

# 第31章 BPF : BSD 分组过滤程序

## 31.1 引言

BSD分组过滤程序(BPF)是一种软件设备，用于过滤网络接口的数据流，即给网络接口加上“开关”。应用进程打开 /dev/bpf0、/dev/bpf1等等后，可以读取BPF设备。每个应用进程一次只能打开一个BPF设备。

因为每个BPF设备需要8192字节的缓存，系统管理员一般限制 BPF设备的数目。

如果open返回EBUSY，说明该设备已被使用，应用进程应该试着打开下一 BPF设备，直到open成功为止。

通过若干ioctl命令，可以配置 BPF设备，把它与某个网络接口相关连，并安装过滤程序，从而能够选择性地接收输入的分组。 BPF设备打开后，应用进程通过读写设备来接收分组，或将分组放入网络接口队列中。

我们将一直使用“分组”，尽管“帧”可能更准确一些，因为 BPF工作在数据链路层，在发送和接收的数据帧中包含了链路层的首部。

BPF设备工作的前提是网络接口必须能够支持 BPF。第3章中提到以太网、SLIP和环回接口的驱动程序都调用了 bpfattach，用于配置读取 BPF设备的接口。本节中，我们将介绍 BPF设备驱动程序是如何组织的，以及数据分组在驱动程序和网络接口之间是如何传递的。

BPF一般情况下用作诊断工具，查看某个本地网络上的流量，卷 1附录A 介绍的tcpdump程序是此类工具中最好的一个。通常情况下，用户感兴趣的是一组指定主机间交互的分组，或者某个特定协议，甚至某个特定 TCP连接上的数据流。BPF设备经过适当配置，能够根据过滤程序的定义丢弃或接受输入的分组。过滤程序的定义类似于伪机器指令， BPF的细节超出了本书的讨论范围，感兴趣的读者请参阅 bpf(4)和[McCanne and Jacobson 1993]。

## 31.2 代码介绍

下面将要介绍的有关 BPF设备驱动程序的代码，包括两个头文件和一个 C文件，在图31-1中给出。

文 件	描 述
net/bpf.h	BPF常量
net/bpfdesc.h	BPF结构
net/bpf.c	BPF设备支持

图31-1 本章讨论的文件

### 31.2.1 全局变量

本章用到的全局变量在图 31-2中给出。

变 量	数 据 类 型	描 述
bpf_iflist	struct bpf_if *	支持BPF的接口组成的链表
bpf_dtab	struct bpf_d []	BPF描述符数组
bpf_bufsize	int	BPF缓存大小默认值

图31-2 本章用到的全局变量

### 31.2.2 统计量

图31-3列出了bpf\_d结构中为每个活动的BPF设备维护的两个统计量。

bpf_d成员变量	描 述
bd_rcount	从网络接口接收的分组的数目
bd_dcount	由于缓存空间不足而丢弃的分组的数目

图31-3 本章讨论的统计值

本章的其余内容分为4个部分：

- BPF接口结构；
- BPF设备描述符；
- BPF输入处理；和
- BPF输出处理。

### 31.3 bpf\_if结构

BPF维护一个链表，包括所有支持BPF的网络接口。每个接口都由一个bpf\_if结构描述，全局指针bpf\_iflist指向表中的第一个结构。图31-4给出了BPF接口结构。

```
67 struct bpf_if {  
68     struct bpf_if *bif_next;    /* list of all interfaces */  
69     struct bpf_d *bif_dlist;    /* descriptor list */  
70     struct bpf_if **bif_driverp; /* pointer into softc */  
71     u_int    bif_dlt;           /* link layer type */  
72     u_int    bif_hdrlen;        /* length of header (with padding) */  
73     struct ifnet *bif_ifp;      /* correspodng interface */  
74 };  
bpfdesc.h
```

图31-4 bpf\_if 结构

67-79 bif\_next指向链表中的下一个BPF接口结构。bif\_dlist指向另一个链表，包括所有已打开并配置过的BPF设备。

70 如果某个网络接口已配置了BPF设备，即被加上了开关，则bif\_driverp将指向ifnet结构中的bpf\_if指针。如果网络接口还未加上开关，\*bif\_driverp为空。为某个网络接口配置BPF设备时，\*bif\_driverp将指向bif\_if结构，从而告诉接口可以开始向BPF传递分组。

71 接口类型保存在bif\_dlt中。图31-5中列出了前面提到的几个接口所分别对应的常量值。

bif_dlt	描 述
<i>DLT_EN10MB</i>	10 Mb以太网接口
<i>DLT_SLIP</i>	SLIP接口
<i>DLT_NULL</i>	环回接口

