Configure (010b)
  - Multicast Setup (011b)
  - Transmit (100b)
  - Load Microcode (101b)
  - Dump (110b)
  - Diagnose (111b)

在这次实验中我们在CU模块的操作只会遇到两个: NOP和Transmit, 严格的说只要使用Transmit就好了, 但是在写的时候我们可以用NOP来进行一些测试, 来检测程序的正确性。NOP的作用就是使CU在处理到该CB时什么事情也不作, **但是相应的状态位S、C等等都会对CU的执行有作用**。所以我们可以利用NOP指令设置S位**使其停止在某个CB上**, 然后利用debug的输出信息进行调试, 而不用去理会Transmit指令中其他那些乱七八糟的参数的设置。

看完CB的相关字段说明, 我们可以将这些状态定义到头文件里方便后续的设置和读取判断了:

```
kern/e100.h

1 // Control Block Command
2 #define CBF_EL      0x8000
3 #define CBF_S       0x4000
4 #define CBF_I       0x2000
5
6 #define CBC_NOP      0x0
7 #define CBC_IAS      0x1
8 #define CBC_CONFIG   0x2
9 #define CBC_MAS      0x3
10 #define CBC_TRANSMIT 0x4
11 #define CBC_LOADMC   0x5
```



```

12 #define CBC_DUMP      0x6
13 #define CBC_DIAGNOSE  0x7
14
15 // Control Block Status
16 #define CBS_F          0x0800
17 #define CBS_OK         0x2000
18 #define CBS_C          0x8000

```

### 2.3.1 C Structures

上面描述的是CB的一些通用概念，适用所有命令格式，具体到发送指令Transmit的时候，我们将这样一连串的CB称为TCBs (Transmit Command Blocks)。在MIT资料中已经有一个形象的图为我们展示出了TCB的串联结构，更详细的信息可以在手册中的6.4.2.5 Transmit找到。这个是TCB的一个布局图。

Figure 19. Transmit Command Format

Offset	Command Word Bits 31:16								Status Word Bits 15:0				
00h	EL	S	I	CID	000	NC	SF	100	C	X	OK	U	XXXXXXXXXXXX
04h	Link Address (A31:A0)												
08h	Transmit Buffer Descriptor Array Address												
	TBD Number				Transmit Threshold				EOF	0	Transmit Command Block Byte Count		

可以看到，虽然TCB的结构比CB更细化了，但是我们现在还不需要关注它的设置的细节，第一步我们先考虑在内存中分配空间和建立起相应的TCB结构。这就需要结合CU的工作方式来阐述它是如何读取和使用TCB的了。以下列出了使用TCBs发送数据包的全过程：

1. 首先在内存中建立起TCBs的环状结构
2. 将其中第一个TCB的物理地址写入CU的General Pointer，通知其要操作的TCB所在的内存位置
3. 给CU一个Start指令，CU开始工作
4. 如果环状的TCB中所有的S位都是0，那么CU将根据Link Address无限循环的处理这些TCB，每次执行完成以后，以中断和修改TCB中C状态位来表明任务已完成(后面我们会说明到底是以中断还是轮询查看C状态位来处理消息，现在先不用管)
5. 如果CU碰到了某个TCB的S被设置为1，那么处理完该TCB之后，CU进入Suspend状态，等待用户发送Resume命令恢复运行，那么CU将从当前TCB的下个TCB开始执行

这只是硬件的执行机制，我们现在要考虑在TCBs上实现一个支持多数据包发送的排队系统，应该注意些什么？

1. 首先要记录当前TCB中哪些块是等待被发送的，哪些是可以被重新利用的

2. 其次要做好相应S位设置，让CU在发送完需要发送的包之后就停下来

考虑到上述两点，我们设计出的系统是这样描述TCB的：

```

                                kern/e100.h
1  #define TCB_MAXSIZE      1518
2  #define CB_MAX_NUM       10
3
4  // Transmit Command Blocks
5  struct tcb {
6      uint32_t tcb_tbd_array_addr;
7      uint16_t tcb_byte_count;
8      uint8_t tcb_thrs;
9      uint8_t tcb_tbd_count;
10     char tcb_data[TCB_MAXSIZE];
11 };
12
13 // Control Blocks
14 struct cb {
15     volatile uint16_t cb_status;
16     uint16_t cb_control;
17     uint32_t cb_link;
18
19     union cb_cmd_spec_data {
20         struct tcb tcb;
21     } cb_cmd_spec;
22
23     struct cb *prev, *next;
24     physaddr_t phy_addr;
25 };
26
27 // Control Block List
28 struct cbl {
29     int cb_avail;
30     int cb_wait;
31
32     struct cb *start;
33     struct cb *front, *rear;
34 };

```

解释一下这些结构：

- tcb: 直接按照MIT给出的结构定义的，没什么好说
- cb: 这里有两点值得注意：
  1. union结构cb\_cmd\_spec是根据不同指令的需要说明的，其实这里我们只用到了Transmit指令，所以union中只有一个成员看着比较别扭，但是这样写更具有**扩展性和维护性**，当需要使用到其他指令时，直接在cb\_cmd\_spec里添加其他的数据结构即可
  2. 在cb\_cmd\_spec的后面我们增加了三个成员prev, next和phy\_addr，这个和手册上关于CB的定义是不符的，但是并不影响CU的执行，主要是为了我们自己在后续操作中的方便，
- cbl: 这里维护了几个值，分别说明一下：
  - cb\_avail: 表示当前有多少个闲置的TCB可以用来放置数据以发送

- `cb_wait`: 表示当前有多少个TCB正在处于等待发送状态
- 很明显上面两者相加应该等于所有TCB的总数, 在这里应该是我们定义的`CB_MAX_NUM = 10`
- `start`: 所有TCB中的第一个TCB, 用于开始的时候初始化CU用
- `front`和`rear`: 表示当前正在等待发送的TCB的起始和结尾

那么这个系统是如何根据cbl中定义的结构工作的呢?

1. 初始的时候`front = start, rear = start → prev`, 表示当前等待发送的数据包为空, 并且`cb_avail = CB_MAX_NUM, cb_wait = 0`
2. 当需要发送一个数据包时, 将其添加到`rear`后面, 并且移动`rear`指针, 同时增加`cb_wait`, 减少`cb_avail`
3. 当确认等待数据包发送时, 检查`front`指向TCB的C状态是否为1, 可以的话则将`front`向后移动, 表示`front`指向的数据包已经被CU发送。同时减少`cb_wait`, 增加`cb_avail`

那么相应的边界状态比如队列空或者队列满就可以很容易的通过`cb_wait`和`cb_avail`检测出来了。我们在这里还没有结合控制位中S的设置来说明CU的工作方式, 这个在讲发送数据包的时候会讲到移动指针时会如何进行Suspend位的设置, 具体请参考2.3.1

好了到这里我们已经清楚了一个TCB结构是如何被建立起来并且在后续过程中维护的详细过程, 现在我们可以考虑建立起这样的结构了:

**Exercise 5.** Construct a control DMA ring for the CU to use. You do not need to worry about configuring the device because the default setting are fine. You also do not need to worry about setting up the device MAC address because the emulated E100 has one already configured.

首先为了程序的易于管理, 我们新定义了一个结构`nic`:

```

                                kern/e100.h
1  // Network Interface Card
2  struct nic {
3      uint32_t io_base;
4      uint32_t io_size;
5
6      struct cbl cbl;
7  };

```

`nic`是用于编写E100网卡驱动中所有过程中一个记录需要使用到的资源的工具, 管理我们用到的I/O端口和CBL、RFA等等。

注意在前面的初始化硬件过程中添加上`nic`的初始化过程:

```
kern/e100.c
1 struct nic nic;
2
3 int
4 e100_attach(struct pci_func *pcif)
5 {
6     pci_func_enable(pcif);
7     e100.bus = pcif->bus;
8     e100.dev_id = pcif->dev_id;
9     e100.dev_class = pcif->dev_class;
10    int i;
11    for (i = 0; i < 6; i++) {
12        e100.reg_base[i] = pcif->reg_base[i];
13        e100.reg_size[i] = pcif->reg_size[i];
14    }
15    e100.irq_line = pcif->irq_line;
16
17    // Initialize NIC
18    nic.io_base = pcif->reg_base[E100_IO];
19    nic.io_size = pcif->reg_size[E100_IO];
20
21    e100_init ();
22
23    return 0;
24 }
```

在初始化的第一步我们先需要为TCB分配**物理空间**，为什么要强调是物理地址呢？

因为一开始我犯了一个错误，我特别在系统**虚拟地址空间**里在KERNBASE以上找了一块连续的区域（我找的是KERNBASE + PGSIZE开始），专门映射这些TCB。但是后来就发现这样映射是完全没有必要的：

