



结题报告

基于xv6内核的网络子系统实现

姓 名	熊恪峥 肖凯欣 谢哲 涵 赵淇奥
学 号	22920202204622
日 期	2022年12月22日
学 院	信息学院
课程名称	计算机网络与通信

目录

第一章 分工	1
第二章 项目介绍	2
一、 选题介绍	2
1、 xv6简介	2
2、 lwIP简介	2
二、 项目内容	2
第三章 项目实施	4
一、 总体架构	4
二、 设备驱动程序的实现	5
1、 PCI总线的驱动	5
2、 e1000网卡的驱动程序	5
三、 基于LWIP的网络协议栈	7
1、 LWIP的配置	7
2、 LWIP依赖函数的实现	7
3、 LWIP的初始化和工作循环	8
四、 Socket API	10
1、 架构设计	10
2、 Socket API实现（以UDP为例）	11
五、 POSIX NetDB API和DNS解析	14
1、 架构设计	14
2、 基本DNS查询的实现和测试	15
六、 内核改进	16

1、	内存分配改进	16
2、	文件系统改进	18
3、	进程管理改进	19
4、	常用数据结构实现	20
七、	用户态库	20
八、	应用程序	21
1、	nslookup	21
2、	ping和traceroute	22
3、	httpget	24
4、	RFC 867授时服务客户端daytime	24
5、	TCP和UDP的Echo Server	25
第四章	项目管理	29
第五章	项目感想	30
一、	熊恪峥	30
二、	赵淇奥	30
三、	肖凯欣	30
四、	谢哲涵	30
参考文献		31
附录：代码清单		31

第一章 分工

表格 1.1: 成员及分工

成员	分工
熊恪崢	项目总体设计、PCI和网卡驱动程序、Socket API和POSIX NetDB API编写、用户库支持、更先进的内存管理和文件系统支持的实现、daytime程序编写
肖凯欣	ping/traceroute命令实现
赵淇奥	nslookup命令实现
谢哲涵	httpget命令实现

第二章 项目介绍

一、选题介绍

由6.828 Lab 6 [?]得到的灵感，我们将为xv6内核实现完整的网络子系统以及相应的API。

6.828课程中讲授xv6内核 [?]的实现，并在实验中要求学生在称之为JOS的实验性微内核操作系统中添加相应的功能。但事实上，xv6内核本身的结构和现代网络服务器中常常使用的Linux和FreeBSD更为类似。因此，我们好奇，我们是否可以对xv6内核进行修改，来添加对网络功能的支持，而不是在JOS上进行实现。这样，我们能够更好地通过实现相应的功能和适配相应的库来加深对各种网络协议以及Socket API本身的理解。

1、xv6简介

xv6是对Dennis Ritchie's and Ken Thompson的Unix Version 6 (v6)的重新实现。xv6与v6有大致相同的结构和代码风格，但是使用了ANSI C实现，并且为现代的处理器的功能，如对称多核心等做出了相应的适配。它主要为MIT的操作系统教学而设计。

2、lwIP简介

lwIP [?]是一个开源的、轻量级的TCP/IP协议栈，它是一个可移植的、可裁剪的、可靠的、高效的、低成本的TCP/IP协议栈，它可以用于嵌入式系统，如微控制器、单片机、DSP等。lwIP最初是由Swedish Institute of Computer Science的Adam Dunkels开发，现在由开源社区开发和维护。并且被多个技术公司用于商业产品中，如Intel/Altera, Analog Devices, Xilinx, TI, ST 和 Freescale。

二、项目内容

为了实现完整的网络子系统，我们需要实现以下内容：

- 网卡驱动程序
- IP协议栈
- Socket API等应用编程接口支持、DNS等辅助功能的支持
- 用户库支持、网络应用程序
- 操作系统内核其他组件的增强

因此，我们首先实现了Intel e1000网卡的驱动程序。该设备被模拟器QEMU较好地支持，因此易于测试。为了使该设备能够正常运行，需要实现PCI总线的支持。然后，我们移植了LWIP作为IP协议栈，它提供包括IP、ICMP、UDP和TCP协议的支持。为了支持该库的正常运行，我们实现了相应的时钟中断处理程序及其各组件的初始化。基于该协议栈，我们首先提供了以Socket API为代表的网络编程接口的支持，然后实现了以符合POSIX标准的`gethostbyaddr`和`gethostbyname` 为接口DNS的支持。最后，

我们实现了包括基础的TCP/UDP Echo Server、例如nslookup、ping、traceroute、httpget等常见的实用程序，以及符合RFC 867的日期时间协议的授时客户端。

此外，在开发过程中，我们严格采用了版本控制系统Git，以便于团队成员之间的协作开发。在开发流程上，我们坚持通过fork/pull request的方式进行代码的提交和合并，以保证代码的质量、防止不当更改损坏项目的其他部分。我们的工作流程和现代的软件开发工作流程高度一致、完全接轨。实现了有效的项目管理。在开发过程中，我们严格遵循了谷歌开源项目风格指南的代码规范，以保证代码的可读性和可维护性。

我们在开发结束后贯彻“从开源中来、到开源中去”的理念，将我们的项目以MIT协议开源，以供其他人参考和使用。项目地址是https://github.com/SmartPolarBear/xv6_enhanced

第三章 项目实施

一、总体架构

总体的项目架构如图 3.1所示。

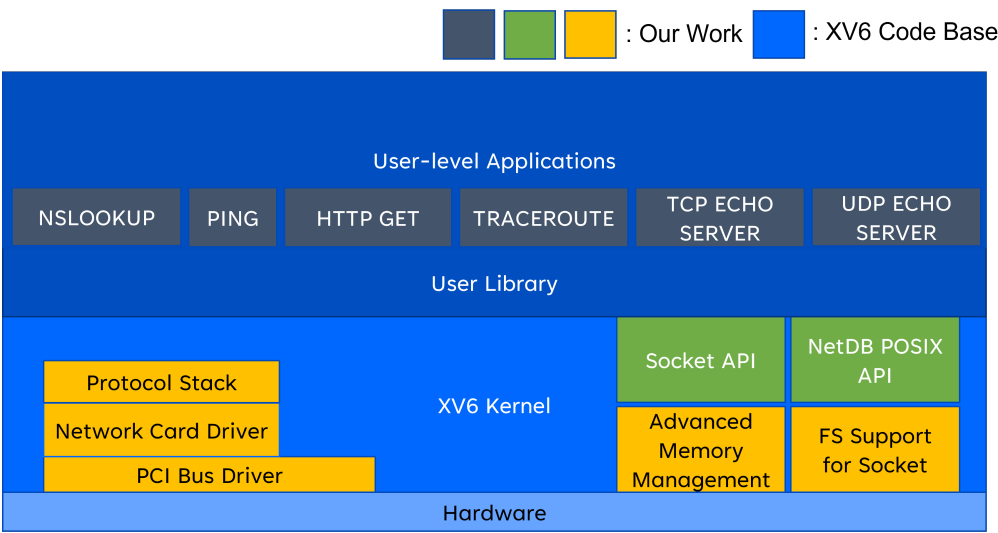


图 3.1: 项目架构

在内核态中的部分按功能分，可以分为三个部分：

一是负责处理实际数据收发、通讯的部分，自下而上分别是PCI总线的驱动、e1000网卡的驱动程序、基于LWIP的网络协议栈。二是作为向上层服务的部分，包括提供一套数据收发原语的Socket API、提供一套使用DNS服务的接口的POSIX NetDB API部分。三是围绕网络子系统内存分配频繁、动态内存分配要求高的需求、Socket API与类UNIX操作系统的文件系统部分结合紧密的特点，对内核进行的改进。例如对内存分配的改进、对文件系统的改进、引入了与Linux内核中实现方式类似的链表等基础数据结构。

在内核和用户之间，用户态库也需要进行相应接口的实现。我们实现了Socket API和用于查询DNS的接口，例如 `gethostbyaddr`和`gethostbyname`的用户态库接口。这些接口完全符合POSIX标准。

在用户态，我们实现了两类的应用程序：

第一类是用来调试、检验基础功能的应用程序。包括TCP、UDP的Echo Server、测试和远程主机实际通信以及DNS的 RFC 867日期时间服务客户端。第二类是常用的网络实用程序，包括`nslookup`、`ping`、`traceroute`、`httpget`等。

二、设备驱动程序的实现

1、PCI总线的驱动

通用串行总线（PCI）是一种用于确保高性能、低开销的传输的本地总线。PCI总线的组件和卡片接口时处理器无关的，因此该总线结构可以被多种处理器架构利用。

为了初始化PCI总线，需要驱动程序访问PCI配置空间，从而获取PCI设备的基本信息，例如设备的ID、设备的基地址等。访问PCI配置空间有两种方式，但只有访问I/O端口的方式是符合标准的。这种方式使用两个32位的I/O端口，一个用于指定配置地址，一个用于读写配置数据。配置地址的格式如图3.2所示。据此，我们实现的PCI驱动程序在代码`kernel/pci.c`中，它会访问PCI配置空间，获取PCI设

图 3.2: PCI配置地址的格式

Bit 31	Bits 30-24	Bits 23-16	Bits 15-11	Bits 10-8	Bits 7-0
Enable Bit	Reserved	Bus Number	Device Number	Function Number	Register Offset ¹

备的基本信息，然后加入设备表中。为了保证该驱动程序的可扩展性，我们的实现将不同驱动程序实现的相应初始化程序注册到数组中，如代码3.1。

代码 3.1: PCI驱动程序的初始化程序注册

```
1 struct pci_driver pci_attach_vendor[] = {
2     {0x8086, 0x100e, &e1000_nic_attach},
3     {0x1af4, 0x1000, &virtio_nic_attach},
4     {0, 0, NULL},
5 };
```