图31-5 bif\_dlt 值

72-74 BPF接受的所有分组都有一个附加的 BPF首部。bif\_hdrlen等于首部大小。最后，bif\_ifp指向对应接口的ifnet结构。

图31-6给出了每个输入分组中附加的 bpf\_hdr结构。

```

122 struct bpf_hdr {
123     struct timeval bh_tstamp; /* time stamp */
124     u_long bh_caplen; /* length of captured portion */
125     u_long bh_datalen; /* original length of packet */
126     u_short bh_hdrlen; /* length of bpf header (this struct plus
127                        alignment padding) */
128 };

```

bpf.h

图31-6 bpf\_hdr 结构

122-128 bh\_tstamp记录了分组被捕捉的时间。bh\_caplen等于BPF保存的字节数，bh\_datalen等于原始分组中的字节数。bh\_headlen等于bpf\_hdr的大小加上所需填充字节的长度。它用于解释从BPF设备中读取的分组，应该等同于接收接口的 bif\_hdrlen。

图31-7给出了bpf\_if结构是如何与前述3个接口(le\_softc[0]、sl\_softc[0]和loif)的ifnet结构建立连接的。

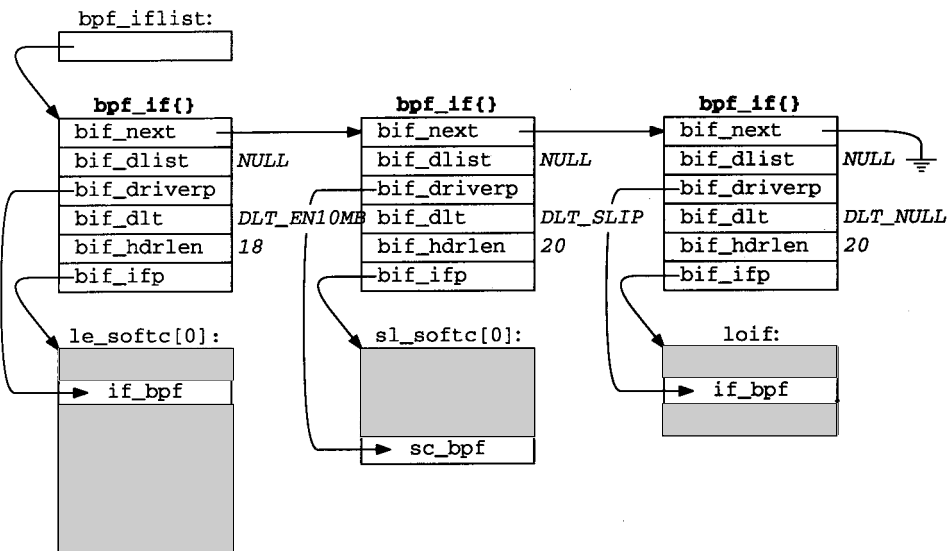


图31-7 bpf\_if 和ifnet 结构

注意，bif\_driverp指向网络接口的if\_bpf和sc\_bpf指针，而不是接口结构。

SLIP设备使用sc\_bpf，而不是if\_bpf。这可能是因为在SLIP BPF代码完成时，

if\_bpf成员变量还未加入到ifnet结构中。Net/2中的ifnet结构不包括if\_bpf成员。

按照各接口驱动程序调用 bpfattach时给出的信息，对3个接口初始化链路类型和首部长度成员变量。

第3章介绍了bpfattach被以太网、SLIP和环回接口的驱动程序调用。每个设备驱动程序初始化调用bpfattach时，将构建BPF接口结构链表。图31-8给出了该函数。

```

1053 void
1054 bpfattach(driverp, ifp, dlt, hdrlen)
1055 caddr_t *driverp;
1056 struct ifnet *ifp;
1057 u_int dlt, hdrlen;
1058 {
1059     struct bpf_if *bp;
1060     int i;
1061     bp = (struct bpf_if *) malloc(sizeof(*bp), M_DEVBUF, M_DONTWAIT);
1062     if (bp == 0)
1063         panic("bpfattach");
1064     bp->bif_dlist = 0;
1065     bp->bif_driverp = (struct bpf_if **) driverp;
1066     bp->bif_ifp = ifp;
1067     bp->bif_dlt = dlt;
1068     bp->bif_next = bpf_iflist;
1069     bpf_iflist = bp;
1070     *bp->bif_driverp = 0;
1071     /*
1072     * Compute the length of the bpf header. This is not necessarily
1073     * equal to SIZEOF_BPF_HDR because we want to insert spacing such
1074     * that the network layer header begins on a longword boundary (for
1075     * performance reasons and to alleviate alignment restrictions).
1076     */
1077     bp->bif_hdrlen = BPF_WORDALIGN(hdrlen + SIZEOF_BPF_HDR) - hdrlen;
1078     /*
1079     * Mark all the descriptors free if this hasn't been done.
1080     */
1081     if (!D_ISFREE(&bpf_dtab[0]))
1082         for (i = 0; i < NBPFILTER; ++i)
1083             D_MARKFREE(&bpf_dtab[i]);
1084     printf("bpf: %s%d attached\n", ifp->if_name, ifp->if_unit);
1085 }