- 首先如果将其映射到虚拟地址上是为了方便内核里使用指针访问这些TCB的话，那么我们利用page\_alloc() 得到这些物理页面对应的struct Page 的时候，就可以利用page2kva (Page)的虚拟地址访问他们了
- 其次，前面我们使用page2kva (Page)访问的前提，是因为在内存管理的Lab中，KERNBASE以上的空间全部被**静态映射**成一一对应的物理内存（注意回顾静态映射的概念，就是被映射到的物理页**没有改变其引用数**，方便我们可以重用这些物理页，静态映射的意义就只是单纯提供一个顺序访问所有物理地址的转换而已）
- 如果我们设置一个KERNBASE上的区域专门映射我们分配到的TCB物理页面，那么该虚拟地址原来对应的真实物理页就要被卸载，如果在未来的某个时候这个真实物理页被分配后需要用page2kva这样的地址进行访问，那么映射就会失败，其访问到的是我们现在安装上的TCB的物理页，这样就会造成系统的崩溃

**所以为TCB的物理页分配虚拟地址空间不仅是没必要的，在有的时候更可能造成系统的崩溃。感谢张顺廷湿胸指出了我理解上的错误！**

真实分配的时候，我们为CB\_MAX\_NUM个TCB每个分配一个物理页（实际

上一个TCB大概只占用半页的样子，因为tcb\_data最大才1518个bytes而已，但是这样比较方便操作），于是初始化的程序如下：

```

kern/e100.c:  cbl_alloc()
1  /**
2   * Allocate CB_MAX_NUM pages, each page for a control block
3   */
4  static void
5  cbl_alloc () {
6      int i, r;
7      struct Page *p;
8      struct cb *prevcb = NULL;
9      struct cb *currcb = NULL;
10
11     // Allocate physical page for Control block
12     for (i = 0; i < CB_MAX_NUM; i++) {
13
14         if ((r = page_alloc (&p)) != 0)
15             panic ("cbl_init: _run_out_of_physical_memory!_%e\n", r);
16
17         p->pp_ref++;
18         memset (page2kva (p), 0, PGSIZE);
19
20         currcb = (struct cb *)page2kva (p);
21         currcb->phy_addr = page2pa (p);
22
23
24         if (i == 0)
25             nic.cbl.start = currcb;
26         else {
27             prevcb->cb_link = currcb->phy_addr;
28             prevcb->next = currcb;
29             currcb->prev = prevcb;
30         }
31
32         prevcb = currcb;
33     }
34
35     prevcb->cb_link = nic.cbl.start->phy_addr;
36     nic.cbl.start->prev = prevcb;
37     prevcb->next = nic.cbl.start;
38
39     nic.cbl.cb_avail = CB_MAX_NUM;
40     nic.cbl.cb_wait = 0;
41
42     nic.cbl.front = nic.cbl.start;
43     nic.cbl.rear = nic.cbl.start->prev;
44 }

```

通过page\_alloc() 拿到物理页之后，第一记得增加它的ref数，第二记得清空所有的数据。

TCB结构建立完毕以后，我们可以考虑使用CU在TCB上跑一跑了。从这里开始就涉及到对于CU的操作字的问题。首先回顾一下我们在2.3.1中提到的建立TCB以后可能需要进行的操作：

1. 将其中第一个TCB的物理地址写入CU的General Pointer，通知其要操作的TCB所在的内存位置

2. 给CU一个Start指令，CU开始工作
3. 如果环状的TCB中所有的S位都是0，那么CU将根据Link Address无限循环的处理这些TCB，每次执行完成以后，以中断和修改TCB中C状态位来表明任务已完成(后面我们会说明处理这些消息的方式，现在先不用管)
4. 如果CU碰到了某个TCB的S被设置为1，那么处理完该TCB之后，CU进入Suspend状态，等待用户发送Resume命令恢复运行，那么CU将从当前TCB的下个TCB开始执行

在这里首先回答一下关于CU**执行完成后的响应问题**。如果让CU使用中断的方式提醒我们，效率会比较高，但是比较麻烦的是要去处理中断响应。而如果使用轮询的方式的话，根据我们CBL的结构，只需要每次测试front头的C状态位是否完成即可知道任务的完成状态，非常方便。

我们带着这些需求来看看具体操作方法。控制CU是通过向SCB中写入相应的控制字决定的，具体规定在手册的6.3.2 System Control Block (SCB)中规定：

Table 12. System Control Block

31	16 15	0
Upper Word	Lower Word	Offset
SCB Command Word	SCB Status Word	Base + 00h
SCB General Pointer		Base + 04h

最主要的两个部分是控制字和状态字：

控制字：

Figure 10. SCB Command Word

31	26	25	24	23	20	19	18	16
Specific Interrupt Mask Bits			SI	M	CU Command	0	RU Command	

我们主要需要关注三个字段：

- M: 当控制字的这个位被设置成1，那么CU将不会发出任何中断，这个位正是当我们屏蔽中断时第一个需要给出的命令
- CUC (CU Command): 这里对CU的控制命令有以下几种：
  1. 0000 NOP
  2. 0001 CU Start
  3. 0010 CU Resume
  4. 0100 Load Dump Counters Address
  5. 0101 Dump Statistical Counters
  6. 0110 Load CU Base
  7. 0111 Dump and Reset Statistical Counters

## 8. 1010 CU Static Resume

我们要用到的主要是**CU Start**和**CU Resume**，他们对应的需求为：

- 在建立好TCB并且给CU设置好TCB的地址以后，发出一个CU Start指令，CU开始工作
- 当CU在某个TCB因为S位被挂起以后，发出一个CU Resume指令可以让其恢复工作

手册读到这里我们可以把相应的设置场位定义到kern/e100.h中了：

```
kern/e100.h

1 // CU Command Word
2 #define CUC_NOP      0x00
3 #define CUC_START    0x10
4 #define CUC_RESUME   0x20
5 #define CUC_LD_COUNTER 0x40
6 #define CUC_DUMP_SCNT 0x50
7 #define CUC_LOAD_BASE 0x60
8 #define CUC_DUMP_RSCNT 0x70
9 #define CUC_SRESUME  0xa0
```

- RUC (RU Command): 这个我们在后面3.1会详细讲述

看到这里，我们已经知道**如何关闭CU的中断**、使其**开始**和**恢复**执行，但是在开始前需要将TCB的地址告诉CU使其能够开始运行，这个如何设置？这里就涉及到SCB中General Pointer的设置了：

Table 15. SCB General Pointer for the CU Command

RUC Field	RU Command	SCB General Pointer	Added to
0	NOP	Don't care	
1	CU Start	Pointer to first command block in the command block list	CU Base
2	CU Resume	Don't care	
3	CU HPQ Start	Pointer to first command block in the HPQ command block list	CU Base
4	Load Dump Counters Address	Absolute address written to by Dump Counters and Dump & Reset Counters commands	
5	Dump Counters	Don't care	
6	Load CU Base	32-bit Base Register for CU data structures	
7	Dump & Reset Counters	Don't care	
10	CU Static Resume	Don't care	
11	CU HPQ Resume	Don't care	

General Pointer是SCB中根据不同的命令设置的一个场位，**提供某些命令执行时需要的数据**，比如这里我们只需要关注CU Start命令，在执行该命令前，General Pointer里必须写入开始执行的TCB的物理地址，那么CU Start时就可以从该TCB开始执行。而CU Resume就不需要，因为其下一次运行的TCB地址已经在被Suspend的时候被写入了内部寄存器。

设置General Pointer只需要像写入SCB的状态字和控制字一样直接像相应的I/O端口写入值即可。

Figure 9. SCB Status Word

15	8	7	6	5	2	1	0
STAT / ACK				CUS		RUS	
						0	0

状态字：

这里我们只需要关注CU的状态字即可，它主要可能有以下状态：

1. 00: Idle
2. 01: Suspend
3. 10: LPQ Active
4. 11: HQP Active

后面两个我都不知道是干吗的，这次实验里只需要知道前两个即可，因为如果不是前面两个的停止状态，那么当前CU肯定是在工作状态，然后将该状态值也定义到头文件里：

```

kern/e100.h
1 // CU Status Word
2 #define CUS_MASK      0xc0
3 #define CUS_IDLE      0x00
4 #define CUS_SUSPENDED 0x40
5 #define CUS_LPQ_ACTIVE 0x80
6 #define CUS_HQP_ACTIVE 0xc0

```

CUS\_MASK是由于要从SCB中读取其状态但CUS又是在中间的位置，所以先要用CUS\_MASK取出CUS相应场位出来

到这里我们已经将所有的控制指令都了解完毕了，在开始编写真正的驱动之前，还有一个东西我们需要明确，就是发出指令的状态控制，在手册的6.3.2.2 SCB Command Word中提到了这么一段话：

When software wants to issue an action command, it should write to the Command byte. The CUC and RUC fields of the Command byte specify the actions to be performed by the 8255x. The command is ready for acceptance by the device as soon as it is written into the CUC or RUC field. The actual command execution may not start instantaneously and will depend on current receive and transmit DMA activity. The Command byte is set by the CPU and cleared by the 8255x indicating command acceptance.