当驱动程序初始化的时候，每当注册新设备，都会在其中查找相应的初始化程序，然后调用该初始化程序。查找的方式是读取设备类和设备ID，分别位于配置空间的偏移0x0b和0x0a中。如果找到相应的初始化程序，就调用该初始化程序。

当PCI驱动程序完成后，操作系统启动时，会调用`pci_init`函数，该函数会遍历所有的PCI设备，完成上述功能，并输出设备的基本信息。输出的信息如图3.3所示。

图 3.3: PCI设备的基本信息

```
Booting from Hard Disk...
PCI: 0:0:0: 8086:1237: class: 6.0 (Bridge device) irq: 0
PCI: 0:1:0: 8086:7000: class: 6.1 (Bridge device) irq: 0
PCI: 0:1:1: 8086:7010: class: 1.1 (Mass storage controller) irq: 0
PCI: 0:1:3: 8086:7113: class: 6.80 (Bridge device) irq: 9
PCI: 0:2:0: 1234:1111: class: 3.0 (Display controller) irq: 0
PCI: 0:3:0: 8086:100e: class: 2.0 (Network controller) irq: 11
    mem region 0: 131072 bytes at 0xfeb80000
    io region 1: 64 bytes at 0xc000
PCI function 0:3:0 (8086:100e) enabled
nic_register: e1000_0 registered
```

2、e1000网卡的驱动程序

有了上述PCI驱动程序和相应扩展架构的实现。我们就可以实现e1000网卡的驱动程序了。e1000网卡的驱动程序的重要组成部分是网卡的初始化程序，该程序会初始化网卡的寄存器，设置相应信息，使网卡进入可用状态。

网卡的初始化

有了完善的PCI驱动程序，e1000网卡的初始化在操作上就是对相应寄存器地址的读写。方法是调用PCI驱动程序的`pciw`和`pcir`函数，分别用于写和读。初始化大致分为以下几个步骤：

1. 从PCI配置空间中读取网卡的基地址寄存器，得到网卡的基地址。
2. 从e1000的eeprom中读取MAC地址。
3. 设置多播表为全0。
4. 设置接受和发送的缓冲区等相关设置。

由于本项目并没有要求大批数据的收发，只需要基本功能，因此e1000网卡工作在轮询模式下，不使用中断。e1000网卡的收、发分别使用64、16个物理内存页作为缓冲区，其中每两个Receive Descriptor和Transmission Descriptor 共用一个缓冲区，因此每个缓冲区的大小是2KB。这样划分可以较好地保证缓冲区的对齐。

驱动程序完成后，会在启动时输出网卡相关的硬件信息，如图 3.4所示。52:54:00:12:34:56是QEMU默

图 3.4: 网卡的硬件信息

```

NIC register: e1000-0 registered
e1000: base: 0xfeb80000 status reg: 80080783
e1000: macaddr: 52:54:0:12:34:56
e1000: init done

```

认的虚拟网卡MAC地址，可见实现是正确的。

网卡的收发

从网卡收发数据是通过写相应的Read/Transmission Descriptor而实现的。默认情况下，这些Descriptor的Status字段为0，表示准备好接收或发送数据。当网卡收到数据时，会将数据写入相应的缓冲区，并将Status字段置为0，命令字段设置RS和EOP，表示数据已经正确写入，要求网卡回报状态。然后再在相应的寄存器中写入使用的Descriptor 的编号。

之后，驱动程序轮询网卡的状态寄存器，当网卡的状态寄存器中的RD和TD字段都为1时，表示数据已经正确接收或发送。其中发送部分的代码如代码 3.2所示。

代码 3.2: 数据发送

```

1  int e1000_net_send(void *state, const void *data, int len)
2  {
3      netdev_t *card = (netdev_t *)state;
4      e1000_t *e1000 = (e1000_t *)card->priv;
5      uint32 tail = pcir(e1000, E1000_TDT);
6      struct TD *next_desc = &e1000->tx_descs[tail];
7
8      if (!(next_desc->status & E1000_TXD_STAT_DD) && (next_desc->cmd & E1000_TXD_CMD_RS))
9      {
10         return -1;
11     }
12
13     if (len > PKTSIZE)
14     {
15         len = PKTSIZE;
16     }
17
18     memmove(e1000->tbuf[tail]->buf, data, len);
19     next_desc->length = len;
20     next_desc->status = 0;
21     next_desc->cmd = (E1000_TXD_CMD_RS | E1000_TXD_CMD_EOP);

```

```
22
23     pciw(e1000, E1000_TDT, (tail + 1) % E1000_TXDESC_LEN);
24     while (!(next_desc->status & 0x0f))
25     {
26         microdelay(1);
27     }
28     return len;
29 }
```

三、基于LWIP的网络协议栈

由于网络协议栈的实现工作量巨大、相当复杂。我们采用了常用的LWIP作为协议栈。该协议栈实现了常见用的TCP/IP协议，如IP、TCP、UDP、ICMP、ARP、DHCP、DNS等。一方面，该协议栈较为轻量级，一方面，配置该库只需要在一头文件中定义相应的宏。因此我们选择了该实现。

1、LWIP的配置

LWIP的配置主要是在lwipopts.h中进行的。该文件中定义了各种宏，用于控制LWIP的功能。

由于协议栈运行在内核态。因此在LWIP库的视角下，LWIP的代码工作在没有操作系统的模式下，因此，我们配置文件中主要的定义如代码 3.3所示。

代码 3.3: LWIP配置

```
1 #define NO_SYS 1
2
3 #define LWIP_NETCONN 0
4 #define LWIP_SOCKET 0
5
6 #define LWIP_ARP 1
7 #define LWIP_DHCP 1
8 #define LWIP_DNS 1
9 #define LWIP_ETHERNET 1
10 #define LWIP_ICMP 1
11 #define LWIP_RAW 1
```

首先，我们关闭了LWIP的socket和netconn接口，因为它们不能工作在无操作系统模式（NO_SYS）下，因此，我们需要禁用这些功能，仅保留基础的功能，然后用基础的API实现我们的Socket等接口。

我们启用了ARP、DHCP、DNS、Ethernet、ICMP和RAW协议。这些协议对于Socket API和POSIX NetDB API的实现是不可或缺的。

2、LWIP依赖函数的实现

NO_SYS下的LWIP只依赖极少数的函数，这些函数主要是用于时间管理、和随机数生成。这些函数实现在kernel/lib/lwip/os/中的文件中。分别是sys_now和lwip_rand。由于LWIP仅利用时间戳进行时间差的计算，不用于任何时间戳的标注。因此可以直接反馈系统时钟中断中记录的Ticks。而随机数的生成，我们采用的是简单的线性同余法。这种方法能使用位运算快速实现，有利于内核态代码的高效工作。随机数生成的代码如代码 3.4所示。

代码 3.4: LWIP随机数生成

```
1 unsigned long lwip_rand(void)
2 {
3     uint32 time = sys_now();
4     static unsigned int z1 = 12345, z2 = 12345, z3 = 12345, z4 = 12345;
```

```

5      unsigned int b;
6      b = ((z1 << 6) ^ z1) >> 13;
7      z1 = ((z1 & time) << 18) ^ b;
8      b = ((z2 << 2) ^ z2) >> 27;
9      z2 = ((z2 & time) << 2) ^ b;
10     b = ((z3 << 13) ^ z3) >> 21;
11     z3 = ((z3 & time) << 7) ^ b;
12     b = ((z4 << 3) ^ z4) >> 12;
13     z4 = ((z4 & time) << 13) ^ b;
14     return (z1 ^ z2 ^ z3 ^ z4);
15 }

```

3、LWIP的初始化和工作循环

LWIP的主要工作都在`kernel/net/net.c`中完成。首先，在NO_SYS模式下，要使用LWIP，需要用`netif->input` 函数将数据传递到LWIP中。需要注意的是，LWIP不原生支持并发，而XV6支持对称多处理器。因此，为了防止LWIP内部的数据结构出现不一致的情况，我们对任何使用LWIP的上下文都使用同一自旋锁保护。这意味着只有一个进程能同时访问网络功能。这虽然一定程度上降低了效率，但是是不可或缺的。

其次，为了驱动LWIP不断收取、分发、处理数据，需要不断调用`sys_check_timeouts`函数。因此，在xv6的IRQ0中断，即时钟中断中，我们调用该函数。这样，LWIP就能初步正常工作了。

为了实现正常的收发数据，我们需要做好网卡驱动程序和LWIP的对接。在初始化时，我们为每一个网卡都创建一个`netif`结构体。然后在其中设置网卡的MAC地址、IP地址、子网掩码、网关地址、MTU等信息。最后，我们调用`netif_add`函数将该网卡加入LWIP 的网络接口列表中。这样，LWIP就能够使用该网卡了。

但在收发数据前，首先，我们需要获取合适的IP地址，而非要求用户手动输入，这符合现代操作系统的实际。因此我们使用了DHCP协议。首先，调用LWIP的`dhcp_start`函数，然后若未获取IP地址，就反复试图从网卡中读入数据并交给LWIP处理。这样，LWIP 最终就能正确地设置网卡的IP地址了。这时，正在启动的操作系统内核会输出获取的IP地址等配置信息。这部分代码如代码??所示。