```

图31-8 bpfattach 函数

1053-1063 每个支持BPF的设备驱动程序都将调用 bpfattach。第一个参数是保存在bif\_driverp的指针(图31-4给出)，第二个参数指向接口的ifnet结构，第三个参数确认数据链路层类型，第四个参数传递分组中的数据链路首部大小，为接口分配一个新的 bpf\_if结构。

#### 1. 初始化bpf\_if结构

1064-1070 `bpf_if` 结构根据函数的参数进行初始化, 并插入到 BPF接口链表, `bpf_iflist`, 的表头。

## 2. 计算BPF首部大小

1071-1077 设定 `bif_hdrlen` 大小, 强迫网络层首部 (如IP首部) 从一个长字的边界开始。这样可以提高性能, 避免为 BPF 加入不必要的对齐限制。图 31-9列出了在前述3个接口上, 各自捕捉到的BPF分组的总体结构。

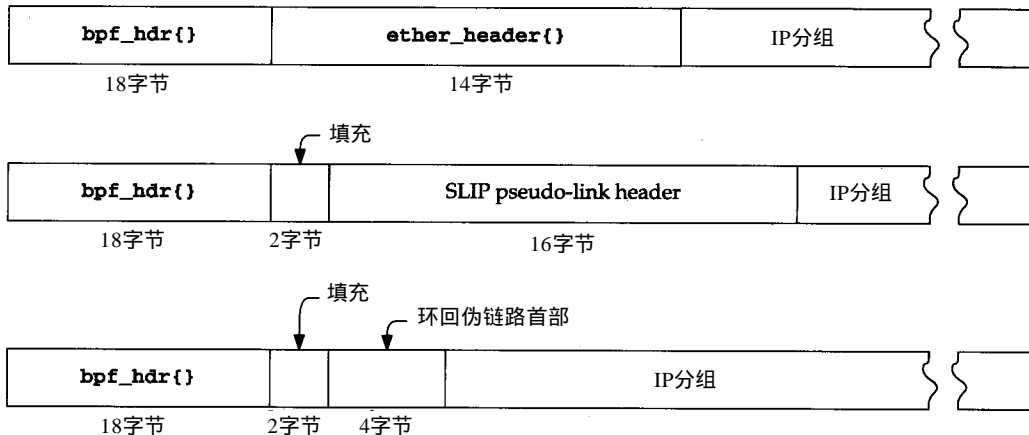


图31-9 BPF分组结构

`ether_header` 结构在图 4-10 中给出, `SLIP` 伪链路首部在图 5-14 中给出, 而环回接口伪链路首部在图 5-28 中给出。

请注意, `SLIP` 和环回接口分组需要填充 2 字节, 以强迫 IP 首部按 4 字节对齐。

## 3. 初始化 `bpf_dtab` 表

1078-1083 代码初始化图 31-10 中给出的 BPF 描述符表。注意, 仅在第一次调用 `bpfattach` 时进行初始化, 后续调用将跳过初始化过程。

## 4. 打印控制台信息

1084-1085 系统向控制台输出一条短信息, 宣告接口已配置完毕, 可以支持 BPF。

## 31.4 `bpf_d` 结构

为了能够选择性地接收输入报文, 应用进程首先打开一个 BPF 设备, 调用若干 `ioctl` 命令规定 BPF 过滤程序的条件, 指明接口、读缓存大小和超时时限。每个 BPF 设备都有一个相关的 `bpf_d` 结构, 如图 31-10 所示。

45-46 如果同一网络接口上配置了多个 BPF 设备, 与之相应的 `bpf_d` 结构将组成一个链表。`bd_next` 指向链表中的下一个结构。

### 分组缓存

47-52 每个 `bpf_d` 结构都有两个分组缓存。输入分组通常保存在 `bd_sbuf` 所对应的缓存 (存储缓存) 中。另一个缓存要么对应于 `bd_fbuf` (空闲缓存), 意味着缓存为空; 或者对应于 `bd_hbuf` (暂留缓存), 意味着缓存中有分组等待应用进程读取。`bd_slens` 和 `bd_hlens` 分别记

录了保存在存储缓存和暂留缓存中的字节数。

```

45 struct bpf_d {
46     struct bpf_d *bd_next;        /* Linked list of descriptors */
47     caddr_t bd_sbuf;              /* store slot */
48     caddr_t bd_hbuf;              /* hold slot */
49     caddr_t bd_fbuf;              /* free slot */
50     int     bd_slenn;              /* current length of store buffer */
51     int     bd_hlen;              /* current length of hold buffer */
52     int     bd_bufsize;           /* absolute length of buffers */
53     struct bpf_if *bd_bif;        /* interface descriptor */
54     u_long  bd_rtout;             /* Read timeout in 'ticks' */
55     struct bpf_insn *bd_filter;   /* filter code */
56     u_long  bd_rcount;            /* number of packets received */
57     u_long  bd_dcount;            /* number of packets dropped */
58     u_char  bd_promisc;           /* true if listening promiscuously */
59     u_char  bd_state;             /* idle, waiting, or timed out */
60     u_char  bd_immediate;         /* true to return on packet arrival */
61     u_char  bd_pad;              /* explicit alignment */
62     struct selinfo bd_sel;        /* bsd select info */
63 };