因为硬件控制的原因，所以在我们对于SCB写入相应的控制指令时，**并不会马上开始执行**，硬件要过一段时间以后才会接受，那么如何知道是否接受呢？上面的材料提到当CU接受命令以后，**SCB中的命令字会被硬件清除**。进一步的，在手册中的6.5 Starting and Completing Control Commands中也提到：



```
* Software must wait for this byte to be cleared before the next control
command can be issued.
* CU and RU control commands must never be issued together in the same SCB
write cycle.
```

通过修改SCB对CU发出一条指令之后，**我们必须等待其命令字被清空以后**，才能继续下面的指令。

好了我们可以正式开始编程了，首先是对CU发出命令的基本模块：

```
kern/e100.c: e100_exec_cmd()

1 static void
2 e100_exec_cmd (int csr_comp, uint8_t cmd)
3 {
4     int scb_command;
5
6     outb(nic.io_base + csr_comp, cmd);
7     do {
8         scb_command = inb(nic.io_base + CSR_COMMAND);
9     } while (scb_command != 0);
10 }
```

csr\_comp是SCB命令字的一个字段，cmd是需要执行的命令。发出命令后，我们通过轮询等待命令字被清空确认指令被接受。

初始化的第一步，屏蔽所有的中断：

```
kern/e100.c: e100_init()

1 static void
2 e100_init ()
3 {
4     // Software Reset E100
5     e100_sw_reset(e100);
6
7     // disable all interrupts
8     e100_exec_cmd (CSR_INT, 1);
9
10    cbl_init ();
11 }
```

我们把刚才在e100\_attach()中调用的软重启放到了e100\_init()中，并将e100\_attach()换成了调用e100\_init()完成E100的所有初始化。

然后我们看看对于CBL进行初始化的cbl\_init()：

```
kern/e100.c

1 static int
2 cbl_append_nop (uint16_t flag)
3 {
4     if (nic.cbl.cb_avail == 0)
5         return -E_CBL_FULL;
6
7     nic.cbl.cb_avail --;
8     nic.cbl.cb_wait ++;
9
10    nic.cbl.rear = nic.cbl.rear->next;
```

```

11     nic.cbl.rear->cb_status = 0;
12     nic.cbl.rear->cb_control = CBC_NOP | flag;
13
14     return 0;
15 }
16
17 static void
18 cbl_init ()
19 {
20     cbl_alloc ();
21
22     cbl_append_nop (0);
23     cbl_append_nop (0);
24     cbl_append_nop (0);
25     cbl_append_nop (CBF_S);
26
27     outl(nic.io_base + CSR_GP, nic.cbl.front->phy_addr);
28     e100_exec_cmd (CSR_COMMAND, CUC_START);
29 }
30

```

cbl\_append\_nop() 是在CBL待发送队列中添加一个NOP指令，用于我们会查看网卡输出判断我们在cbl\_alloc() 建立的结构是否被CU正确找到。

具体的话它的工作就是在待发送队列的末尾添加了一个NOP指令TCB。然后设置其flag为我们需要的状态，一般来说就是S=0或者S=1的区别。这里特别需要注意的是12行**对SCB的状态字清空**，因为前面提到过，CU通过设置状态字来提醒我们发送执行的状态。

同样这里使用到了出错状态，我一共定义了四种边界的错误：

kern/e100.h

```

1 // Error CODE
2 #define E_CBL_FULL 1
3 #define E_CBL_EMPTY 2
4 #define E_RFA_FULL 3
5 #define E_RFA_EMPTY 4

```

在cbl\_init() 第23行，添加完以后我们往SCB的General Pointer里写入了当前CBL里的第一个TCB的物理地址，然后发送了一条CU Start指令。这个逻辑产生的效果应该是，网卡启动后执行了4条NOP指令，然后在最后一条执行完后被挂起。

我们通过make qemu QEMUEXTRA="-debug-e100"启动JOS，其打印出的网卡记录为：

```

1 PCI: 00:00:0: 8086:1237: class: 6.0 (Bridge device) irq: 0
2 PCI: 00:01:0: 8086:7000: class: 6.1 (Bridge device) irq: 0
3 PCI: 00:01:1: 8086:7010: class: 1.1 (Storage controller) irq: 0
4 PCI: 00:01:3: 8086:7113: class: 6.80 (Bridge device) irq: 9
5 PCI: 00:02:0: 1013:00b8: class: 3.0 (Display controller) irq: 0
6 PCI: 00:03:0: 8086:1209: class: 2.0 (Network controller) irq: 11
7 EE100 pci_smio_map region 0, addr=0xf2020000, size=0x00001000, type=8
8 EE100 pci_map region 1, addr=0x0000c040, size=0x00000040, type=1
9 EE100 pci_smio_map region 2, addr=0xf2040000, size=0x00020000, type=0
10 PCI function 00:03:0 (8086:1209) enabled
11 EE100 eepro100_write4 addr=Port+0 val=0x00000000
12 EE100 nic_reset 0x92d8008
13 EE100 nic_selective_reset checksum=0xb34
14 EE100 eepro100_writel addr=Command/Status+3 val=0x01
15 EE100 eepro100_readl addr=Command/Status+2 val=0x00
16 EE100 eepro100_write_pointer val=0x040bb000
17 EE100 eepro100_writel addr=Command/Status+2 val=0x10
18 EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040ba000

```

```

19 EE100 action_command CU list with at least one more entry
20 EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b9000
21 EE100 action_command CU list with at least one more entry
22 EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b8000
23 EE100 action_command CU list with at least one more entry
24 EE100 action_command val=0x10 (cu start), status=0x0000, command=0x4000, link=0x040b7000
25 EE100 action_command CU list empty
26 EE100 eepr0100_readl addr=Command/Status+2 val=0x00
27 FS is running
28 FS can do I/O

```

从13行nic\_selective\_reset记录后开始，产生的log分别为：

1. 14-15行，向Command/Status+3写入0x01，即写入SCB Command字的M位，用于屏蔽中断，然后读取Command/Status+2即SCB Command字，查看是否被清空，这两句是我们在e100\_exec\_cmd()中一起执行的两条命令
2. 16行，向General Pointer写入TCB的物理地址
3. 17行，向Command/Status+2写入0x10，即发出CU Start指令
4. 18-25行，CU开始执行TCB中指明的指令了，到了第四个以后执行完毕停止
5. 26行，这句是对应17行中CU Start指令的验证，也是同在e100\_exec\_cmd()的

现在看起来还挺不错的，但是我们还没有验证以下几个方面：

- TCB的环状结构是否正常
- TCB的发送指令是否正常
- CU的Resume命令是否正常

所以我们将初始化过程修改成这样：

```

                                kern/e100.c
1 static int
2 cbl_append_transmit (const char *data, uint16_t l, uint16_t flag)
3 {
4     if (nic.cbl.cb_avail == 0)
5         return -E_CBL_FULL;
6
7     nic.cbl.cb_avail--;
8     nic.cbl.cb_wait++;
9
10    nic.cbl.rear = nic.cbl.rear->next;
11
12    nic.cbl.rear->cb_status = 0;
13    nic.cbl.rear->cb_control = CBC_TRANSMIT | flag;
14
15    nic.cbl.rear->cb_cmd_spec.tcb.tcb_tbd_array_addr = 0xFFFFFFFF;
16    nic.cbl.rear->cb_cmd_spec.tcb.tcb_byte_count = 1;
17    nic.cbl.rear->cb_cmd_spec.tcb.tcb_thrs = 0xE0;
18    nic.cbl.rear->cb_cmd_spec.tcb.tcb_tbd_count = 0;
19
20    memmove (nic.cbl.rear->cb_cmd_spec.tcb.tcb_data, (void *)data, l);
21
22    return 0;

```

```

23 }
24 static void
25 cbl_init ()
26 {
27     cbl_alloc ();
28
29     cbl_append_nop (0);
30     cbl_append_nop (0);
31     cbl_append_nop (0);
32     cbl_append_nop (CBF_S);
33     cbl_append_nop (0);
34     cbl_append_nop (0);
35     cbl_append_nop (0);
36     cbl_append_nop (0);
37     cbl_append_nop (0);
38
39     cbl_append_transmit ("aaaax", 5, 0);
40
41     outl(nic.io_base + CSR_GP, nic.cbl.front->phy_addr);
42     e100_exec_cmd (CSR_COMMAND, CUC_START);
43
44     e100_exec_cmd (CSR_COMMAND, CUC_RESUME);
45 }
46

