

Lab6 networking

Lab6 networking

前置知识

网卡接收与传输

实验内容

Your Job(hard)

任务

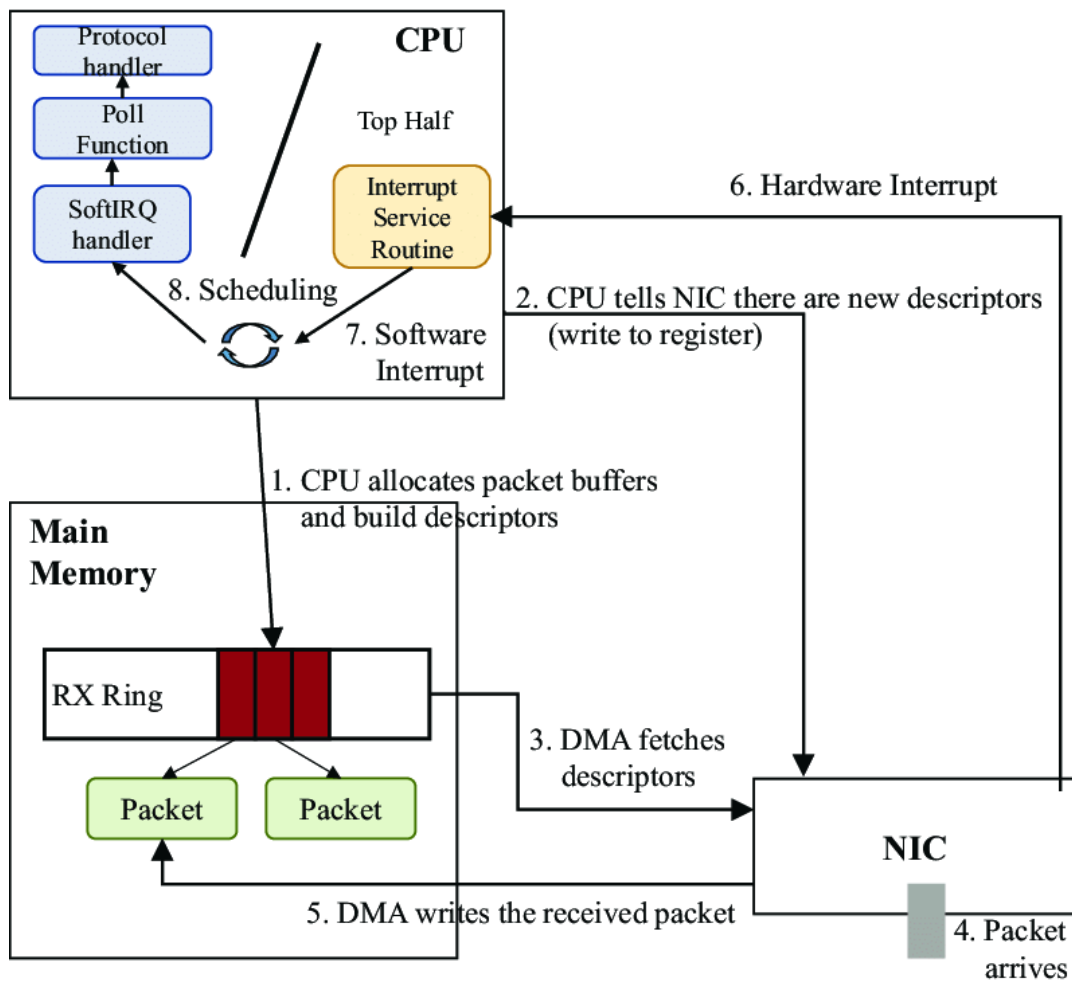
函数实现

socket 实现过程

实验得分

前置知识

网卡接收与传输



- 由这张图我们可以梳理下关于网卡收发包的细节，首先内核需要分配 `rx_ring` 和 `tx_ring` 两块环形缓冲区的内存用来接收和发送报文。其中以太网控制器的寄存器记录了关于 `rx_ring` 和 `tx_ring` 的详细信息。接收 `packet` 的细节如下：
 - a. 内核首先在主存中分配内存缓冲区和环形缓冲区，并由 CPU 将 `rx_ring` 的详细信息写入以太网控制器
 - b. 随后 NIC (Network Interface Card) 通过 DMA 获取到下一个可以写入的缓冲区的地址，当 `packet` 从硬件收到的时候外设通过 DMA 的方式写入对应的内存地址中
 - c. 当写入内存地址后，硬件将会向 CPU 发生中断，操作系统检测到中断后会调用网卡的异常处理函数
 - d. 异常处理函数可以通过由以太网控制寄存器映射到操作系统上的内存地址访问寄存器获取到下一个收到但未处理的 `packet` 的描述符，根据该描述符可以找到对应的缓冲区地址进行读取并传输给上层协议。