```

bpfdesc.h

bpfdesc.h

图31-10 bpf\_d 结构

如果存储缓存已满，它将被连接到 bd\_hbuf，而空闲缓存将被连接到 bd\_sbuf。当暂留缓存清空时，它会被连接到 bd\_fbuf。宏 ROTATE\_BUFFERS 负责把存储缓存连接到 bd\_hbuf，空闲缓存连接到 bd\_sbuf，并清空 bd\_fbuf。存储缓存满或者应用进程不想再等待更多的分组时调用该宏。

bd\_bufsize 记录与设备相连的两个缓存的大小，其默认值等于 4096 (BPF\_BUFSIZE) 字节。修改内核代码可以改变默认值大小，或者通过 BIOCSBLEN 命令改变某个特定 BPF 设备的 bd\_bufsize。BIOCGBLEN 命令返回 bd\_bufsize 的当前值，其最大值不超过 32768 (BPF\_MAXBUFSIZE) 字节，最小值为 32 (BPF\_MINBUFSIZE) 字节。

53-57 bd\_bif 指向 BPF 设备所对应的 bpf\_if 结构。BIOCSETIF 命令可指明设备。bd\_rtout 是等待分组时，延迟的滴答数。bd\_filter 指向 BPF 设备的过滤程序代码。两个统计值，应用进程可通过 BIOCGSTATS 命令读取，分别保存在 bd\_rcount 和 bd\_dcount 中。

58-63 bd\_promisc 通过 BIOCPROMISC 命令置位，从而使接口工作在混杂 (promiscuous) 状态。bd\_state 未使用。bd\_immediate 通过 BIOCIMMEDIATE 命令置位，促使驱动程序收到分组后即返回，不再等待暂留缓存填满。bd\_pad 填充 bpf\_d 结构，从而与长字边界对齐。bd\_sel 保存的 selinfo 结构，可用于 select 系统调用。我们不准备介绍如何对 BPF 设备使用 select 系统调用，16.13 节已介绍了 select 的一般用法。

### 31.4.1 bpfopen 函数

应用进程调用 open，试图打开一个 BPF 设备时，该调用将被转到 bpfopen (图 31-11)。

256-263 系统编译时，BPF 设备的数目受到 NBPFILTER 的限制。如果设备的最小设备号大

于NBPFILTER，则返回 ENXIO，这是因为系统管理员创建的 /dev/bpf<sub>x</sub>项数大于NBPFILTER的值。

```

256 int
257 bpfopen(dev, flag)
258 dev_t dev;
259 int flag;
260 {
261     struct bpf_d *d;
262     if (minor(dev) >= NBPFILTER)
263         return (ENXIO);
264     /*
265      * Each minor can be opened by only one process. If the requested
266      * minor is in use, return EBUSY.
267      */
268     d = &bpf_dtab[minor(dev)];
269     if (!D_ISFREE(d))
270         return (EBUSY);
271     /* Mark "free" and do most initialization. */
272     bzero((char *) d, sizeof(*d));
273     d->bd_bufsize = bpf_bufsize;
274     return (0);
275 }

```

bpf.c

bpf.c

图31-11 bpfopen 函数

#### 分配bpf\_d结构

264-275 同一时间内，一个应用进程只能访问一个 BPF设备。如果bpf\_d结构已被激活，则返回EBUSY。应用程序，如tcpdump，收到此返回值时，会自动寻找下一个设备。如果该设备已存在，最小设备号所指定的bpf\_dtab表中的项被清除，分组缓存大小复位为默认值。

#### 31.4.2 bpfiocctl函数

设备打开后，可通过iocctl命令进行配置。图31-12总结了与BPF设备有关的iocctl命令。图31-13给出了bpfiocctl函数，只列出BIOCSETF和BIOCSETIF的处理代码，其他未涉及到的iocctl命令则被忽略。