```

这里添加了一个函数cbl.append.transmit() 添加一个发送指令的TCB到队尾，不再赘述。

然后初始化过程中我们添加发送了CBL中所有的10个TCB，并且最后一次发送transmit不是suspend，那么在我们发出的第一次Start指令后，应该会在第四个NOP指令停下来，接下来发出Resume指令后，应该会循环执行CBL中的所有TCB后在同样在第四个NOP停下来，启动JOS后打印的信息如下：

```

PCI function 00:03.0 (8086:1209) enabled
EE100 eepro100_write4 addr=Port+0 val=0x00000000
EE100 nic_reset 0x8a8e008
EE100 nic_selective_reset checksum=0xbe34
EE100 eepro100_writel addr=Command/Status+3 val=0x01
EE100 eepro100_readl addr=Command/Status+2 val=0x00
EE100 eepro100_write_pointer val=0x040bb000
EE100 eepro100_writel addr=Command/Status+2 val=0x10
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040ba000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b9000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b8000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x4000, link=0x040b7000
EE100 action_command CU list empty
EE100 eepro100_readl addr=Command/Status+2 val=0x00
EE100 eepro100_writel addr=Command/Status+2 val=0x20
EE100 eepro100_cu_command CU resuming
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b6000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b5000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b4000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b3000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0000, link=0x040b2000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0x0000, command=0x0004, link=0x040bb000
EE100 action_command transmit, TBD array address 0xffffffff, TCB byte count 0x0005, TBD count 0
EE100 action_command TBD (simplified mode): buffer address 0x040b2010, size 0x0005
EE100 action_command 0x8a8e008 sending frame, len=5, 61 61 61 61 78
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0xa000, command=0x0000, link=0x040ba000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0xa000, command=0x0000, link=0x040b9000
EE100 action_command CU list with at least one more entry
EE100 action_command val=0x10 (cu start), status=0xa000, command=0x0000, link=0x040b8000
EE100 action_command CU list with at least one more entry

```

```
EE100 action_command val=0x10 (cu start), status=0xa000, command=0x4000, link=0x040b7000
EE100 action_command CU list empty
EE100 eepro100_readl addr=Command/Status+2 val=0x00
FS is running
FS can do I/O
```

可以看到发送的流程和我们的预期是正常的，但是这里看到发送的时候打印的信息不是很给力：

```
EE100 action_command transmit, TBD array address 0xffffffff, TCB byte count 0x0005, TBD count 0
EE100 action_command TBD (simplified mode): buffer address 0x040b2010, size 0x0005
EE100 action_command 0x8a8e008 sending frame, len=5, 61 61 61 61 78
EE100 action_command CU list with at least one more entry
```

尤其是在发送大量数据的话在大量包记录里很难找到有效信息，我们可以使用提供的包拦截参数将其记录到文件，使用`make qemu QEMUEXTRA="-debug-e100-pcap slirp.cap"`打印调试信息的同时将网卡的包抓取后放入slirp.cap文件中，这个文件是一个二进制文件，人无法直接读取，需要tcpdump来为我们解析，使用`tcpdump -XXr slirp.cap`命令打印其中的内容，得到的输出如下：

```
zhangchi@zhangchi-vostro1400:~/lab$ tcpdump -XXr slirp.cap
reading from file slirp.cap, link-type EN10MB (Ethernet)
15:05:39.002011 [|ether]
0x0000: 6161 6161 78 aaaax
zhangchi@zhangchi-vostro1400:~/lab$
```

可以看到我们发送的aaaax消息，这样看着就直观多了。

最后正常的初始化过程应该是这样：

```
kern/e100.c: cbl_init()

1 static void
2 cbl_init ()
3 {
4     cbl_alloc ();
5
6     cbl_append_nop (CBF_S);
7
8     outl(nic.io_base + CSR_GP, nic.cbl.front->phy_addr);
9     e100_exec_cmd (CSR_COMMAND, CUC_START);
10 }
```

执行一个NOP指令后停在第一个TCB上，等待Resume

到这里我们已经做完了所有的实验，到现在已经可以实现出一个完整的系统调用了：

**Exercise 6.** Create a system call for transmitting packets. The interface is up to you. As described in the Device Driver Organization section the send system call should add the packet to the transmit DMA ring and restart or resume the CU if it is idle or suspended. If your design requires it, you can take this opportunity to reclaim any buffers which have been marked as transmitted by the E100 in order to free up space in the transmit DMA ring.

我设计的接口如下：

```

                                kern/e100.c
1  static void
2  cbl_validate ()
3  {
4      while (nic.cbl.cb_wait > 0 && (nic.cbl.front->cb_status & CBS_C) != 0) {
5          nic.cbl.front = nic.cbl.front->next;
6          nic.cbl.cb_avail ++;
7          nic.cbl.cb_wait --;
8      }
9  }
10
11 int
12 e100_transmit (const char *data, uint16_t len)
13 {
14     cbl_validate ();
15
16     if (nic.cbl.cb_avail == 0)
17         return -E_CBL_FULL;
18
19     nic.cbl.rear->cb_control &= ~CBF_S;
20     cbl_append_transmit (data, len, CBF_S);
21
22     int scb_status = inb(nic.io_base + CSR_STATUS);
23     if ((scb_status & CUS_MASK) == CUS_SUSPENDED)
24         e100_exec_cmd (CSR_COMMAND, CUC_RESUME);
25
26     return 0;
27 }

```

**cbl\_validate() :**

将在待发送队列中已经发送完毕的TCB回收，并维护CBL相关的状态值

**e100\_transmit() :**

- 发送前先进行回收，如果回收后TCB仍然没有空闲，那么则返回错误信息，交给用户程序重新发送。
- 否则将当前CBL的最后一个的控制位中的S清0，并且将要发送的数据以TRANSMIT命令形式添加到队尾，并设置S位为1，即指示CU每次在最后一个待发送包发送完后挂起
- 恢复CU运行

然后根据e100\_transmit() 我向系统添加了一条系统调用：

```

                                kern/syscall.c: sys_net_try_send()
1  static int
2  sys_net_try_send (const char *data, uint16_t len)
3  {
4      return e100_transmit (data, len);
5  }

```

将命令注册为系统调用需要修改很多文件，有：

- kern/syscall.c, kern/syscall.h

- lib/syscall.c, lib/lib.h
- inc/syscall.h

然后我写了一个用户程序调用网络发送字符串：

```

                                user/testsend.c
1  #include <inc/x86.h>
2  #include <inc/lib.h>
3
4  void
5  umain(void)
6  {
7      sys_net_try_send ("Hello, World", 12);
8  }

```

将用户程序加入内核镜像记得修改kern/Makefrag，然后就可以使用make run-testsend QEMUEXTRA="-debug-e100 -pcap slirp.cap" 运行该程序了，运行完成后使用tcpdump分析包，得到预期的结果：

```

zhangchi@zhangchi-vostro1400:~/lab$ tcpdump -XXr slirp.cap
reading from file slirp.cap, link-type EN10MB (Ethernet)
23:25:17.333149 [|ether]
                0x0000:  4865  6c6c  6f72  2057  6f72  6c64           Hello, World
zhangchi@zhangchi-vostro1400:~/lab$

```

## 2.4 Transmitting Packets: Network Server

这里我们就可以实现前面2.0.2中的Output Helper Environment了，写的时候我对于IPC的细节又忘的差不多了，可以参考Lab5中文件服务器的代码fs/serv.c的serve() 很有帮助，这里就直接贴代码了：

```

                                net/output.c
1  #include "ns.h"
2
3  extern union Nsipc nsipcbuf;
4
5  void
6  output(envid_t ns_envid)
7  {
8      binaryname = "ns_output";
9
10     uint32_t req, whom;
11     int perm, r;
12
13     while (1) {
14         perm = 0;
15         req = ipc_recv((int32_t *) &whom, &nsipcbuf, &perm);
16
17         // All requests must contain an argument page
18         if (!(perm & PTE_P))
19             panic ("output: Invalid request from %08x: no argument page\n", whom);
20
21         if (req != NSREQ_OUTPUT)
22             panic ("output: Invalid IPC Request type: %d\n", req);
23     }

```

```

24 while ((r = sys_net_try_send (nsipcbuf.pkt.jp_data, nsipcbuf.pkt.jp_len))
25         != 0);
26     sys_page_unmap(0, &nsipcbuf);
27 }
28 }

```

启动JOS后可以看到网络服务器发送的包：

```

ns: 52:54:00:12:34:56 bound to static IP 10.0.2.15
NS: TCP/IP initialized.
EE100 eepro100_writel      addr=Command/Status+2 val=0x20
EE100 eepro100_cu_command  CU resuming
EE100 action_command      val=0x10 (cu start), status=0x0000, command=0x4004, link=0x040bb000
EE100 action_command      transmit, TBD array address 0xffffffff, TCB byte count 0x002a, TBD count 0
EE100 action_command      TBD (simplified mode), buffer address 0x040bb010, size 0x002a
EE100 action_command      0xa052008 sending frame, len=42, ff ff ff ff ff 52 54 00 12 34 56 08 06 00 01
EE100 nic_receive         0xa052008 received frame for me, len=42
EE100 nic_receive         no resources, state=0
EE100 action_command      CU list empty
EE100 eepro100_readl      addr=Command/Status+2 val=0x00
file flush is good
file truncate is good
file rewrite is good

```

测试make grade可以顺利通过output测试。

## 3 Receiving packets and the web server

### 3.1 Receiving Packets

在E100处理数据接收的时候，和前面处理发送几乎是一样的，只是RU处理的是RFD而CU处理的是CB，而且CU可能有多个命令，每个命令对应的CB结构是不一样的，发送对应的CB就是TCB。