实验内容

Your Job(hard)

任务

- 你的任务是完成 `kernel/e1000.c` 文件中的 `e1000_transmit()` 和 `e1000_recv()` 函数，使驱动程序能够发送和接收数据包。目标是通过运行 `make grade` 来确保你的解决方案通过所有测试。

函数实现

```
1  int
2  e1000_transmit(struct mbuf *m)
3  {
4      //
5      // Your code here.
6      //
7      // the mbuf contains an ethernet frame; program it into
8      // the TX descriptor ring so that the e1000 sends it. Stash
9      // a pointer so that it can be freed after sending.
10     //
```

```

11 // 获取 ring position
12 acquire(&e1000_lock);
13
14 uint64 tdt = regs[E1000_TDT];
15 uint64 index = tdt % TX_RING_SIZE;
16 struct tx_desc send_desc = tx_ring[index];
17 if(!(send_desc.status & E1000_TXD_STAT_DD)) {
18     release(&e1000_lock);
19     return -1;
20 }
21
22 if(tx_mbufs[index] != 0){
23     // 如果该位置的缓冲区不为空则释放
24     mbuf_free(tx_mbufs[index]);
25 }
26
27 tx_mbufs[index] = m;
28 tx_ring[index].addr = (uint64)tx_mbufs[index]->head;
29 tx_ring[index].length = (uint16)tx_mbufs[index]->len;
30 tx_ring[index].cmd = E1000_TXD_CMD_RS | E1000_TXD_CMD_EOP;
31 tx_ring[index].status = 0;
32
33 tdt = (tdt + 1) % TX_RING_SIZE;
34 regs[E1000_TDT] = tdt;
35 __sync_synchronize();
36
37 release(&e1000_lock);
38 return 0;
39 }
40
41 static void
42 e1000_recv(void)
43 {
44     //
45     // Your code here.
46     //
47     // Check for packets that have arrived from the e1000
48     // Create and deliver an mbuf for each packet (using
49     net_rx()).
50     //
51     // 获取接收 packet 的位置
52     uint64 rdt = regs[E1000_RDT];

```

```

52     uint64 index = (rdt + 1) % RX_RING_SIZE;
53
54     // acquire(&e1000_lock); // 锁定以进行安全的并发访问
55
56     if(!(rx_ring[index].status & E1000_RXD_STAT_DD)){
57         // 查看新的 packet 是否有 E1000_RXD_STAT_DD 标志, 如果没有, 则直接
        返回
58         return;
59     }
60     while(rx_ring[index].status & E1000_RXD_STAT_DD){
61         // 使用 mbufput 更新长度并将其交给 net_rx() 处理
62         struct mbuf* buf = rx_mbufs[index];
63         mbufput(buf, rx_ring[index].length);
64
65         // 分配新的 mbuf 并将其写入到描述符中并将状态码设置成 0
66         rx_mbufs[index] = mbufalloc(0);
67         rx_ring[index].addr = (uint64)rx_mbufs[index]->head;
68         rx_ring[index].status = 0;
69         rdt = index;
70         regs[E1000_RDT] = rdt;
71         __sync_synchronize();
72
73         // 将数据包传递给net_rx()处理
74         net_rx(buf);
75
76         // 更新index, 继续处理下一个接收到的数据包
77         index = (regs[E1000_RDT] + 1) % RX_RING_SIZE;
78     }
79
80     // release(&e1000_lock); // 在循环结束后释放锁
81 }

```

socket实现过程

- 在类 Unix 操作系统上面, 设备、pipe 和 socket 都要当做文件来处理, 但在操作系统处理的时候需要根据文件描述符来判断是什么类型的文件并对其进行分发给不同的部分进行出来, 我们需要实现的就是操作系统对于 socket 的处理过程。