命 令	第三个参数	函 数	描 述
FIONREAD	u_int	bpfiocctl	返回暂留缓存和存储缓存中的字节数
BIOCGBLEN	u_int	bpfiocctl	返回分组缓存大小
BIOCSBLEN	u_int	bpfiocctl	设定分组缓存大小
BIOCSETF	struct bpf_program	bpf_setf	安装BPF程序
BIOCFLUSH		reset_d	丢弃挂起分组
BIOCPROMISC		ifpromisc	设定混杂方式
BIOCGDLT	u_int	bpfiocctl	返回bif_dlt
BIOCGETIF	struct ifreq	bpf_ifname	返回所属接口的名称
BIOCSETIF	struct ifreq	bpf_setif	为网络接口添加设备
BIOCSRTIMEOUT	struct timeval	bpfiocctl	设定“读”操作的超时时限
BIOCGRTIMEOUT	struct timeval	bpfiocctl	返回“读”操作的超时时限
BIOCGSTATS	struct bpf_stat	bpfiocctl	返回BPF统计值
BIOCIMMEDIATE	u_int	bpfiocctl	设定立即方式
BIOCVERSION	struct bpf_version	bpfiocctl	返回BPF版本信息

图31-12 BPF iocctl 命令

```

501 int
502 bpfioctl(dev, cmd, addr, flag)
503 dev_t dev;
504 int cmd;
505 caddr_t addr;
506 int flag;
507 {
508     struct bpf_d *d = &bpf_dtab[minor(dev)];
509     int s, error = 0;
510     switch (cmd) {
511         /*
512          * Set link layer read filter.
513          */
514         case BIOCSETF:
515             error = bpf_setf(d, (struct bpf_program *) addr);
516             break;
517         /*
518          * Set interface.
519          */
520         case BIOCSETIF:
521             error = bpf_setif(d, (struct ifreq *) addr);
522             break;
523
524         /* other ioctl commands from Figure 31.12 */
525
526         default:
527             error = EINVAL;
528             break;
529     }
530     return (error);
531 }

```

图31-13 bpfioctl 函数

501-509 与bpfopen类似，通过最小设备号从bpf\_dtab表中选取相应的bpf\_d结构。整个命令处理是一个大的switch/case语句。我们给出了两个命令，BIOCSETF和BIOCSETIF，以及default子句。

510-522 bpf\_setf函数安装由addr指向的过滤程序，bpf\_setif建立起指定名称接口与bpf\_d结构间的对应关系。本书中没有给出bpf\_setf的实现代码。

668-673 如果命令未知，则返回EINVAL。

图31-14的例子中，bpf\_setif已把bpf\_d结构连接到LANCE接口上。

图中，bif\_dlist指向bpf\_dtab[0]，以太网接口描述符链表中的第一个也是仅有的一个描述符。在bpf\_dtab[0]中，bd\_sbuf和bd\_hbuf成员分别指向存储缓存和暂留缓存。两个缓存大小都等于4096(bd\_bufsize)字节。bd\_bif回指接口的bpf\_if结构。

ifnet结构(le\_softc[0])中的if\_bpf也指回bpf\_if结构。如图4-19和图4-11所示，如果if\_bpf非空，则驱动程序开始调用bpf\_tap，向BPF设备传递分组。