代码 3.5: 网卡初始化

```

1 void
2 netadd(int n)
3 {
4     struct netif *new = &netif[n];
5     netdev_t *card = find_card_by_id(n);
6
7     int i;
8     char addr[IPADDR_STRLEN_MAX], netmask[IPADDR_STRLEN_MAX], gw[IPADDR_STRLEN_MAX];
9
10    if (!netif_add_noaddr(new, (void *)card, linkinit, netif_input))
11    {
12        panic("netadd");
13    }
14
15    new->name[0] = 'e';
16    new->name[1] = 'n';
17
18    netif_set_link_up(new);
19    netif_set_up(new);
20 }

```

```

21     cprintf("net %d: mac ", n);
22     for (i = 0; i < ETH_HWADDR_LEN; ++i)
23     {
24         if (i)
25         {
26             cprintf(":");
27         }
28         if (new->hwaddr[i] < 0x10)
29         {
30             cprintf("0");
31         }
32         cprintf("%x", new->hwaddr[i]);
33     }
34     cprintf("\n");
35
36     if (!card->opts->open || card->opts->open(card))
37     {
38         panic("netadd: cannot open device\n");
39     }
40
41     dhcp_start(new);
42     /* wait until DHCP succeeds */
43     while (!dhcp_supplied_address(new))
44     {
45         linkinput(new);
46         sys_check_timeouts();
47     }
48
49     ipaddr_ntoa_r(netif_ip_addr4(new), addr, sizeof(addr));
50     ipaddr_ntoa_r(netif_ip_netmask4(new), netmask, sizeof(netmask));
51     ipaddr_ntoa_r(netif_ip_gw4(new), gw, sizeof(gw));
52     cprintf("net %d: addr %s netmask %s gw %s\n", n, addr, netmask, gw);
53
54 }

```

在实际运行中，QEMU虚拟机下的xv6输出如图3.5a所示。使用QEMU参数保存相应的数据包，如图3.5b所示。可见，一个正确的、典型的DHCP过程正常完成了。首先，没有配置好IP地址的网卡发送了一个DHCP Discover数据包，源地址是0.0.0.0。然后位于10.0.2.2的DHCP服务器以广播的方式返回了一个DHCP Offer数据包，这时，网卡继续以0.0.0.0作为源地址发送了一个DHCP Request数据包，请求DHCP服务器分配IP地址。最后，DHCP服务器返回了一个DHCP ACK数据包，这时，网卡就能正确地获取到IP地址、子网掩码和网关地址等重要信息。在此例中，分别为 10.0.2.15、255.255.255.0和10.0.2.2。

图 3.5: DHCP初始化

net 0: mac 52:54:00:12:34:56	1 0.0.0.0	0.0.0.0	255.255.255.255	338 DHCP	DHCP Discover - Transaction ID 8a1049be11
net 0: addr 10.0.2.15 netmask 255.255.255.0 gw 10.0.2.2	2 0.000329	10.0.2.2	255.255.255.255	318 DHCP	DHCP Offer - Transaction ID 8a1049be11
	3 0.000083	0.0.0.0	255.255.255.255	358 DHCP	DHCP Request - Transaction ID 8a1049be11
	4 0.000123	10.0.2.2	255.255.255.255	318 DHCP	DHCP ACK - Transaction ID 8a1049be11

(a) DHCP获取的内容

(b) 收发数据包

由于时钟中断不断驱动LWIP进行工作，当DHCP初始化完后，网卡立即发送、接受ARP数据包，并且能够正常地进行网络通信。所捕捉的数据包如图3.6所示。此时，所有驱动程序和网络协议栈都可以正常地协同工作。为实现重要的用户编程接口打下了良好的基础。

图 3.6: ARP数据包

5	1.004813	RealtekU_12:34:56	Broadcast	42 ARP	Who has 10.0.2.15? (ARP Probe)
6	1.592322	RealtekU_12:34:56	Broadcast	42 ARP	Who has 10.0.2.15? (ARP Probe)
7	1.592576	RealtekU_12:34:56	Broadcast	42 ARP	Who has 10.0.2.15? (ARP Probe)
8	1.592813	RealtekU_12:34:56	Broadcast	42 ARP	ARP Announcement for 10.0.2.15
9	2.594221	RealtekU_12:34:56	Broadcast	42 ARP	ARP Announcement for 10.0.2.15
10	3.595736	RealtekU_12:34:56	Broadcast	42 ARP	ARP Announcement for 10.0.2.15
11	4.594379	RealtekU_12:34:56	Broadcast	42 ARP	ARP Announcement for 10.0.2.15

四、Socket API

Socket API是一组用于网络编程的接口，它是POSIX标准的一部分。它提供了一组函数，来支持主要的用户程序间的通信功能，例如：

- 设置和建立与其它网络中用户的连接；
- 向其它网络中用户发送数据、从其它网络中用户接收数据；
- 关闭连接、释放资源。

此外，它还能够提供一些高级的功能，例如：

- 询问网络系统中的主机和状态；
- 按照其他用户的要求，执行系统功能和控制功能；

在我们的项目中，Socket API处于核心位置，是耗费了最长时间、完成了最多代码的部分。

Internet Socket主要有三类：数据报套接字、流套接字和原始套接字。数据报套接字提供了无连接的Socket服务，主要用于UDP协议；流套接字提供了面向连接的Socket服务，常见的协议有TCP、SCTP等；原始套接字提供了对IP协议的完全访问，它允许用户构造任何类型的IP数据包，包括IP头部和数据部分。前两种对典型的网络应用程序来说十分重要，而原始套接字则是一种高级的功能，对我们在章节八中实现的网络应用程序nslookup、ping和traceroute来说也相当重要，因此我们需要恰当的实现这三类套接字。

LWIP也提供了三种不同的API，分别为RAW IP API、RAW TCP API、RAW UDP API。它们提供了一层较为“薄”的封装，提供了构造这三种类型数据包的基础接口。我们将借助它们实现Socket API。

注意到Socket类型和使用的协议不一定是一一对应的，例如，对于流套接字，用户可以使用TCP协议也可以使用STCP。这样一来，我们并不应该在实现中将利用不同协议的Socket实现耦合在一起，而应当进行逻辑上和实现上的分离，这样才能保证这一部分的可扩展性和可维护性。

1、架构设计

首先，我们将Socket的基本操作定义在结构体sockopt中，如代码 ??所示。

代码 3.6: sockopts结构体

```
1 typedef struct sockopts
2 {
3     struct
4     {
5         int type, protocol;
6     } meta;
7
8     int (*alloc)(struct socket *s);
```

```
9      int (*connect)(struct socket *s, struct sockaddr *addr, int addrlen);
10     int (*bind)(struct socket *s, struct sockaddr *addr, int addrlen);
11     int (*listen)(struct socket *s, int backlog);
12     int (*accept)(struct socket *s, struct socket *newsock, struct sockaddr *addr, int *addrlen)
13         ;
14     int (*send)(struct socket *s, void *buf, int len, int flags);
15     int (*recv)(struct socket *s, void *buf, int len, int flags);
16     int (*sendto)(struct socket *s, void *buf, int len, int flags, struct sockaddr *addr, int
17         addrlen);
18     int (*recvfrom)(struct socket *s, void *buf, int len, int flags, struct sockaddr *addr, int
19         *addrlen);
20     int (*close)(struct socket *s);
21
22     int (*getsockopt)(struct socket *s, int level, int optname, void *optval, int *optlen);
23     int (*setsockopt)(struct socket *s, int level, int optname, void *optval, int optlen);
24 } sockopts_t;
```

该结构体存储了一种针对特定“协议-Socket类型”二元组的Socket实现，它包含相应的元数据，以及一组函数指针。它们将指向具体的Socket实现，例如，对于TCP流套接字，它们将指向TCP流套接字的实现。

每支持一种新的“协议-Socket类型”二元组，我们只需要实现相应的Socket实现，并将其注册到全局数组sockopts中，它定义在文件kernel/net/socket.c中，如代码??所示。

代码 3.7: sockopts数组

```
1 sockopts_t *sockopts[] = {
2     &tcp_opts,
3     &udp_opts,
4     &raw_opts,
5 };
```

它显示了我们支持的三种“协议-Socket类型”二元组，分别是TCP流套接字、UDP数据报套接字和原始套接字。

每一个用户创建的Socket在内核中都会用一个socket结构体记录相应的参数，其中就有一个指针指向 sockopts数组中的某一项，这个指针在创建Socket时通过遍历上述数组来确定，如果无法找到，就可以向用户汇报错误。kernel/net/socket.c中任何其他函数实现的主要部分都是根据该指针来调用相应的Socket实现。并不实际处理通信的过程，而是将这些过程交给具体的Socket实现来处理，它们位于kernel/net/tcpsock.c、kernel/net/udpsock.c、kernel/net/rawsock.c等具体的实现文件中。

这样，就完美实现了具有高度可扩展性的Socket实现。实现了不同协议、不同类型Socket实现在逻辑上和代码构成上的分离。使得调试、维护、扩展都变得非常方便。

上述sockopt结构体中包含了几乎全部的Socket API的定义，但不是所有具体的Socket实现都能够实现全部的API。对于那些不支持的操作，具体的实现只需要提供一个简单的、返回错误代码的函数实现即可。

下面，我们以UDP协议的数据报套接字为例，来介绍一种我们已经支持的协议和类型的Socket API的具体实现。

2、Socket API实现（以UDP为例）

UDP协议的套接字是数据报套接字，它的Socket实现定义在文件kernel/net/udpsock.c中。

首先，数据报Socket有一些操作不能执行，比如listen、accept、connect等操作，根据上述分析，

我们只需要在`udpsock.c`中提供一个简单的、返回错误代码的函数实现即可。例如，`connect`操作的实现如代码 3.8所示，它只是简单地返回一个错误代码。