socket 的读取过程需要根绝给定的 socket 从所有 sockets 中找到并读取 mbuf, 当对应的 socket 的缓冲区为空的时候则需要进行 sleep 从而将 CPU 时间让给调度器, 当对应的 socket 收到了 packet 的时候再唤醒对应的进程:

```

1 sockrecvudp(struct mbuf *m, uint32 raddr, uint16 lport, uint16
  rport)
2 {
3     //
4     // Your code here.
5     //
6     // Find the socket that handles this mbuf and deliver it,
  waking
7     // any sleeping reader. Free the mbuf if there are no sockets
8     // registered to handle it.
9     //
10    acquire(&lock);
11    struct sock* sock = sockets;
12    // 首先找到对应的 socket
13    while(sock != 0){
14        if(sock->lport == lport && sock->raddr == raddr && sock-
  >rport == rport){
15            break;
16        }
17        sock = sock->next;
18        if(sock == 0){
19            printf("[kernel] sockrecvudp: can't find socket.\n");
20            return;
21        }
22    }
23    release(&lock);
24    acquire(&sock->lock);
25    // 将 mbuf 分发到 socket 中
26    mbufq_pushtail(&sock->rxq, m);
27    // 唤醒可能休眠的 socket
28    release(&sock->lock);
29    wakeup((void*)sock);
30 }
31
32 int sock_read(struct sock* sock, uint64 addr, int n){
33     acquire(&sock->lock);
34     while(mbufq_empty(&sock->rxq)) {
35         // 当队列为空的时候, 进入 sleep, 将 CPU
36         // 交给调度器
37         if(myproc()->killed) {
38             release(&sock->lock);
39             return -1;

```

```

40     }
41     sleep((void*)sock, &sock->lock);
42 }
43 int size = 0;
44 if(!mbufq_empty(&sock->rxq)){
45     struct mbuf* recv_buf = mbufq_pophead(&sock->rxq);
46     if(recv_buf->len < n){
47         size = recv_buf->len;
48     }else{
49         size = n;
50     }
51     if(copyout(myproc()->pagetable, addr, recv_buf->head, size)
    != 0){
52         release(&sock->lock);
53         return -1;
54     }
55     // 或许要考虑一下读取的大小再考虑是否释放, 因为有可能
56     // 读取的字节数要比 buf 中的字节数少
57     mbuffree(recv_buf);
58 }
59 release(&sock->lock);
60 return size;
61 }

```

- 写 socket 的过程:

```

1  int sock_write(struct sock* sock, uint64 addr, int n){
2      acquire(&sock->lock);
3      struct mbuf* send_buf = mbufalloc(sizeof(struct udp) +
    sizeof(struct ip) + sizeof(struct eth));
4      if (copyin(myproc()->pagetable, (char*)send_buf->head, addr,
    n) != 0){
5          release(&sock->lock);
6          return -1;
7      }
8      mbufput(send_buf, n);
9      net_tx_udp(send_buf, sock->raddr, sock->lport, sock->rport);
10     release(&sock->lock);
11     return n;
12 }

```

- 关闭 `socket` 的过程:

```
1 void sock_close(struct sock* sock){
2     struct sock* prev = 0;
3     struct sock* cur = 0;
4     acquire(&lock);
5     // 遍历 sockets 链表找到对应的 socket 并将其
6     // 从链表中移除
7     cur = sockets;
8     while(cur != 0){
9         if(cur->lport == sock->lport && cur->raddr == sock->raddr &&
10            cur->rport == sock->rport){
11             if(cur == sockets){
12                 sockets = sockets->next;
13                 break;
14             }else{
15                 sock = cur;
16                 prev->next = cur->next;
17                 break;
18             }
19             prev = cur;
20             cur = cur->next;
21         }
22         // 释放 sock 所有的 mbuf
23         acquire(&sock->lock);
24         while(!mbufq_empty(&sock->rxq)){
25             struct mbuf* free_mbuf = mbufq_pophead(&sock->rxq);
26             mbuffree(free_mbuf);
27         }
28         // 释放 socket
29         release(&sock->lock);
30         release(&lock);
31         kfree((void*)sock);
32     }
```

实验得分

```
chengyu@chengyu-virtual-machine: ~/os-labs/xv6-labs-2023
proc.o kernel/swtch.o kernel/trampoline.o kernel/trap.o kernel/syscall.o kernel/
sysproc.o kernel/bio.o kernel/fs.o kernel/log.o kernel/sleeplock.o kernel/file.o
kernel/pipe.o kernel/exec.o kernel/sysfile.o kernel/kernelvec.o kernel/plic.o k
kernel/virtio_disk.o kernel/e1000.o kernel/net.o kernel/sysnet.o kernel/pci.o ker
nel/start.o kernel/console.o kernel/printf.o kernel/uart.o kernel/spinlock.o
riscv64-linux-gnu-objdump -S kernel/kernel > kernel/kernel.asm
riscv64-linux-gnu-objdump -t kernel/kernel | sed '1,/SYMBOL TABLE/d; s/ .* / /;
/^$/d' > kernel/kernel.sym
make[1]: 离开目录“/home/chengyu/os-labs/xv6-labs-2023”
== Test running nettests ==
$ make qemu-gdb
(17.7s)
== 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
chengyu@chengyu-virtual-machine:~/os-labs/xv6-labs-2023$
```