图31-15接着图31-10，给出了打开第二个BPF设备，并连接到同一个以太网网络接口后的各结构变量的状态。



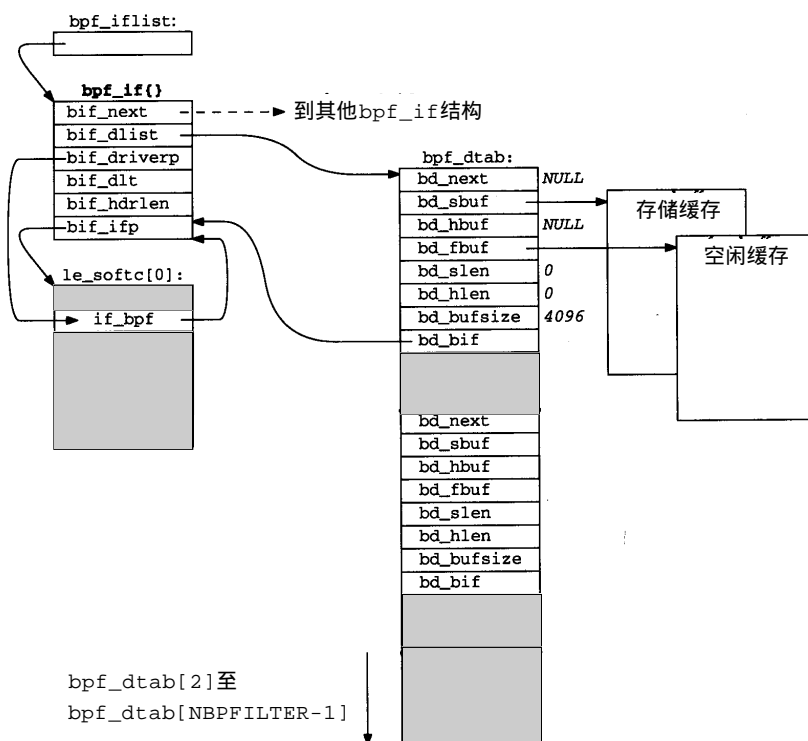


图31-14 连接到以太网接口的BPF设备

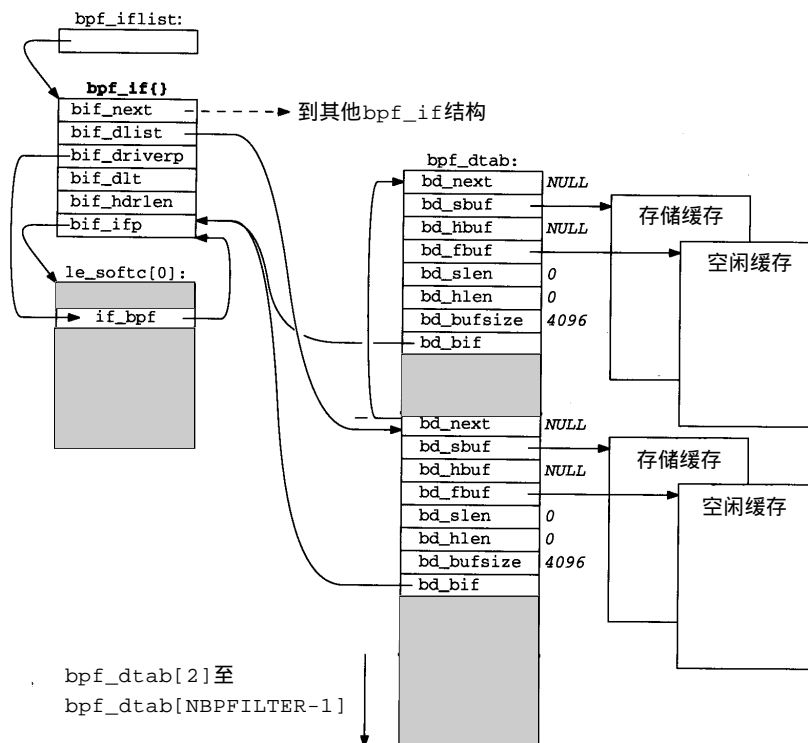


图31-15 连接到以太网接口的两个BPF设备

第二个BPF设备打开时，在 `bpf_dtab`表中分配一个新的 `bpf_d`结构，本例中为 `bpf_dtab[1]`。因为第二个BPF设备也连接到同一个以太网接口，`bif_dlist`指向 `bpf_dtab[1]`，并且 `bpf_dtab[1].bd_next`指向 `bpf_dtab[0]`，即以太网上对应的第一个BPF描述符。系统为新的描述符结构分别分配存储缓存和暂留缓存。

### 31.4.3 `bpf_setif`函数

`bpf_setif`函数，负责建立BPF描述符与网络接口间的连接，如图31-16所示。

```

721 static int
722 bpf_setif(d, ifr)
723 struct bpf_d *d;
724 struct ifreq *ifr;
725 {
726     struct bpf_if *bp;
727     char *cp;
728     int unit, s, error;

729     /*
730      * Separate string into name part and unit number. Put a null
731      * byte at the end of the name part, and compute the number.
732      * If the a unit number is unspecified, the default is 0,
733      * as initialized above. XXX This should be common code.
734      */
735     unit = 0;
736     cp = ifr->ifr_name;
737     cp[sizeof(ifr->ifr_name) - 1] = '\0';
738     while (*cp++) {
739         if (*cp >= '0' && *cp <= '9') {
740             unit = *cp - '0';
741             *cp++ = '\0';
742             while (*cp)
743                 unit = 10 * unit + *cp++ - '0';
744             break;
745         }
746     }
747     /*
748      * Look through attached interfaces for the named one.
749      */
750     for (bp = bpf_iflist; bp != 0; bp = bp->bif_next) {
751         struct ifnet *ifp = bp->bif_ifp;

752         if (ifp == 0 || unit != ifp->if_unit
753             || strcmp(ifp->if_name, ifr->ifr_name) != 0)
754             continue;
755         /*
756          * We found the requested interface.
757          * If it's not up, return an error.
758          * Allocate the packet buffers if we need to.
759          * If we're already attached to requested interface,
760          * just flush the buffer.
761          */
762         if ((ifp->if_flags & IFF_UP) == 0)
763             return (ENETDOWN);
    
```

图31-16 `bpf_setif` 函数

```

764         if (d->bd_sbuf == 0) {
765             error = bpf_alloccbufs(d);
766             if (error != 0)
767                 return (error);
768         }
769         s = splimp();
770         if (bp != d->bd_bif) {
771             if (d->bd_bif)
772                 /*
773                  * Detach if attached to something else.
774                  */
775                 bpf_detachd(d);
776             bpf_attachd(d, bp);
777         }
778         reset_d(d);
779         splx(s);
780         return (0);
781     }
782     /* Not found. */
783     return (ENXIO);
784 }