代码 3.8: UDP的`connect`实现

```
1 int udpconnect(struct socket *s, struct sockaddr *addr, int addrlen)
2 {
3     return -EPROTONOSUPPORT;
4 }
```

按照POSIX标准，失败的操作返回-1，然后设置全局的`errno`变量。因此，我们在这里返回的是详细的错误代码，而系统调用的实现函数会据此设置`errno`变量，然后返回-1。例如，`connect`系统调用的实现如代码 3.9所示。

代码 3.9: `connect`系统调用实现

```
1 int
2 sys_connect(void)
3 {
4     struct file *f;
5     int addrlen;
6     struct sockaddr *addr;
7
8     if (argfd(0, NULL, &f) < 0 || argint(2, &addrlen) < 0 || argptr(1, (void *)&addr, addrlen) <
9         0)
10     {
11         return -1;
12     }
13     if (f->type != FD_SOCKET)
14     {
15         seterror(-ENOTSOCK);
16         return -1;
17     }
18     int err = socketconnect(f->socket, addr, addrlen);
19     if (err < 0)
20     {
21         seterror(err);
22         return -1;
23     }
24     return 0;
25 }
```

而对于那些可以执行的操作，我们需要提供一个具体的实现。这些实现会借助相应的LWIP协议栈的API来完成。以UDP的`recvfrom`操作为例。为了实现该接口，有一些协议无关的操作和一些协议相关的操作。协议无关的操作包括：

1. 检查参数的合法性；
2. 在调用之后、确实收到数据之前，将当前进程阻塞；
3. 从LWIP的`pbuf`缓冲区中拷贝数据到用户空间；

协议相关的操作包括：

1. 为LWIP的UDP协议栈设置相关参数；

2. 实现LWIP的接收回调函数：

为了实现以上的操作，我们进行了如下的设计考量：

- 由于网络操作等待时间较长，不应该进行循环式的“忙等”，因此应当将进程设为睡眠；
- 等待时间不能无限延长，同时Socket API也支持设置超时参数，因此应当处理超时；
- LWIP允许为所有事件设置一个处理函数，或者为特定的事件设置单独的回调函数。为了隔离不同的逻辑，我们选择后者；

协议相关的操作由UDP的Socket实现完成，代码如代码 3.10所示。

代码 3.10: UDP Socket实现

```
1 int udprecvfrom(struct socket *s, void *buf, int len, int flags, struct sockaddr *addr, int *addrlen
2 )
3 {
4     struct udp_pcb *pcb = (struct udp_pcb *)s->pcb;
5     if (!pcb)
6     {
7         return -ECONNRESET;
8     }
9
10    netbegin_op();
11    udp_recv(pcb, lwip_udp_recv_callback, s);
12    netend_op();
13
14    return 0;
15 }
```

而上述的协议无关的操作则由Socket API的顶层封装完成，代码如代码 5.2所示。该段代码较长，我们将其放在附录：代码清单中。

可见，该函数首先检查了参数，然后调用`skt->opts->recvfrom`，对于一个UDP Socket，该函数指针指向了UDP Socket的操作函数。该函数会调用`udp_recv`注册回调函数，并将Socket作为参数传递给回调函数。然后将当前进程置于休眠状态。

在回调函数中，会将接收到的数据包保存到Socket的接收缓冲区中，然后唤醒等待的进程。如代码 3.11所示。

代码 3.11: LWIP的接受回调

```
1 static inline void lwip_udp_recv_callback(void *arg, struct udp_pcb *pcb, struct pbuf *p,
2                                             const ip_addr_t *
3                                             addr, u16_t
4                                             port)
5 {
6     socket_t *socket = (socket_t *)arg;
7     socket->recvfrom_params.recv_addr.sin_addr.addr = addr->addr;
8     socket->recvfrom_params.recv_addr.sin_port = port;
9     socket->recvfrom_params.recv_addr.sin_family = AF_INET;
10
11    if (!p || p->tot_len > socket->recvfrom_params.recv_len)
12    {
13        if (p)
14        {
15            // ...
16        }
17    }
18 }
```



```
13         netbegin_op();
14         pbuf_free(p);
15         netend_op();
16     }
17     socket->recv_closed = TRUE;
18     return;
19 }
20 // buffer hasn't been received
21 if (socket->recv_buf)
22 {
23     return;
24 }
25
26 socket->recv_buf = p;
27 socket->recv_offset = 0;
28
29 socket->wakeup_retcode = 0;
30 wakeup(&socket->recv_chan);
31 }
```

这样，可以灵活扩展的Socket `recvfrom`就实现了。其它的操作实现也类似。其中，带休眠时间限制的休眠函数`sleepddl`可参见章节六 中的相关介绍。

其它协议Socket实现也是类似的，可参见源代码中的`kernel/net/tcpsock.c`、`kernel/net/udpsock.c`、`kernel/net/rawsock.c`等具体的实现文件，这里就不再赘述了。所有Socket API的测试都是与TCP/UDP Echo Server 配合完成的，可参见章节 八。

五、POSIX NetDB API和DNS解析

为了实现网络应用程序的方便，仅有Socket API是不够的。因为在现实场景中，用户极少会指定IP进行操作，而是使用域名。因此，需要实现一个域名解析的功能，即DNS解析。

现代类UNIX系统中，经常使用的接口是`getaddrinfo`和`getnameinfo`，然而这组API的实现比较复杂，功能比较多，不利于我们形成可以正常运行的系统。因此，我们选择实现较为传统的`gethostbyname`和`gethostbyaddr`接口。这两个接口虽然已经被POSIX标准弃用，但是实现简单、功能更为正交，因此更适合我们的实现。

1、架构设计

近年来，传统DNS解析的问题已经引发了广泛的关注。因此，新的DNS标准也正在普及，例如DNS over TLS、DNS over HTTPS等。这些新的标准都是为了解决DNS的安全性问题，已经被Windows、Linux等主流操作系统所支持。因此，考虑到未来的发展，我们的DNS实现的架构也应当具有可扩展性。

另一个需要考虑的问题是，`gethostbyname`和`gethostbyaddr`的返回值是连续的内存块，其中变长内容由相对起始的偏移量表示，需要进行计算。因此，不把查询逻辑和构造返回值的逻辑分离开来，就会导致代码难以维护、难以扩展。

因此，我们将DNS的基本操作抽象为一个`netdb_proto`结构体，该结构体包含了DNS的基本操作，如代码 3.12所示。

代码 3.12: DNS操作结构体

```
1 typedef struct netdb_proto
2 {
3     void (*proto_init)();
```

```
4     int (*proto_query)(char *name, netdb_query_type_t type, netdb_answer_t **answer);
5     int (*proto_free_answer)(netdb_answer_t *answer);
6     int (*proto_dump_answer)(const netdb_answer_t *answer);
7     void (*proto_shutdown)();
8 } netdb_proto_t;
```

其中，`proto_init`用于初始化DNS操作，`proto_query`用于查询，`proto_free_answer`用于释放查询结果，`proto_dump_answer`用于打印查询结果，`proto_shutdown`用于操作系统内核关闭时的清理。

为了分离构造返回值的逻辑，我们底层的DNS实现将DNS的查询结果构造成一个`netdb_answer`结构体组成的链表，该结构体包含了DNS的查询结果，而系统调用的实现会进一步将该链表转换为合理的`hostent`返回值。

这样，架构设计就能实现高度的可扩展性，同时达到了模块中高内聚、模块间低耦合的目的。

2、基本DNS查询的实现和测试

普通DNS使用明文传输、基于UDP协议。我们利用LWIP的UDP RAW API实现了DNS数据的传输。

而构造DNS报文的过程则比较复杂，因为DNS将字符串部分进行了压缩，为了构造和解析DNS报文，需要实现相应的解析逻辑。RFC 1035中对DNS报文的格式进行了详细的描述：

In order to reduce the size of messages, the domain system utilizes a compression scheme which eliminates the repetition of domain names in a message. In this scheme, an entire domain name or a list of labels at the end of a domain name is replaced with a pointer to a prior occurrence of the same name.

简而言之，DNS报文中的字符串部分会被压缩，压缩的方式是将字符串部分的每个子串都用一个指针指向之前出现过的子串，高位的两位用于标识该指针是否指向字符串的开始，如果是，则该指针指向的是字符串的开始。代码 3.13给出了字符串转化为正常字符串的实现。

代码 3.13: 字符串压缩解压缩

```
1 static inline char *parse_dns_name(char *dst, char *name, char *ans)
2 {
3     char *p = name;
4     char *d = dst;
5     char *s = ans;
6     int len = 0;
7     int i = 0;
8     int offset = 0;
9     int first = 1;
10
11     while (*p != 0)
12     {
13         if ((*p & 0xc0) == 0xc0)
14         {
15             offset = (*p & 0x3f) << 8;
16             p++;
17             offset |= *p;
18             p = s + offset;
19             continue;
20         }
21     }
```

```
22         len = *p;
23         p++;
24
25         if (first)
26         {
27             first = 0;
28         }
29         else
30         {
31             *d = ',';
32             d++;
33         }
34
35         for (i = 0; i < len; i++)
36         {
37             *d = *p;
38             d++;
39             p++;
40         }
41     }
42
43     *d = 0;
44
45     return p + 1;
46 }
```

然后，我们需要实现DNS报文的构造。该代码较长，因此我们将其放在了附录中，见附录：代码清单 中的代码 5.1。接着，在LWIP UDP接收数据的回调函数中，我们实现了DNS回报的解析。这些部分大体就是对协议内容的简单实现。

`proto_query`的查询既要支持普通的DNS查询，也要支持从IP地址到域名的反向查询。但该功能不需要实现额外的功能，RFC 1035中定义了一种特殊的域，即`in-addr.arpa`，该域中的域名用于IP到域名的反向查询。因此，我们只需要将IP地址转化为`in-addr.arpa`域中的域名，然后进行正常的DNS查询，并取结果中的PTR记录即可。RFC中的相关描述如下：

The Internet uses a special domain to support gateway location and Internet address to host mapping. Other classes may employ a similar strategy in other domains. The intent of this domain is to provide a guaranteed method to perform host address to host name mapping, and to facilitate queries to locate all gateways on a particular network in the Internet.