而RU只有RFD一种结构，其布局如下：

Figure 25. Receive Frame Descriptor Format

Offset	Command Word Bits 31:16						Status Word Bits 15:0			
00h	EL	S	00000000	H	SF	000	C	0	OK	Status Bits
04h	Link Address (A31:A0)									
08h	Reserved									
0Ch	0	0	Size				EOF	F	Actual Count	

这里控制字和状态字中需要我们注意的只有S和C两位，没有CMD。

- 如果S被设置成1，那么RU在接受完这个RFD的数据以后，挂起
- 如果C被设置成1，那么代表接收完成，注意**接收之前用户应该完成这个位的清除工作**

前面的都和TCB是差不多的，但是下面的有两个场位是不同的：



- Size: 表示data buffer的大小，在我们的数据中，数据包大小最大为1518，这个值在接收数据之前要由用户来设置好。
- Actual Count: 如果C被置1表示接收完成，那么Actual Count表示的就是data buffer中总共的数据bytes数。这个是由RU接收完以后填好

根据这些字段说明，我们把他们定义到头文件里方便后面使用和设置：

```
kern/e100.h

1 // Recieve Frame Descriptor Command
2 #define RFDF_EL      0x8000
3 #define RFDF_S       0x4000
4 #define RFDF_H       0x10
5 #define RFDF_SF      0x8
6
7
8 // Recieve Frame Descriptor Status
9 #define RFDS_C       0x8000
10 #define RFDS_OK      0x2000
11 #define RFDS_MASK    0x1fff
12
13
14 // Recieve Frame Descriptor Data
15 #define RFD_SIZE_MASK 0x3fff
16 #define RFD_AC_MASK  0x3fff
17 #define RFD_EOF      0x8000
18 #define RFD_F        0x4000
```

了解完RFD的结构，类似的，我们需要在其上面实现一个支持多数据包等待接收的排队系统，同样也要注意两方面

1. 首先要记录的是当前RFA环状结构中哪些RFD是已经被RU收取完毕，正等待被用户接收的，哪些是可以被RU重新利用来接受新数据的
2. 其次要做好相应的S位设置，让RU在接收到不能接收的时候就停下来，不要让它将还没被用户接收的RFD覆盖了

```
kern/e100.h

1 #define RFD_MAXSIZE  1518
2 #define RFD_MAX_NUM   10
3
4 // Receive Frame Descriptor
5 struct rfd {
6     volatile uint16_t rfd_status;
7     uint16_t rfd_control;
8     uint32_t rfd_link;
9
10    uint32_t rfd_reserved;
11    uint16_t rfd_actual_count;
12    uint16_t rfd_size;
13
14    char rfd_data[RFD_MAXSIZE];
15
16    struct rfd *prev, *next;
17    physaddr_t phy_addr;
18 };
19
```

```
20 // Receive Frame Area
21 struct rfa {
22     int rfd_avail;
23     int rfd_wait;
24
25     struct rfd *start;
26     struct rfd *front, *rear;
27 };
```

struct rfd的结构没什么好说的，按照手册的规定设计即可，这里关键的是struct rfa，虽然和struct cbl的成员结构几乎一模一样，但是他们两者成员的意义是不一样的。

- rfd\_avail: 表示当前有多少个闲置的RFD可以用来**接收RU新收到的数据**
- rfd\_wait: 表示当前有多少个RFD正在处于**等待被用户接收**的状态
- 很明显上面两者相加应该等于所有RFD的总数，在这里应该是我们定义的RFD\_MAX\_NUM = 10
- start: 所有RFD中的第一个RFD，用于开始的时候初始化RU用
- front和rear: 表示当前**正在等待被用户接收**的RFD的起始和结尾

那么这个系统是如何根据rfa中定义的结构工作的呢？

1. 初始的时候front=start, rear=start→prev，表示当前等待被用户接收的RFD队列为空，同时将start→prev的Suspend控制位设成1，有两个目的：
  - 在第一次让RU开始执行START指令时，让RU可以一直接收数据直到最后一个RFD才停下。
  - 一旦RU开始运行，保证任何时候队列里只有front→prev的S位为1，即RU一直接受数据**直到front之前就要停止**

并且rfd\_avail = RFD\_MAX\_NUM, rfd\_wait = 0

2. 当用户需要收取一个数据包时，将front的内容取出，作为数据返回。并且清空front→prev的S位，将front的S位置为1，表示该位置**可以让RU重新回收使用了**。做完以后front向后移动一个位置。同时增加rfd\_avail, 减少rfd\_wait
3. 当确认数据包被RU接收完毕时，检查rear的后一个指向TCB的C状态是否为1，可以的话则将rear向后移动，同时减少rfd\_avail, 增加rfd\_wait

好了RFA的结构了解完毕，现在我们可以考虑建立起这样的结构了：

**Exercise 8.** Construct a receive DMA ring and start the RU. If you use interrupts, make sure that 82559ER-generated interrupts are routed to your driver and are handled.

首先将RFA加入nic结构：

```
kern/e100.h
1 // Network Interface Card
2 struct nic {
3     uint32_t io_base;
4     uint32_t io_size;
5
6     struct cbl cbl;
7     struct rfa rfa;
8 };
```

然后为RFA分配物理空间：

```
kern/e100.c: rfa_alloc()
1 static void
2 rfa_alloc () {
3     int i, r;
4     struct Page *p;
5     struct rfd *prevrfd = NULL;
6     struct rfd *currrfd = NULL;
7
8     // Allocate physical page for Control block
9     for (i = 0; i < RFD_MAX_NUM; i++) {
10         if ((r = page_alloc (&p)) != 0)
11             panic ("rfa_init: _run_out_of_physical_memory!_%e\n", r);
12
13         p -> pp_ref ++;
14         memset (page2kva (p), 0, PGSIZE);
15
16         currrfd = (struct rfd *)page2kva (p);
17         currrfd->phy_addr = page2pa (p);
18         currrfd->rfd_control = 0;
19         currrfd->rfd_status = 0;
20         currrfd->rfd_size = RFD_MAXSIZE;
21
22         if (i == 0)
23             nic.rfa.start = currrfd;
24         else {
25             prevrfd->rfd_link = currrfd->phy_addr;
26             prevrfd->next = currrfd;
27             currrfd->prev = prevrfd;
28         }
29
30         prevrfd = currrfd;
31     }
32
33     prevrfd->rfd_link = nic.rfa.start->phy_addr;
34     nic.rfa.start->prev = prevrfd;
35     prevrfd->next = nic.rfa.start;
36
37     nic.rfa.rfd_avail = RFD_MAX_NUM;
38     nic.rfa.rfd_wait = 0;
39
40     nic.rfa.front = nic.rfa.start;
41     nic.rfa.rear = nic.rfa.start->prev;
42     nic.rfa.rear->rfd_control |= RFD_S;
43 }
```

注意上面程序和2.3.1的区别，因为RU一启动就开始等待接收数据，所以我们一开始就要设置好RFD的命令字、状态字和data buf大小等等，请关注18行到20行。最后将RFA中的末尾RFD的S位置成1。

然后是对RU的控制，前面在2.3.1介绍SCB时只介绍了针对CU的控制字和状态字，现在我们把RU的补充上来：

控制字：

Figure 10. SCB Command Word

31	26	25	24	23	20	19	18	16
Specific Interrupt Mask Bits			SI	M	CU Command		0	RU Command

RUC (RU Command): 这里对RU的控制命令有以下几种：

1. 000 NOP
2. 001 RU Start
3. 010 RU Resume
4. 011 Receive DMA Redirect
5. 100 RU Abort
6. 101 Load Header Data Size (HDS)
7. 110 Load RU Base

我们要用到的主要是**RU Start**和**RU Resume**，把相应的设置场位定义到kern/e100.h中：

```

kern/e100.h
1 // RU Command Word
2 #define RUC_NOP      0x0
3 #define RUC_START    0x1
4 #define RUC_RESUME   0x2
5 #define RUC_REDIR    0x3
6 #define RUC_ABORT    0x4
7 #define RUC_LOADHDS  0x5
8 #define RUC_LOAD_BASE 0x6

```

状态字：

Figure 9. SCB Status Word

15	8	7	6	5	2	1	0
STAT / ACK			CUS	RUS		0	0

RU的状态字主要可能有以下状态：

1. 0000: Idle
2. 0001: Suspend

3. 0010: No resources

4. 0011: Reserved

我们只需要知道Suspend用于给出恢复命令时判断RU的当前状态，相关头文件定义：

```
kern/e100.h  
1 // RU Status Word  
2 #define RUS_MASK      0x3c  
3 #define RUS_IDLE      0x0  
4 #define RUS_SUSPEND   0x4  
5 #define RUS_NORES     0x8  
6 #define RUS_READY     0x10
```