```

bpf.c

图31-16 (续)

721-746 bpf\_setif的第一部分完成ifreq结构(图4-23)中接口名的正文与数字部分的分离, 数字部分保存在unit中。例如, 如果ifr\_name的头4字节为“sl1\0”, 代码执行完毕后, 将等于“sl\0\0”, 且unit等于1。

#### 1. 寻找匹配的ifnet结构

747-754 for循环用于在支持BPF的接口(bpf\_iflist中)中查找符合ifreq定义的接口。

755-768 如果未找到匹配的接口, 则返回 ENETDOWN。如果接口存在, bpf\_allocate为bpf\_d分配空闲缓存和存储缓存, 如果它们还未被分配的话。

#### 2. 连接bpf\_d结构

769-777 如果BPF设备还未与网络接口建立连接关系, 或者连接的网络接口不是 ifreq中指定的接口, 则调用bpf\_detachd丢弃原先的接口(如果存在), 并调用bpf\_attachd将其连接到新的接口上。

778-784 reset\_d复位分组缓存, 丢弃所有在应用进程中等待的分组。函数返回 0, 说明处理成功; 或者ENXIO, 说明未找到指定接口。

### 31.4.4 bpf\_attachd函数

图31-17给出的bpf\_attachd函数, 建立起BPF描述符与BPF设备和网络接口间的对应关系。

bpf.c

```

189 static void
190 bpf_attachd(d, bp)
191 struct bpf_d *d;
192 struct bpf_if *bp;
193 {
194     /*
195      * Point d at bp, and add d to the interface's list of listeners.

```

图31-17 bpf\_attachd 函数

```

196      * Finally, point the driver's bpf cookie at the interface so
197      * it will divert packets to bpf.
198      */
199      d->bd_bif = bp;
200      d->bd_next = bp->bif_dlist;
201      bp->bif_dlist = d;

202      *bp->bif_driverp = bp;
203 }

```

bpf.c

图31-17 (续)

189-203 首先，令bd\_bif指向网络接口的BPF接口结构。接着，bpf\_d结构被插入到与设备对应的bpf\_d结构链表的头部。最后，改变网络接口中的BPF指针，指向当前BPF结构，从而促使接口向BPF设备传递分组。

## 31.5 BPF的输入

一旦BPF设备打开并配置完毕，应用进程就通过 read系统调用从接口中接收分组。BPF过滤程序复制输入分组，因此，不会干扰正常的网络处理。输入分组保存在与BPF设备相连的存储缓存和暂留缓存中。

### 31.5.1 bpf\_tap函数

下面列出了图4-11中LANCE设备驱动程序调用 bpf\_tap的代码，并利用这一调用介绍 bpf\_tap函数。图4-11中的调用如下：

```
bpf_tap(le->sc_if.if_bpf, buf, len + sizeof(struct ether_header));
```

图31-18给出了bpf\_tap函数。

```

869 void
870 bpf_tap(arg, pkt, pktlen)
871 caddr_t arg;
872 u_char *pkt;
873 u_int  pktlen;
874 {
875     struct bpf_if *bp;
876     struct bpf_d *d;
877     u_int  slen;
878     /*
879      * Note that the ipl does not have to be raised at this point.
880      * The only problem that could arise here is that if two different
881      * interfaces shared any data. This is not the case.
882      */
883     bp = (struct bpf_if *) arg;
884     for (d = bp->bif_dlist; d != 0; d = d->bd_next) {
885         ++d->bd_rcount;
886         slen = bpf_filter(d->bd_filter, pkt, pktlen, pktlen);
887         if (slen != 0)
888             catchpacket(d, pkt, pktlen, slen, bcopy);
889     }
890 }

```

bpf.c

bpf.c

图31-18 bpf\_tap 函数

869-882 第一个参数是指向 bpf\_if 结构的指针, 由 bpfattach 设定。第二个参数是指向进入分组的指针, 包括以太网首部。第三个参数等于缓存中包含的字节数, 本例中, 等于以太网首部(14字节)大小加上以太网帧的数据部分。

向一个或多个 BPF 设备传递分组

883-890 for 循环遍历连接到网络接口的 BPF 设备链表。对每个设备, 分组被递交给 bpf\_filter。如果过滤程序接受了分组, 它返回捕捉到的字节数, 并调用 catchpacket 复制分组。如果过滤程序拒绝了分组, slen 等于 0, 循环继续。循环终止时, bpf\_tap 返回。这一机制确保了同一网络接口上对应了多个 BPF 设备时, 每个设备都能拥有一个独立的过滤程序。

环回驱动程序调用 bpf\_mtap, 向 BPF 传递分组。这个函数与 bpf\_tap 类似, 然而是在 mbuf 链, 而不是在一个内存的连续区域中复制分组。本书中不介绍这个函数。