由于本函数实现过长，故不在报告中完整呈现，请参见源代码中`kernel/net/baredns.c`中的相关内容。

基本DNS功能的测试依赖于用户应用程序`nslookup`的实现，具体数据和内容见应用程序章节的`nslookup`小节。该API连通此命令一同测试，能够正常工作。

六、内核改进

1、内存分配改进

原始的xv6内核中，内存分配只有`kalloc`和`kfree`两个函数，分别用于分配和释放内核内存。该函数的作用是按页分配内存，每次分配的内存大小为一页，即4KB。这种分配方式存在两个问题：

1. 由于每次分配的内存大小为一页，因此分配的内存大小不可控，如果需要分配的内存大小不是一

页的整数倍，那么就会造成内存浪费。

2. 当分配超过一页时，难以保证分配的内存是连续的。

然而，在网络驱动程序、网络协议栈、Socket API等部分，需要频繁地进行内存分配，而且分配的内存大小不固定、分配的大小常常远小于4KB。为了避免耗尽内存，提高内存分配的效率。需要改进内存分配的方式。

常见的类UNIX操作系统，例如Linux有更加复杂的内存分配机制。例如Buddy/Slab分配器，可以分配任意大小的内存块。其中Buddy分配器分配2幂次大小的内存块，而Slab分配任意指定大小的内存块。这种分配方式的优点是可以分配任意大小的内存块，然而完整的实现相当复杂，不利于我们按期完成项目。因此，我们采用了一种简单的内存分配方式，即利用xv6现有的分配器分配页大小的内存。然后使用Slab分配器分配任意大小的内存块。具体的实现在kernel/slab.c中。

Slab分配器的原理是将内存分为多个大小相同的内存块，每个内存块称为一个Slab。Slab分配器维护一个Slab链表，每个Slab中都有一个空闲链表，用于记录空闲的内存块。当需要分配内存时，Slab分配器从空闲链表中取出一个内存块，当需要释放内存时，Slab分配器将内存块插入空闲链表。当Slab中的内存块全部被分配时，Slab分配器从Slab链表中取出一个Slab，当Slab中的内存块全部被释放时，Slab分配器将Slab插入Slab链表。它们操作的内存空间均由xv6的分配器分配。

我们的实现接口大致与Linux相应组件的接口相同。kernel/slab.h中定义了Slab的存储结构，如代码3.14所示。

代码 3.14: Slab的存储结构

```
1 typedef struct slab
2 {
3     kmem_cache_t *cache;
4     size_t in_use;
5     size_t next_free;
6     void *objects;
7     kmem_bufctl *bufctl;
8
9     spinlock_t lock;
10
11     list_head_t link;
12 } slab_t;
```

而这些Slab由称之为Cache的结构管理。Cache的存储结构如代码3.15所示，在kernel/include/slab.h中。

代码 3.15: Cache的存储结构

```
1 typedef struct kmem_cache
2 {
3     list_head_t full, partial, free;
4
5     list_head_t link;
6
7     size_t obj_size, obj_count;
8     uint flags;
9     spinlock_t lock;
10
11     char name[KMEM_CACHE_NAME_MAXLEN];
12 } kmem_cache_t;
```

有了这些结构，我们就可以实现`kmem_cache_create`、`kmem_cache_alloc`、`kmem_cache_free`等接口。

在这些接口之上，我们也实现了`kmalloc`，来方便不创建对应Cache，只指定大小的内存分配。这些接口在驱动程序、Socket API和DNS相关的实现中发挥了关键的作用。

2、文件系统改进

xv6没有常见的VFS结构，但是其文件系统实现包含了对不同类型的文件节点的初步支持。由于Socket API在分配、关闭Socket时需要创建相应类型的文件节点来返回File Descriptor，因此我们在xv6的文件系统实现上做了一些扩展。

`kernel/fs.h`中定义了文件节点的存储结构，如代码 3.16所示。

代码 3.16: 文件节点的存储结构

```
1
2 typedef enum file_type
3 {
4     FD_NONE = 1,
5     FD_PIPE,
6     FD_DEVICE,
7     FD_SOCKET,
8 } file_type_t;
9
10 typedef struct file
11 {
12     file_type_t type;
13     int ref; // reference count
14     char readable;
15     char writable;
16
17     union
18     {
19         struct pipe *pipe;
20         struct inode *ip;
21     };
22     uint off;
23 }file_t;
```

我们在类型中加入了`FD_SOCKET`，并在`union`中加入了`socket`成员。这样就存储了Socket的相关信息。

负责处理文件节点的函数包括`file_alloc`、`file_close`、`file_dup`等。我们在相应函数中加入了对于Socket的处理，调用Socket相关的具体实现，这里举一例说明，如代码 3.17所示。

代码 3.17: `file_alloc`函数的修改

```
1 void
2 fileclose(struct file *f)
3 {
4     struct file ff;
5
6     acquire(&ftable.lock);
7     if (f->ref < 1)
8     {
9         panic("fileclose");
10    }
11    if (--f->ref > 0)
12    {
```

```
13         release(&ftable.lock);
14         return;
15     }
16     ff = *f;
17     f->ref = 0;
18     f->type = FD_NONE;
19     release(&ftable.lock);
20
21     if (ff.type == FD_PIPE)
22     {
23         pipeclose(ff.pipe, ff.writable);
24     }
25     else if (ff.type == FD_INODE)
26     {
27         begin_op();
28         iput(ff.ip);
29         end_op();
30     }
31     else if (ff.type == FD_SOCKET)
32     {
33         socketclose(ff.socket);
34     }
35 }
```

在这段代码中，遇到Socket类型的文件节点时，调用了`socketclose`函数实现Socket专属的关闭逻辑。

3、进程管理改进

本项目中对进程管理的改进主要包含对信号的支持和对可设置超时时间的进程休眠的支持。对信号的支持主要是便于调试。并且实现复杂、代码较长。因此在报告中不做叙述。可参见源代码中的`kernel/proc.c`。这里主要介绍对可设置超时时间的进程休眠的支持。该函数的实现在`kernel/proc.c`中，如代码 3.18所示。

代码 3.18: `sleepddl`的实现

```
1
2 uint sleepddl(void *chan, struct spinlock *lk, uint duration)
3 {
4     if (duration == 0)
5     {
6         sleep(chan, lk);
7         return 0xffffffff;
8     }
9
10    myproc()->sleep.deadline = duration * tickspersecond;
11    sleep(chan, lk);
12    int ret = myproc()->sleep.deadline;
13    myproc()->sleep.deadline = 0;
14    return ret;
15 }
```

它与xv6原有的`sleep`函数非常相似，唯一区别在于，它可以设置一个超时时间。在唤醒之后，会将剩余的时间返回。如果没有超时，返回十六进制值`0xffffffff`。

为了对唤醒时间进行计数，在xv6的调度器实现中，每当下一个时间片到来时，会遍历所有进程，如果某进程使用了带超时时间的休眠，就将其剩余时间减一。如果剩余时间为零，就设置该进程为可运行的进程。在不久的将来，该进程就会被调度、继续执行。

该代码片段如代码 3.19所示。

代码 3.19: 定时唤醒的实现

```

1      if (p->state == SLEEPING && p->sleep.deadline > 0)
2      {
3          p->sleep.deadline--;
4          if (p->sleep.deadline == 0)
5          {
6              p->state = RUNNABLE;
7          }
8      }

```

该功能对Socket API的实现非常有用。在Socket API中，我们常常需要为休眠设置一个超时时间，来避免进程因为等待数据而无限期地休眠。

4、常用数据结构实现

设备驱动程序、DNS的查询过程中，我们需要一些能够动态地插入和删除元素的数据结构。因此，我们实现了Linux类似的侵入式链表。

侵入式（Intrusive）链表是一种链表的实现方式，它将链表的节点单独定义为一个结构体，并将链表节点的指针作为结构体的成员，这是链表中相应的指针都处于数据的内部，所以称之为侵入式链表。

这种链表实现的好处在于可以用同一套链表操作函数来操作不同的具体数据。当需要从链表节点获取数据时，可以通过结构体中内存偏移的计算来获取数据的指针，从而实现链表节点和数据的互相获取。这种计算的方法在Linux内核中被称为container_of。代码 3.20给出了这种计算的实现。

代码 3.20: container_of宏的实现

```

1  // offset of a struct member
2  #define offsetof(type, member) __builtin_offsetof (type, member)
3
4  // get struct pointer from member
5  #define container_of(ptr, type, member) ({
6      const typeof( ((type *)0)->member ) *__mptr = (ptr);
7      (type *) ( (char *)__mptr - offsetof(type,member) );})