下面尝试初始化RFA。首先是两种对于RFA的操作：

- 第一种是确认其所在RFD的数据被RU接收完成
- 第二种是从队头取出一个数据准备返回给用户

这两种操作的情况我们在前面1已经说过了。后面给出的程序中第一种操作对应rfa\_validate()，第二种操作对应rfa\_retrieve\_data()。

为了保证我们后面编写系统调用的正确性，我们在初始化的时候就尽量对前面的RFA初始化以及指令操控作了一系列的测试，保证不要将错误在外层调用它时才显现出来，以免陷入混乱的调试。所以我们在程序中打印了一些消息用于查看：

```
kern/e100.c  
1 static void  
2 rfa_validate ()  
3 {  
4     while (nic.rfa.rfd_avail > 0 && (nic.rfa.rear->next->rfd_status & RFDS_C) !=  
5         0) {  
6         nic.rfa.rear = nic.rfa.rear->next;  
7  
8         nic.rfa.rfd_avail --;  
9         nic.rfa.rfd_wait ++;  
10        cprintf ("zhangchi: validate, _avail_=%d, _wait_=%d, _slot_=%x\n",  
11            nic.rfa.rfd_avail, nic.rfa.rfd_wait, nic.rfa.rear);  
12    }  
13 }  
14 static int  
15 rfa_retrieve_data (char* data)  
16 {  
17     if (nic.rfa.rfd_wait == 0)  
18         return -E_RFA_EMPTY;  
19  
20     nic.rfa.rfd_avail ++;  
21     nic.rfa.rfd_wait --;
```

```

22     cprintf ("zhangchi:_retrieve,_avail_=%d,_wait_=%d,_slot_=%x\n",
23             nic.rfa.rfd_avail, nic.rfa.rfd_wait, nic.rfa.front);
24
25     nic.rfa.front->prev->rfd_control &= ~RDFD_S;
26     nic.rfa.front->rfd_control = RDFD_S;
27     nic.rfa.front->rfd_status = 0;
28
29     int r = nic.rfa.front->rfd_actual_count & RFD_AC_MASK;
30     memmove (data, nic.rfa.front->rfd_data, r);
31
32     nic.rfa.front = nic.rfa.front->next;
33
34     return r;
35 }

```

其中rfa\_retrieve\_data() 接受一个buffer作为参数, 如果当前RFA为空, 那么返回一个负数错误码, 否则将数据填入buffer后返回传输的数据字节数。

然后主过程是这么写的:

```

                                kern/e100.c:  rfa_init()
1  static void
2  rfa_init ()
3  {
4      cprintf ("\n\nRFA_Initialization_started!\n");
5      rfa_alloc ();
6
7      outl(nic.io_base + CSR_GP, nic.rfa.front->phy_addr);
8      e100_exec_cmd (CSR_COMMAND, RUC_START);
9
10     while (nic.rfa.rfd_avail > 0)
11         rfa_validate ();
12
13     int scb_status = inb(nic.io_base + CSR_STATUS);
14     cprintf ("zhangchi:_rfd_slot_is_full,_current_RU_state_=%02x\n",
15             scb_status & RUS_MASK);
16
17     char s[1518];
18     while (rfa_retrieve_data (s) >= 0);
19
20     e100_exec_cmd (CSR_COMMAND, RUC_RESUME);
21
22     while (nic.rfa.rfd_avail > 0)
23         rfa_validate ();
24 }

```

可以看到在为RFA分配物理内存以后, 马上会进行一系列的测试:

1. 第7行, 给RU发送START命令, 使其开始接收数据包, 因为初始化的时候只有最后一个RFD设置了S标志, 所以开始以后只有当RU将所有RFD都接收了数据才会停下来
2. 第10行, 不断检查系统中可用RFD的个数, 当其减少到0的时候, 就证明整个RFA已经被RU收到的数据塞满了。在rfa\_validate() 检查过程中一旦发现有新收到的数据, 将会打印出一条validate信息
3. 第18行, 被塞满以后数据应该被发送到用户, 这里我们只是单纯的把数据提取出来, 不作处理。rfa\_retrieve.data() 执行过程中也会打印出一条信息。

## 4. 第20行，开始让RU继续接收数据

测试的时候我们使用 `make qemu QEMUEXTRA="-debug-e100-pcap slirp.cap"` 启动JOS，那么初始JOS应该会停滞在第10行轮询检查可用RFD的个数，那么这个时候如果开始给它发消息，就会有相应信息打印出来了，新开一个terminal使用 `make nc-7` 命令即可向网卡发送数据包，下面是我们得到的输出（为了打印长度的考虑，我把RFD\_MAX\_NUM改成了4，使我们能更快看到程序执行到边界条件）。

```
RFA Initialization started!
EE100 eepro100_write_pointer val=0x040b1000
EE100 eepro100_writel addr=Command/Status+2 val=0x01
EE100 eepro100_rx_command val=0x01 (rx start)
EE100 eepro100_readl addr=Command/Status+2 val=0x00
EE100 nic_can_receive 0x8adc008
EE100 nic_receive 0x8adc008 received broadcast, len=42
EE100 nic_receive command 0x0000, link 0x040b0000, addr 0x00000000, size 1518
zhangchi: validate, avail = 3, wait = 1, slot = f40b1000
EE100 nic_can_receive 0x8adc008
EE100 nic_receive 0x8adc008 received broadcast, len=42
EE100 nic_receive command 0x0000, link 0x040af000, addr 0x00000000, size 1518
zhangchi: validate, avail = 2, wait = 2, slot = f40b0000
EE100 nic_can_receive 0x8adc008
EE100 nic_receive 0x8adc008 received broadcast, len=42
EE100 nic_receive command 0x0000, link 0x040ae000, addr 0x00000000, size 1518
zhangchi: validate, avail = 1, wait = 3, slot = f40af000
EE100 nic_can_receive 0x8adc008
EE100 nic_receive 0x8adc008 received broadcast, len=42
EE100 nic_receive command 0x4000, link 0x040b1000, addr 0x00000000, size 1518
zhangchi: validate, avail = 0, wait = 4, slot = f40ae000
EE100 eepro100_readl addr=Command/Status+0 val=0x44
zhangchi: rfd slot is full, current RU state = 04
zhangchi: retrieve, avail = 1, wait = 3, slot = f40b1000
zhangchi: retrieve, avail = 2, wait = 2, slot = f40b0000
zhangchi: retrieve, avail = 3, wait = 1, slot = f40af000
zhangchi: retrieve, avail = 4, wait = 0, slot = f40ae000
EE100 eepro100_writel addr=Command/Status+2 val=0x02
EE100 eepro100_readl addr=Command/Status+2 val=0x00
EE100 nic_can_receive 0x8adc008
EE100 nic_receive 0x8adc008 received broadcast, len=42
EE100 nic_receive command 0x0000, link 0x040b0000, addr 0x00000000, size 1518
zhangchi: validate, avail = 3, wait = 1, slot = f40b1000
EE100 nic_can_receive 0x8adc008
EE100 nic_receive 0x8adc008 received broadcast, len=42
EE100 nic_receive command 0x0000, link 0x040af000, addr 0x00000000, size 1518
zhangchi: validate, avail = 2, wait = 2, slot = f40b0000
EE100 nic_can_receive 0x8adc008
EE100 nic_receive 0x8adc008 received broadcast, len=42
EE100 nic_receive command 0x0000, link 0x040ae000, addr 0x00000000, size 1518
zhangchi: validate, avail = 1, wait = 3, slot = f40af000
```

可以看到输出是符合我们的预期的，通过slot可以看到每次操作的RFD编号，可以观察RU操作这些RFD的顺序。测试成功以后，我们需要将初始化改回来：

```
kern/e100.c: rfa_init()

1 static void
2 rfa_init ()
3 {
4     rfa_alloc ();
5
6     outl(nic.io_base + CSR_GP, nic.rfa.front->phy_addr);
7     e100_exec_cmd (CSR_COMMAND, RUC_START);
8 }
```

实验很顺利，可以考虑编写给系统调用的接口了：

**Exercise 9.** Create a system call for receiving packets. As described in the Device Driver Organization section, the system call will read a packet out of the receive DMA ring, mark the DMA buffer as empty (so that the E100 can reuse it), resume the RU if necessary, and pass the packet to the calling user environment.