```

而其余的链表操作函数实现定义在include/list.h中。包括插入、删除、遍历等操作。这些操作函数的实现都是基于list_head结构体的。

七、用户态库

用户态库是用户态和系统调用实现间的“胶水”。它对用户提供了POSIX兼容的编程接口。将相关参数按照系统调用的约定存入内存，然后调用INT 0x80中断，将控制权交给内核。用户态库的实现在lib/ulib目录下。

由于在32位环境下，传参数不会用到寄存器，因此已经正确地将参数存入内存中，只需要使用int指令调用相应终端即可，这部分代码实现在lib/ulib/usys.S中。如代码 3.21所示。

代码 3.21: usys.S中的系统调用实现

```

1  #define SYSCALL(name) \

```



```
2  .globl name; \  
3  name: \  
4      movl $SYS_ ## name, %eax; \  
5      int $T_SYSCALL; \  
6      ret
```

这定义了一个宏，用于定义系统调用的实现。宏中的SYS_name是系统调用的编号，T_SYSCALL是中断号。

这样就可以简单地用系统调用的名称迅速定义大量的系统调用库函数接口了。

八、应用程序

1、nslookup

程序介绍

nslookup是一个用来查询因特网域名服务器的程序。程序允许使用者向域名服务器查询一个域名或地址的信息，通常被用来诊断域名解析系统的基本结构信息，查询域名解析是否正常工作，在网络故障时排查网络问题。

功能分析

一个完整的nslookup程序拥有两种使用模式：交互式与 交互式。交互式的nslookup允许用户交互式查询多个主机/域名信息，并且具有多种交互命令用于处理复杂的需求、提升用户体验。而 交互式的nslookup只允许用户一次性查询单个主机/域名信息。 交互式的nslookup也 持通过命令 指定参数。本项 的 标是实现一个简易的nslookup程序，因此此程序只实现了nslookup的核 功能： 交互地 对域名或地址进 查询，并向用户显示查询结果。

程序实现

在本项目开始时，我尝试通过如下方法实现程序的查询功能：

1. 通过socket编程与域名服务器通信
2. 构造DNS QUERY报文，向域名服务器发起域名查询请求
3. 解析DNS RESPONSE报文，解析由域名服务器返回的主机信息

但是该方法 对的最大的障碍是，由于DNS复杂的压缩规则，通过编程解析DNS REPLY报文实现难度较 。

因此我最后通过调用内核中已经实现的函数：gethostbyname与gethostbyaddr，获得主机信息。内核中，netdb.h中实现了一组用来获取主机信息的函数，其中包括gethostbyname与 gethostbyaddr。这组函数会返回一个保存有主机信息的hostent结构体。其定义如下代码 3.22所示。

代码 3.22: hostent结构体定义

```
1  struct hostent {  
2      char    *h_name;          /* official name of host */  
3      char    **h_aliases;      /* alias list */  
4      int     h_addrtype;       /* host address type */  
5      int     h_length;         /* length of address */
```



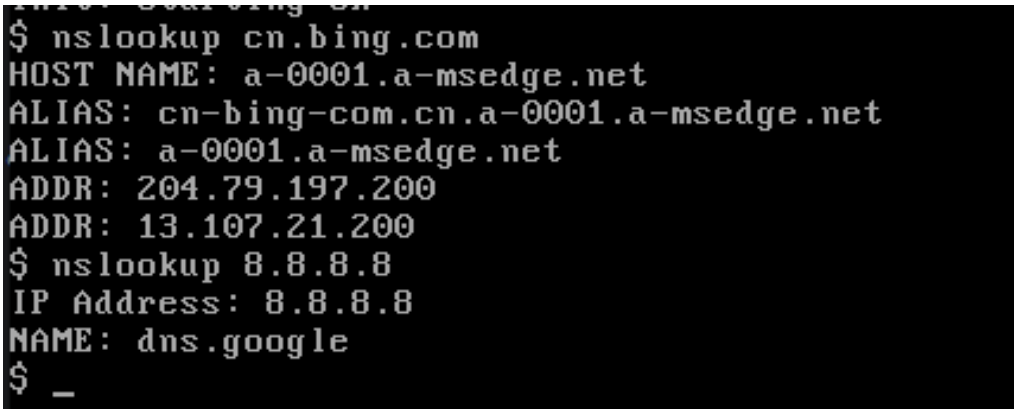
```
6     char    **h_addr_list; /* list of addresses from name
7 server */
8     };
```

程序具体实现如下：

- 1. 首先通过命令行接受用户输入参数，并对用户输入作初步的输入检查
- 2. 对于用户的输入，使用inet_pton判断输入类型，对地址和域名分别处理
- 3. 对于输入的域名，通过调用gethostbyname获取描述该主机的hostent结构体
- 4. 对于输入的地址，通过调用gethostbyaddr获取描述该主机的hostent结构体
- 5. 对于获取的hostent结构体，使用用户友好的形式输出显示。

具体的运行效果如图 3.7所示。首先，对域名cn.bing.com进行查询，程序输出了该域名的所有IP地

图 3.7: nslookup运行效果



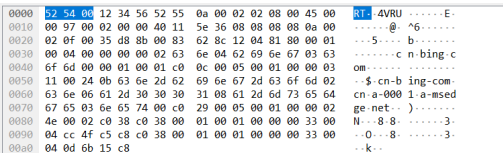
址，以及该域名的的所有别名。然后，对IP地址8.8.8.8反向查询域名。该IP是谷歌的DNS服务器，程序返回的域名为dns.google。可见，该程序是正确的。这其中产生的包如图 3.8所示。 经过检查，以上

图 3.8: nslookup产生的包

14 5.607333	10.0.2.15	8.8.8.8	71 DNS	Standard query 0x1204 A cn.bing.com
15 5.630922	8.8.8.8	10.0.2.15	165 DNS	Standard query response 0x1204 A cn.bing.com CNAME cn-bing-com.cn.a-0001.a-msedge.net CNAME
16 8.703350	10.0.2.15	8.8.8.8	80 DNS	Standard query 0x1204 PTR 8.8.8.8.in-addr.arpa
17 8.725395	8.8.8.8	10.0.2.15	104 DNS	Standard query response 0x1204 PTR 8.8.8.8.in-addr.arpa PTR dns.google

程序的输出符合对应包的内容。尤其是DNS的回复，如图 3.9，可以发现根据DNS规则进行了比较多的压缩，然而我们的程序依然正常地还原了域名。

图 3.9: nslookup过程的回复



2、ping和traceroute

ping的实现

利用原始套接字发送ICMP报文，将ICMP报文头的消息类型设置为8（ICMP request）。ICMP_HDR报
文头如代码 3.23所示。

代码 3.23: ICMP_HDR结构体定义

```

1 struct icmp_hdr {
2     unsigned char type;           /* message type */
3     unsigned char code;          /* type sub-code */
4     unsigned short checksum;
5     unsigned short id;
6     unsigned short sequence;
7     unsigned int timestamp;      /* not part of ICMP, but we need it */
8 };

```

然后，通过原始套接字实现的ICMP报文的发送，如代码 3.24所示。

代码 3.24: ICMP报文发送函数

```

1 struct sockaddr_in dest;
2 dest.sin_family = AF_INET;
3 dest.sin_port = htons(0);
4 dest.sin_addr.s_addr = ip;
5 char buff[sizeof(ICMP_HDR) + 32];
6 ICMP_HDR *pIcmp = (ICMP_HDR *) buff;
7 pIcmp->icmp_type = 8;
8 pIcmp->icmp_code = 0;
9 pIcmp->icmp_id = (uint16_t) pID;
10 pIcmp->icmp_checksum = 0;
11 pIcmp->icmp_sequence = 0;
12 memset(&buff[sizeof(ICMP_HDR)], 'E', 32);

```

然后解析接收到的ICMP报文，判断是否为ICMP reply，如果是，则计算往返时间，否则丢弃并报错。如代码 3.25所示。

代码 3.25: ICMP报文接收函数

```

1 ICMP_HDR *pRecvIcmp = (ICMP_HDR *) (recvBuf + 20); // (ICMP_HDR*)(recvBuf + sizeof(IPHeader));
2 IPHeader *ipHeader = (IPHeader *) recvBuf;

```

其中IPHeader结构体定义如代码 3.26所示。

代码 3.26: IPHeader结构体定义

```

1 typedef struct IPHeader {
2     uint8_t iphVerLen;
3     uint8_t ipTOS;
4     uint16_t ipLength;
5     uint16_t ipID;
6     uint16_t ipFlags;
7     uint8_t ipTTL;
8     uint8_t ipProtocol;
9     uint16_t ipChecksum;
10    uint32_t ipSource;
11    uint32_t ipDestination;
12 } IPHeader;

```

由此可以显示相应的信息，如代码 3.27所示。

代码 3.27: ping显示信息

```

1 printf(" %d bytes from %s:", (int) nRet, inet_ntoa(from.sin_addr));

```

```
2 printf(" icmp_seq = %d. ", pRecvIcmp->icmp_sequence);
3 printf(" ttl = %d. ", ipHeader->ipTTL);
4 printf(" time: %d ms", (int) nTick - (int) pRecvIcmp->icmp_timestamp);
```

traceroute的实现

与ping类似，traceroute从1开始设置ttl，当报文每经过一个路由节点时，ttl减小1，当ttl为0时，节点丢弃报文，并向主机发送一个超时报文，报文的ip数据报中就有该节点的ip地址；随后增大ttl，获取下一个节点，就可以得知主机和目的地址之间的所有路由节点的ip地址

但是，因为部分节点不会响应icmp报文，发送超时报文，所以要为套接字设置超时。其中设置超时的代码如代码 3.28所示。

代码 3.28: 设置超时

```
1 struct timeval timeout = {5,0};
2 setsockopt(sRaw, SOL_SOCKET, SO_RCVTIMEO, (char *)&timeout, sizeof(struct timeval));
```

而设置ttl的代码如代码 3.29所示。

代码 3.29: 设置TTL

```
1 setsockopt(sRaw, IPPROTO_IP, IP_TTL, &ttl, sizeof(ttl))
```

完成实现后，我们在xv6中运行ping，输出如图 3.10所示。ping首先通过gethostbyname函数获取www.baidu.com的

图 3.10: ping输出

```
init: starting sh
6 ping www.baidu.com
ping www.baidu.com
64 bytes from 14.215.177.38: icmp_seq = 256. ttl = 255. time: 3 ms
64 bytes from 14.215.177.38: icmp_seq = 512. ttl = 255. time: 3 ms
64 bytes from 14.215.177.38: icmp_seq = 768. ttl = 255. time: 3 ms
64 bytes from 14.215.177.38: icmp_seq = 1024. ttl = 255. time: 4 ms
```

址，然后发送ICMP消息进行测试，并输出了四次结果。产生的报文如图 3.11所示。可见该过程是正常

图 3.11: ping产生的报文

10.0.2.15	8.8.8.8	73 DNS	Standard query 0x1204 A www.baidu.com
8.8.8.8	10.0.2.15	112 DNS	Standard query response 0x1204 A www.baidu.com CNAME www.a.shifen.com A 14.215.177.39 A 14.215.177.38
10.0.2.15	14.215.177.39	78 ICMP	Echo (ping) request id=0xa4a9, seq=1/256, ttl=255 (reply in 17)
14.215.177.39	10.0.2.15	78 ICMP	Echo (ping) reply id=0xa4a9, seq=1/256, ttl=255 (request in 16)
10.0.2.15	14.215.177.39	78 ICMP	Echo (ping) request id=0xa4a9, seq=2/512, ttl=255 (reply in 19)
14.215.177.39	10.0.2.15	78 ICMP	Echo (ping) reply id=0xa4a9, seq=2/512, ttl=255 (request in 18)
10.0.2.15	14.215.177.39	78 ICMP	Echo (ping) request id=0xa4a9, seq=3/768, ttl=255 (reply in 21)
14.215.177.39	10.0.2.15	78 ICMP	Echo (ping) reply id=0xa4a9, seq=3/768, ttl=255 (request in 20)
10.0.2.15	14.215.177.39	78 ICMP	Echo (ping) request id=0xa4a9, seq=4/1024, ttl=255 (reply in 23)
14.215.177.39	10.0.2.15	78 ICMP	Echo (ping) reply id=0xa4a9, seq=4/1024, ttl=255 (request in 22)

的。

3、httpget

4、RFC 867授时服务客户端daytime

Daytime Protocol是1983年由RFC 867定义的，用于授时服务，客户端通过TCP连接或者到服务器，服务器使用的端口号是13，服务器返回一个ASCII字符串，表示当前的时间，格式为：Day Mon dd hh:mm:ss yyyy

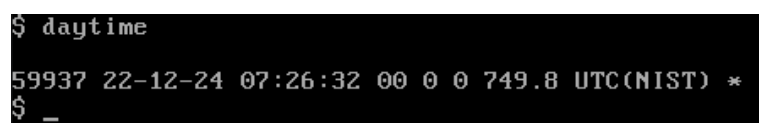
为了测试与远程服务器的连接，我们可以使用这个协议，客户端的代码如代码 3.30所示。该实现使用Socket编程，使用了TCP连接，连接到服务器后，发送一个空的数据包，服务器会返回一个32位的整数，表示从1900年1月1日到现在的秒数，客户端可以通过这个时间戳来输出一个表示时间的用户可读的字符串。

代码 3.30: daytime客户端

```
1 #define SERVER_PORT 13
2 #define SERVER_HOSTIP 0x2c8c8a80
3 char buf[512];
4
5 int main(int argc, char **argv)
6 {
7     // struct hostent *hp;
8     int sockfd, r;
9     struct sockaddr_in addr = {
10         .sin_family = PF_INET, .sin_port = htons(SERVER_PORT),
11     };
12
13     addr.sin_addr.s_addr = htonl(SERVER_HOSTIP);
14
15     sockfd = socket(PF_INET, SOCK_STREAM, IPPROTO_TCP);
16
17     r = connect(sockfd, (struct sockaddr *)&addr, sizeof(struct sockaddr_in));
18     if (r < 0)
19     {
20         printf(1, "daytime: connect failed: %d\n", r);
21         exit(-1);
22     }
23
24     ssize_t n;
25
26     n = recv(sockfd, buf, sizeof(buf));
27     if (n <= 0)
28     {
29         goto end;
30     }
31
32     write(1, buf, n); // stdout
33
34 end:
35     close(sockfd);
36     return 0;
37 }
```

该实现建立了一个TCP连接，连接到服务器后，发送一个空的数据包，服务器会返回一个表示当前日期时间的ASCII字符串，客户端将其输出到标准输出。效果如图3.12所示。位于utcnist.colorado.edu的

图 3.12: daytime客户端输出



```
$ daytime
59937 22-12-24 07:26:32 00 0 0 749.8 UTC(NIST) *
$ _
```

日期时间服务器返回了调用时的日期时间。这其中产生的包如图3.13所示。可见这是一个典型的TCP连接、传输、断开的过程。我们的Socket实现正确处理了TCP连接的全过程。

图 3.13: daytime产生的报文

10.0.2.15	128.138.140.44	58 TCP	60093 → 13 [SYN] Seq=0 Win=2144 Len=0 MSS=536
128.138.140.44	10.0.2.15	60 TCP	13 → 60093 [SYN, ACK] Seq=0 Ack=1 Win=65535 Len=0 MSS=1460
10.0.2.15	128.138.140.44	54 TCP	60093 → 13 [ACK] Seq=1 Ack=1 Win=2144 Len=0
128.138.140.44	10.0.2.15	105 DAYTIME	DAYTIME Response
128.138.140.44	10.0.2.15	60 TCP	13 → 60093 [FIN, ACK] Seq=52 Ack=1 Win=65535 Len=0
10.0.2.15	128.138.140.44	54 TCP	60093 → 13 [ACK] Seq=1 Ack=53 Win=2092 Len=0
10.0.2.15	128.138.140.44	54 TCP	60093 → 13 [FIN, ACK] Seq=1 Ack=53 Win=2092 Len=0
128.138.140.44	10.0.2.15	60 TCP	13 → 60093 [ACK] Seq=53 Ack=2 Win=65535 Len=0

5、TCP和UDP的Echo Server

在实现Socket API的过程中，我们需要简单易行的方式来测试我们的实现是否正确，这里我们使用一个简单的Echo Server来测试我们的实现。

我们实现的TCP和UDP的Echo Server中使用TCP的代码如代码 3.31所示。是一个单进程的简易Echo Server，在用于测试我们的Socket API的实现是否正确的用途中发挥了重要的作用。

代码 3.31: TCP Echo Server

```

1  #include "types.h"
2  #include "user.h"
3  #include "socket.h"
4  #include "inet.h"
5
6  int
7  main(int argc, char *argv[])
8  {
9      int soc, acc, peerlen, ret;
10     struct sockaddr_in self, peer;
11     unsigned char *addr;
12     char buf[2048];
13
14     printf(1, "Starting TCP Echo Server\n");
15     soc = socket(PF_INET, SOCK_STREAM, IPPROTO_TCP);
16     if (soc == 1)
17     {
18         printf(1, "socket: failure\n");
19         exit(soc);
20     }
21     printf(1, "socket: success, soc=%d\n", soc);
22     self.sin_family = AF_INET;
23     self.sin_addr.s_addr = INADDR_ANY;
24     self.sin_port = htons(7);
25     if (bind(soc, (struct sockaddr *)&self, sizeof(self)) == -1)
26     {
27         printf(1, "bind: failure\n");
28         close(soc);
29         exit(soc);
30     }
31     addr = (unsigned char *)&self.sin_addr;
32     printf(1, "bind: success, self=%d.%d.%d.%d\n", addr[3], addr[2], addr[1], addr[0], ntohs(
        (self.sin_port)));
33     listen(soc, 100);
34     printf(1, "waiting for connection...\n");
35     peerlen = sizeof(peer);
36     acc = accept(soc, (struct sockaddr *)&peer, &peerlen);
37     if (acc == -1)
38     {

```

```

39         printf(1, "accept: failure\n");
40         close(soc);
41         exit(soc);
42     }
43     addr = (unsigned char *)&peer.sin_addr;
44     printf(1, "accept: success, peer=%d.%d.%d.%d\n", addr[3], addr[2], addr[1], addr[0],
        ntohs(peer.sin_port));
45     while (1)
46     {
47         memset(buf, 0, sizeof(buf));
48         ret = recv(acc, buf, sizeof(buf));
49         if (ret <= 0)
50         {
51             printf(1, "EOF\n");
52             break;
53         }
54         printf(1, "recv: %d bytes data received\n", ret);
55         hexdump(buf, ret);
56
57         if (send(acc, buf, ret) != ret)
58         {
59             printf(1, "send: failure\n");
60             break;
61         }
62     }
63     close(soc);
64     return 0;
65 }