我设计的接口如下：

```

kern/e100.c: e100_receive()
1 int
2 e100_receive (char *data)
3 {
4     rfa_validate ();
5
6     if (nic.rfa.rfd_wait == 0)
7         return -E_RFA_EMPTY;
8
9     int r = rfa_retrieve_data (data);
10
11     int scb_status = inb(nic.io_base + CSR_STATUS);
12     if ((scb_status & RUS_MASK) == RUS_SUSPEND)
13         e100_exec_cmd (CSR_COMMAND, RUC_RESUME);
14
15     return r;
16 }

```

如果没有数据可以获取，返回负数错误，否则将数据填入data代表的buffer，然后将数据的字节数作为返回值返回。

记得修改相关代码将系统调用安装到JOS中去。我添加的系统调用如下：

```

kern/syscall.c: sys_net_try_recv()
1 static int
2 sys_net_try_recv (char *data)
3 {
4     return e100_receive (data);
5 }

```

数据接口以一个数据buffer为参数，然后以接受数据的长度为返回值。

## 3.2 Receiving Packets: Network Server

**Exercise 10.** Implement net/input.c.

这个Exercise是我本次实验花费的时间最多的地方之一，为什么，可以看看我的代码：



net/input.c

```

1 #include "ns.h"
2 #include <inc/lib.h>
3
4 extern union Nsipc nsipcbuf;
5
6 void
7 input(envid_t ns_envid)
8 {
9     binaryname = "ns_input";
10
11     while (1) {
12         while ((nsipcbuf.pkt.jp_len
13                 = sys_net_try_recv (nsipcbuf.pkt.jp_data)) < 0);
14
15         ipc_send(ns_envid, NSREQ_INPUT, &nsipcbuf, PTE_U|PTE_W|PTE_P);
16     }
17 }

```

这段代码在测试的时候发现make grade在testinput死都不过，于是我尝试自己单独运行testinput，在grade-lab6.sh中可以看到testinput是这么调用的：

```

runtest1 -tag "testinput_[5_packets]" -dir net testinput -DTEST_NO_NS \
-check check_testinput 5

```

注意DTEST\_NO\_NS是不启动网络服务器的宏，这就让我很纳闷了，Input和Output都是网络服务器的一个子模块，如果不启动的话怎么测试Input和Output模块的？结果查看了一下net/testinput.c代码后才发现，它自己在进程中创建了Input环境和Output环境，单独对这两个模块进行测试：

net/testinput.c: umain()

```

1 void
2 umain(void)
3 {
4     environ_t ns_envid = sys_getenvid();
5     int i, r;
6
7     binaryname = "testinput";
8
9     output_envid = fork();
10    if (output_envid < 0)
11        panic("error_forking");
12    else if (output_envid == 0) {
13        output(ns_envid);
14        return;
15    }
16
17    input_envid = fork();
18    if (input_envid < 0)
19        panic("error_forking");
20    else if (input_envid == 0) {
21        input(ns_envid);
22        return;
23    }
24
25    cprintf("Sending_ARP_announcement...\n");
26    announce();
27
28    cprintf("Waiting_for_packets...\n");
29    while (1) {
30        envid_t whom;

```

```

31         int perm;
32
33         int32_t req = ipc_recv((int32_t *)&whom, pkt, &perm);
34         if (req < 0)
35             panic("ipc_recv: %e", req);
36         if (whom != input_envid)
37             panic("IPC_from_unexpected_environment: %08x", whom);
38         if (req != NSREQ_INPUT)
39             panic("Unexpected_IPC: %d", req);
40
41         hexdump("input: ", pkt->jp_data, pkt->jp_len);
42         cprintf("\n");
43     }
44 }

```

所以单独运行testinput就只需要启动testinput进程就可以了，NS和FS都可以不用启动，反正也没有用到，即把kern/init.c的创建进程一段改成：

```

kern/init.c: i386_init()
1
2     // Should always have an idle process as first one.
3     ENV_CREATE(user_idle);
4
5     // Start fs.
6     //ENV_CREATE(fs_fs);
7
8 #if !defined(TEST_NO_NS)
9     //Start ns.
10    //ENV_CREATE(net_ns);
11 #endif
12
13 #if defined(TEST)
14     // Don't touch -- used by grading script!
15     ENV_CREATE2(TEST, TESTSIZE);
16 #else
17     // Touch all you want.
18     ENV_CREATE(net_testinput);
19     // ENV_CREATE(user_echosrv);
20     // ENV_CREATE(user_httpd);
21 #endif // TEST*
22
23     // Schedule and run the first user environment!
24
25     sched_yield ();

```

然后make qemu发现出错了，内核地址错误：

```

Sending ARP announcement...
Waiting for packets...
TRAP frame at 0xefbffecc
edi 0x00806004
esi 0xf40b1010
ebp 0xefbfff10
oesp 0xefbffeec
ebx 0x0000002a
edx 0x0000000a
ecx 0x0000002a
eax 0x00806004
es 0x----0010
ds 0x----0010
trap 0x0000000e Page Fault
err 0x00000003
eip 0xf0105b14

```

```

cs  0x---0008
flag 0x00003002
esp 0xf0286208
ss  0x---6004
kernel panic at kern/trap.c:313: kernel-mode page faults
Welcome to the JOS kernel monitor!
Type 'help' for a list of commands.

```

很奇怪，然后选择打印栈轨迹查看是从哪里出错的：

```

Welcome to the JOS kernel monitor!
Type 'help' for a list of commands.
K> backtrace
Stack backtrace:
ebp efbffd0 eip f010092a args 00000001 efbffd8 00000000 efbffe64 f0285d60
    kern/monitor.c:419: monitor+274
ebp efbffe30 eip f01000e2 args 00000000 efbffe64 00000139 00000008 efbffecb
    kern/init.c:112: _panic+93
ebp efbffe50 eip f01041f8 args f01086c0 00000139 f01086a8 00000008 00000000
    kern/trap.c:356: page_fault_handler+63
ebp efbffe90 eip f01043d5 args efbffecb 00806000 efbffed0 f0103992 efbffef4
    kern/trap.c:232: trap+198
ebp efbffec0 eip f0104582 args efbffecb 00806004 f40b1010 efbfff10 efbffecb
    kern/trapentry.S:103: <unknown>+0
ebp efbfff10 eip f0106205 args 00806004 f40b1010 0000002a 007ff000 00000000
    kern/e100.c:328: e100_receive+196
ebp efbfff30 eip f0104c1c args 00806004 f02b8564 00807000 00000007 f02b9174
    kern/syscall.c:593: syscall+1535
ebp efbfff80 eip f0104416 args 00000010 00806004 00000000 00000000 00000000
    kern/trap.c:248: trap+263
ebp efbfffb0 eip f0104582 args efbfffb0 00000000 00000000 eebfdef0 efbfffdc
    kern/trapentry.S:103: <unknown>+0
ebp eebfdef0 eip 0080050a args 00806004 00000000 00000000 00001001 00001001
    <unknown>:0: <unknown>+0
ebp eebfdff0 eip 008000ce args 00001001 00000000 00000000 00000000 00000000
    <unknown>:0: <unknown>+0
ebp eebfdff0 eip 00800614 args 00000000 00000000 00000000 00000000 00000000
    <unknown>:0: <unknown>+0

```

看到其中出现了kern/e100.c:328 e100\_receive的信息，然后去源文件找到这行代码：

```

                                kern/e100.c: e100_receive()
323     nic.rfa.front->prev->rfd_control &= ~RFDF_S;
324     nic.rfa.front->rfd_control = RFDF_S;
325     nic.rfa.front->rfd_status = 0;
326
327     int r = nic.rfa.front->rfd_actual_count & RFD_AC_MASK;
328     memmove (data, nic.rfa.front->rfd_data, r);
329
330     nic.rfa.front = nic.rfa.front->next;
331
332     return r;

```

328行是memmove，也就是说往data里写入数据的时候出错了，data是input系统调用的传入参数，最外层的时候应该是由net/input.c中调用的：

```

                                net/input.c
1  #include "ns.h"
2  #include <inc/lib.h>

```

```

3
4 extern union Nsipc nsipcbuf;
5
6 void
7 input(envid_t ns_envid)
8 {
9     binaryname = "ns_input";
10
11     while (1) {
12         while ((nsipcbuf.pkt.jp_len
13                 = sys_net_try_rcv (nsipcbuf.pkt.jp_data)) < 0);
14
15         ipc_send(ns_envid, NSREQ_INPUT, &nsipcbuf, PTE_U|PTE_W|PTE_P);
16     }
17 }

```

第12行的`sys_net_try_rcv` 就是调用的`e100_receive()`，可以看到写入的目标是`nsipcbuf.pkt.jp_data`，为什么写入`nsipcbuf.pkt.jp_data`会错误呢？我百思不得其解，直到张顺廷湿胸提醒我才明白过来，Input和Output都是由`testinput`这个进程fork出来的，因此对于`nsipcbuf`这样的数据，在他们的地址中都是COW页。

如果按照我这样的调用方法传入调用系统调用，那么向`nsipcbuf`中写入值的时候是在`kern/e100.c`中的`e100_receive()`，运行过程中是处在内核态，内核态没有COW恢复机制，所以就发生了内核页错误了。COW恢复只可能是在用户态程序中进行写入时才有可能发生，即只有在Input本身向`nsipcbuf`写入数据，才有可能恢复。

所以处理的办法是，在Input环境中开一片buf专门用于存储从系统调用`sys_net_try_rcv()`接受来的数据，接受完成后，再将buf的内容使用`memmove`复制给`nsipcbuf.pkt.jp_data`，修改后的Input如下：

```

                                net/input.c
1 #include "ns.h"
2 #include <inc/lib.h>
3
4 extern union Nsipc nsipcbuf;
5
6 void
7 input(envid_t ns_envid)
8 {
9     binaryname = "ns_input";
10    char buf[1518];
11    int len;
12
13    while (1) {
14        while ((len = sys_net_try_rcv (buf)) < 0);
15
16        nsipcbuf.pkt.jp_len = len;
17        memmove ((void*) nsipcbuf.pkt.jp_data, (void *) buf, len);
18
19        ipc_send(ns_envid, NSREQ_INPUT, &nsipcbuf, PTE_U|PTE_W|PTE_P);
20    }
21 }