```

在QEMU中的xv6操作系统终端中执行tcpecho命令运行TCP Echo Server，结果如图 3.14所示。可

图 3.14: TCP Echo Server

```

Machine View
init: starting sh
$ tcpecho
Starting TCP Echo Server
socket: success, soc=3
bind: success, self=0.0.0.0:1792
waiting for connection...
accept: success, peer=10.0.2.2:42774
recv: 4 bytes data received
+-----+-----+
| 0000 | 61 61 61 0A | aaa. |
+-----+-----+
recv: 5 bytes data received
+-----+-----+
| 0000 | 66 66 66 66 0A | ffff. |
+-----+-----+
recv: 4 bytes data received
+-----+-----+
| 0000 | 64 64 64 0A | ddd. |
+-----+-----+
recv: 2 bytes data received
+-----+-----+
| 0000 | 61 0A | a. |
+-----+-----+

```

见，从宿主机终端中运行的nc命令向TCP Echo Server发送了的数据都被TCP Echo Server将数据原样返回给客户端。同时TCP Echo Server也将收到的数据打印到了终端中。包括数据包的基本信息。由捕捉的数据也可看出这一来往的过程。如图 3.15所示。

图 3.15: TCP Echo Server和宿主机间往来的数据包

16	16.131967	10.0.2.2	10.0.2.15	60	TCP	42774 → 7 [SYN] Seq=0 Win=65535 Len=0 MSS=1460
17	16.139309	10.0.2.15	10.0.2.2	58	TCP	7 → 42774 [SYN, ACK] Seq=0 Ack=1 Win=2144 Len=0 MSS=536
18	16.139330	10.0.2.2	10.0.2.15	60	TCP	42774 → 7 [ACK] Seq=1 Ack=1 Win=65535 Len=0
19	20.853477	10.0.2.2	10.0.2.15	60	ECHO	Request
20	20.974376	10.0.2.15	10.0.2.2	58	ECHO	Response
21	20.977553	10.0.2.2	10.0.2.15	60	TCP	42774 → 7 [ACK] Seq=5 Ack=5 Win=65535 Len=0
22	22.916356	10.0.2.2	10.0.2.15	60	ECHO	Request
23	23.036487	10.0.2.15	10.0.2.2	59	ECHO	Response
24	23.040474	10.0.2.2	10.0.2.15	60	TCP	42774 → 7 [ACK] Seq=10 Ack=10 Win=65535 Len=0
25	24.130993	10.0.2.2	10.0.2.15	60	ECHO	Request
26	24.252652	10.0.2.15	10.0.2.2	58	ECHO	Response
27	24.254399	10.0.2.2	10.0.2.15	60	TCP	42774 → 7 [ACK] Seq=14 Ack=14 Win=65535 Len=0
28	24.902418	10.0.2.2	10.0.2.15	60	ECHO	Request
29	25.038877	10.0.2.15	10.0.2.2	56	ECHO	Response
30	25.042052	10.0.2.2	10.0.2.15	60	TCP	42774 → 7 [ACK] Seq=16 Ack=16 Win=65535 Len=0

第四章 项目管理

第五章 项目感想

一、熊恪峥

在这次计算机网络课程项目中，我负责了项目的设计、分工、全部内核态部分的实现和少数用户态程序的实现。在这个过程中，我一方面增强了对xv6内核代码的理解，另一方面也学习了现有成熟的网络API是怎么一步一步从0开始构建起来。从驱动程序如何将数据从网卡中读出、从用户给定的位置送给网卡，到内核如何处理不同的网络协议，我自己编码一步步将它们实现出来，这是一次非常有意义的经历。也让我对我们习以为常的“开箱即用”的API接口的底层做什么、怎么做的问题有了更深的理解。

此外，由于我们项目中的代码量级极大，为了有效地管理，我们使用了git进行版本控制，为了处理协议、Socket数据报类型等复杂的组合、灵活地复用代码，我不仅调查和研究了现有操作系统内核相应的处理方式，也加入了自己在实践中积累经验得出的设计，这也让我有很大的收获。

二、赵淇奥

在这次计算机网络课程设计项目中，我负责编写了nslookup程序。这个程序的功能是向dns服务器查询给定的地址或域名。在之前计网课程的实验中我多次接触过linux内置的该命令，而通过本次项目，亲自实现该命令，让我更加了解计算机网络编程、熟悉计算机网络内核协议栈，并且对DNS的工作原理有了更深刻的认识。

三、肖凯欣

在项目中我负责编写ping程序和traceroute程序。在计算机网络的理论课上，这两个程序都作为网络层协议应用的示例。在程序实现中需要我根据程序需要生成、发送相应的ICMP报文，并对收到的ICMP报文进行解析和处理。在过程中我对于计算机网络编程，网络层协议以及ICMP协议的工作原理有了更加深刻的认识。

四、谢哲涵

在项目中我负责编写httpget程序。在计算机网络的实验课中，我们对使用HTTP连接获取互联网上的内容有了初步的理解和认识，而在实现这个命令过程中我对于计算机网络编程，应用层协议有了更加深刻的认识。

附录：代码清单

代码 5.1: DNS报文构造

```
1 static inline struct netdb_answer *make_query(char *name, int type)
2 {
3     assert_holding(&netdb_lock);
4
5     char *p = netdb_qbuf + 2048;
6
7     dns_header_t *header = (dns_header_t *)p;
8     header->xid = byteswap16(0x1204); // should be random though
9     header->flags = byteswap16(0x0100); // Q, RD
10    header->qdcount = byteswap16(1); // 1 question
11
12    p += sizeof(dns_header_t);
13    dns_question_t *question = (dns_question_t *)p;
14    question->dnstype = byteswap16(type); // A or PTR
15    question->dnsclass = byteswap16(1); // IN
16
17    p += sizeof(dns_question_t);
18    question->name = p;
19
20    const uint8 hostlen = (uint8)strlen(name);
21    memmove(question->name + 1, name, hostlen + 1);
22
23    uint8 *prev = (uint8 *)question->name;
24    uint8 count = 0;
25
26    for (uint8 i = 0; i < hostlen; i++)
27    {
28        if (name[i] == '.')
29        {
30            *prev = count;
31            prev = (uint8 *) (question->name + i + 1);
32            count = 0;
33        }
34        else
35        {
36            count++;
37        }
38    }
39    *prev = count;
40
41    p = netdb_qbuf;
42    size_t packet_len = sizeof(dns_header_t) + sizeof(dns_question_t) - sizeof(char *) + hostlen
43        + 2;
44
45    memmove(p, header, sizeof(dns_header_t));
46    p += sizeof(dns_header_t);
47
48    memmove(p, question->name, hostlen + 1);
49    p += hostlen + 2;
50
51    memmove(p, &question->dnstype, sizeof(question->dnstype));
52    p += sizeof(question->dnstype);
```

```

52     memmove(p, &question->dnsclass, sizeof(question->dnsclass));
53
54     netbegin_op();
55     struct pbuf *pbuf = pbuf_alloc(PBUF_TRANSPORT, packet_len, PBUF_RAM);
56     if (!pbuf)
57     {
58         netend_op();
59         return NULL;
60     }
61     memmove(pbuf->payload, netdb_qbuf, packet_len);
62     query_ans = NULL;
63     udp_rcv(dns_pcb, netdb_dns_rcv, NULL);
64     err_t err = udp_send(dns_pcb, pbuf);
65
66     if (err != ERR_OK)
67     {
68         netend_op();
69         pbuf_free(pbuf);
70         return NULL;
71     }
72     netend_op();
73
74     if (!query_ans)
75     {
76         if (sleepddl(&dns_pcb, &netdb_lock, DNS_TIMEOUT) == 0)
77         {
78             return NULL;
79         }
80     }
81     return query_ans;
82 }

```

代码 5.2: Socket API的顶层封装

```

1  int socketrecvfrom(socket_t *skt, char *buf, int len, int flags, struct sockaddr *addr, int *addrlen
   )
2  {
3      if (skt->type == SOCK_STREAM)
4      {
5          if (addr)
6          {
7              return -EISCONN;
8          }
9
10         // as stated in the document, it will have the same effect as rcv()
11         return socketrecv(skt, buf, len, flags);
12     }
13
14     if (!skt->opts->recvfrom)
15     {
16         return -ESOCKTNOSUPPORT;
17     }
18
19     skt->recvfrom_params.recv_len = len;
20

```

```
21     int err = skt->opts->recvfrom(skt, buf, len, flags, addr, addrlen);
22     if (err)
23     {
24         return err;
25     }
26
27     if (skt->recv_buf == NULL)
28     {
29         acquire(&skt->lock);
30         skt->wakeup_retcode = 0;
31         if (sleepddl(&skt->recv_chan, &skt->lock, skt->recv_timeout) == 0)
32         {
33             release(&skt->lock);
34             return -EWOULDBLOCK;
35         }
36         release(&skt->lock);
37
38         if (skt->wakeup_retcode != 0)
39         {
40             return skt->wakeup_retcode;
41         }
42     }
43
44     netbegin_op();
45
46     len = pbuf_copy_partial(skt->recv_buf, buf, len, skt->recv_offset);
47     skt->recv_offset += len;
48     KDEBUG_ASSERT(skt->recv_offset <= skt->recv_buf->tot_len);
49
50     if (skt->recv_offset == skt->recv_buf->tot_len)
51     {
52         pbuf_free(skt->recv_buf);
53         skt->recv_buf = NULL;
54         skt->recv_offset = 0;
55
56         memmove(addr, &skt->recvfrom_params.recv_addr, sizeof(struct sockaddr_in));
57         *addrlen = sizeof(struct sockaddr_in);
58         addrin_byteswap((struct sockaddr_in *)addr);
59
60         memset(&skt->recvfrom_params, 0, sizeof(skt->recvfrom_params));
61     }
62
63     netend_op();
64
65     return len;
66 }
```