```

修改以后再次测试，出现了很奇怪的现象：

```
testinput [5 packets]: WRONG, receiving packet 001/5 (5.0s)
expected input: 0030 2030 3031 0a
got      input: 0030 2030 3032 0a
testinput [100 packets]: WRONG, receiving packet 001/100 (5.2s)
expected input: 0030 2030 3031 0a
got      input: 0030 2030 3032 0a
tcp echo server [echosrv]: OK (5.0s)
```

echosrv过了，但是testinput还是不过，而且跟答案的差距还非常小，在怎么检查都查不出的情况下，我通过比对张磊同学的代码，发现了我忘记看他提示中的一句话：

```
// LAB 6: Your code here:
// - read a packet from the device driver
// - send it to the network server
// Hint: When you IPC a page to the network server, it will be
// reading from it for a while, so don't immediately receive
// another packet in to the same physical page.
```

即会有**这样的情况**，Input将数据放入nsipcbuf中传给NS，NS处理这些数据需要一些时间。因为Input是轮询查询网卡是不是有新数据的，所以可能在NS还没处理完的时候就利用原来放置nsipcbuf的页面装进了新收到的数据，准备发给NS，这个时候就悲剧了，原来NS正在读取的数据被刷新了。

张磊的做法是每次把nsipcbuf以ipc发送给NS以后，进行一次sys\_yield()我也这么试了，让我很郁闷的是，他进行一次就可以了，**尼玛我一定要yield三次才可以!!!**其实事实上yield多少次都不安全，因为没法知道网络服务器端什么时候才读完所有的数据，允许Input可以在nsipcbuf上写新的东西。

**张顺廷湿胸**跟我说了方法，是所有方法里感觉比较靠谱的一个了，就是当使用ipc传输数据的时候，**每个接受的进程在使用完共享的页面时，应该都会将页面卸载**。（但是不卸载也是没有问题的，因为当你下次使用IPC接受页面时，重新映射的时候如果发现已有装入的页面，系统的IPC接口会自动帮用户卸载，但是一般编写的比较好的程序应该都会卸载）。所以我们可以发送以后记录nsipcbuf页面的引用数，在标准库中提供用户访问页面引用数的接口：

lib/pageref.c

```
1 #include <inc/lib.h>
2
3 int
4 pageref(void *v)
5 {
6     pte_t pte;
7
8     if (!(vpd[PDX(v)] & PTE_P))
9         return 0;
10    pte = vpt[VPN(v)];
11    if (!(pte & PTE_P))
12        return 0;
13    return pages[PPN(pte)].pp_ref;
14 }
```

然后我们可以轮询查询该页面的引用数，当NS使用完该页面的时候卸载页面，那么它的引用数一定会减小，那么这个时候我们就可以知道该页面已被用完，可以装入新数据了。

这个做法还有一个要注意的地方，就是前面提到过，由于Input和Output都是网络服务器的一个模块，是由NS在进程里fork出来的进程，所以**nsipcbuf都是COW页面**，COW页面可能会被很多程序使用，当一个程序对页面进行写操作，那么COW恢复程序就会申请新的页面给改程序，那么原来的COW页面的引用数就会减少了。另一方面，如果另外一个程序执行了fork，那么该COW页面的引用数就会增加。

前面这段说明的就是如果nsipcbuf的页面是COW页面，**那么对其读取引用数是没有意义的，它随时可以改变**，为了避免这种情况，可以Input一开始的时候就可以对nsipcbuf进行一次写操作。那么这个时候COW页面的恢复机制会自动为nsipcbuf拷贝一个新的物理页，那么这个物理页就是唯一专属Input模块了，这个时候再进行读写引用数就是唯一的了。

但是我仍然没有采用这样的方法，因为采用这种方式的前提是我们需要知道网络服务器一定遵照了ipc调用的原则卸载了传输的IPC页面，实际情况下我不清楚它是否这么作了，即如果别人的实现比较丑的情况下，这样的机制仍会挂。

所以我的想法是应该让网络服务器接受和处理完毕以后，向Input发送一条消息，通知他可以继续收数据。但是具体服务器代码我没有看，所以也无所谓改不改了。

### 3.3 The Web Server

**Exercise 11.** The web server is missing the code that deals with sending the contents of a file back to the client. Finish the web server by implementing `send_file` and `send_data`.

这段程序主要是查询文件系统中关于Socket的用法，以及将文件读取连接后输出，都不难，直接贴代码了：

```

user/httpd.c: send_file()
1 static int
2 send_file(struct http_request *req)
3 {
4     int r;
5     off_t file_size = -1;
6     int fd;
7
8     // open the requested url for reading
9     // if the file does not exist, send a 404 error using send_error
10    // if the file is a directory, send a 404 error using send_error
11    // set file_size to the size of the file
12
13    // LAB 6: Your code here.
14    char path[MAXPATHLEN];
15    struct Stat stat;
16

```

```

17     memmove(path, req->url, strlen(req->url));
18
19     if ((fd = open(path, O_RDONLY)) < 0) {
20         send_error(req, 404);
21         goto end;
22     }
23
24     if ((r = fstat(fd, &stat)) < 0) {
25         goto end;
26     }
27
28     if (stat.st_isdir) {
29         send_error(req, 404);
30         goto end;
31     }
32
33     file_size = stat.st_size;
34
35
36     if ((r = send_header(req, 200)) < 0)
37         goto end;
38
39     if ((r = send_size(req, file_size)) < 0)
40         goto end;
41
42     if ((r = send_content_type(req)) < 0)
43         goto end;
44
45     if ((r = send_header_fin(req)) < 0)
46         goto end;
47
48     r = send_data(req, fd);
49
50 end:
51     close(fd);
52     return r;
53 }

```

还有

```

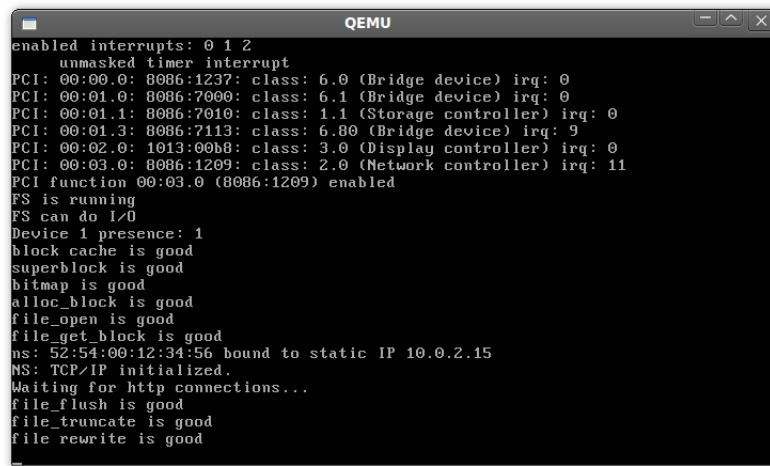
                                user/httpd.c: send_data()
1  static int
2  send_data(struct http_request *req, int fd)
3  {
4      // LAB 6: Your code here.
5      char *buff = malloc(1520 + 10);
6      struct Stat stat;
7      int size, r;
8
9      if (fstat(fd, &stat) < 0) {
10         return 0;
11     }
12     size = stat.st_size;
13
14     if (size > 1520 + 10)
15         return 0;
16     if ((r = readn(fd, buff, stat.st_size)) < 0) {
17         return 0;
18     }
19     if (write(req->sock, buff, stat.st_size) != stat.st_size)
20         cprintf("Failed to send the file\n");
21
22     return 0;
23 }

```

然后按照材料的提示，使用make run-httpd运行服务器，在make which-ports查看QEMU为我们映射的监听端口

```
zhangchi@zhangchi-vostro1400:~/lab$ make which-ports
Local port 26001 forwards to JOS port 7 (echo server)
Local port 26002 forwards to JOS port 80 (web server)
zhangchi@zhangchi-vostro1400:~/lab$
```

于是在浏览器里打开<http://127.0.0.1:26002/index.html>即可看到JOS内运行的Web服务器了：



到这里Lab6终于完全结束了！真的太消耗人了！！！写到这里我想衷心的感谢一些人：

- **李春奇**：春哥提供的程序和讲解帮我从一无所知到终于鼓起勇气开始动手有很大的帮助
- **张磊**：磊牛的代码给了我很多启发，并且提供了非常好的参考资料，省去了我在茫茫信息的大海中自己找寻的麻烦
- **张顺廷**：顺顺湿胸是在整个操统实习课程中对我帮助最大的人，他对我指点纠正了我很多理解上的误区，让我对JOS有了更加深入的认识，同时他也是我见过的对于JOS的见解最到位也最负责的助教了