### 31.5.2 catchpacket 函数

图31-18中, 过滤程序接受了分组后, 将调用 catchpacket, 图31-19给出了这个函数。

bpf.c

```

946 static void
947 catchpacket(d, pkt, pktlen, snaplen, cpf)
948 struct bpf_d *d;
949 u_char *pkt;
950 u_int   pktlen, snaplen;
951 void    (*cpfn) (const void *, void *, u_int);
952 {
953     struct bpf_hdr *hp;
954     int    totlen, curlen;
955     int    hdrlen = d->bd_bif->bif_hdrlen;
956     /*
957      * Figure out how many bytes to move.  If the packet is
958      * greater or equal to the snapshot length, transfer that
959      * much.  Otherwise, transfer the whole packet (unless
960      * we hit the buffer size limit).
961      */
962     totlen = hdrlen + min(snaplen, pktlen);
963     if (totlen > d->bd_bufsize)
964         totlen = d->bd_bufsize;
965     /*
966      * Round up the end of the previous packet to the next longword.
967      */
968     curlen = BPF_WORDALIGN(d->bd_slen);
969     if (curlen + totlen > d->bd_bufsize) {
970         /*
971          * This packet will overflow the storage buffer.
972          * Rotate the buffers if we can, then wakeup any
973          * pending reads.
974          */
975         if (d->bd_fbuf == 0) {
976             /*
977              * We haven't completed the previous read yet,
978              * so drop the packet.
979              */
980             ++d->bd_dcount;
981         }
982     }

```

图31-19 catchpacket 函数

```

981         return;
982     }
983     ROTATE_BUFFERS(d);
984     bpf_wakeup(d);
985     curlen = 0;
986 } else if (d->bd_immediate)
987     /*
988      * Immediate mode is set. A packet arrived so any
989      * reads should be woken up.
990      */
991     bpf_wakeup(d);
992 /*
993  * Append the bpf header.
994  */
995 hp = (struct bpf_hdr *) (d->bd_sbuf + curlen);
996 microtime(&hp->bh_tstamp);
997 hp->bh_datalen = pktlen;
998 hp->bh_hdrlen = hdrlen;
999 /*
1000  * Copy the packet data into the store buffer and update its length.
1001  */
1002 (*cpfn) (pkt, (u_char *) hp + hdrlen, (hp->bh_caplen = totlen - hdrlen));
1003 d->bd_slen = curlen + totlen;
1004 }

```

bpf.c

图31-19 (续)

946-955 `catchpacket`的参数包括：`d`，指向BPF设备结构的指针；`pkt`，指向进入分组的通用指针；`pktlen`，分组被接收时的长度；`snaplen`，从分组中保存下来的字节数；`cpfn`，函数指针，把分组从`pkt`中复制到一块连续内存中。如果分组已经保存连续内存中，则`cptn`等于`bcopy`。如果分组被保存在`mbuf`中(`pkt`指向`mbuf`链表中的第一个`mbuf`，如环回驱动程序)，则`cptn`等于`bpf_mcopy`。

956-964 除了链路层首部和分组，`catchpacket`为每个分组添加`bpf_hdr`。从分组中保存的字节数等于`snaplen`和`pktlen`中较小的一个。处理过的分组和`bpf_hdr`必须能放入分组缓存中(`bd_bufsize`字节)。

### 1. 分组能否放入缓存

965-985 `curlen`等于存储缓存中已有的字节数加上所需的填充字节，以保证下一分组能从长字边界处开始存放。如果进入分组无法放入剩余的缓存空间，说明存储缓存已满。如果空闲缓存不可用(如应用进程正从暂留缓存中读取数据)，则进入分组被丢弃。如是空闲缓存可用，则调用`ROTATE_BUFFERS`宏轮转缓存，并通过`bpf_wakeup`唤醒所有等待输入数据的应用进程。

### 2. 立即方式处理

986-991 如果设备处于立即方式，则唤醒所有等待进程以处理进入分组——内核中没有分组的缓存。

### 3. 添加BPF 首部

992-1004 当前时间(`microtime`)、分组长度和首部长度均保存在`bpf_hdr`中。调用`cptf`所指的函数，把分组复制到存储缓存，并更新存储缓存的长度。因为在把分组从设备缓存传送到某个`mbuf`链表之前，`bpf_tab`已由`leread`直接调用，接收时间戳近似等于实际的接收

时间。

### 31.5.3 bpfread函数

内核把针对BPF设备的read转交给bpfread处理。通过BIOCSRTIMEOUT命令，BPF支持限时读取。这个“特性”也可通过select系统调用来实现，但至少tcpdump还是采用了BIOCSRTIMEOUT，而非select。应用进程提供一个读缓存，能够与设备的暂留缓存大小相匹配。BICOGBLEN命令返回缓存大小。一般情况下，读操作在存储缓存已满时返回。内核轮转缓存，把存储缓存转给暂留缓存，后者在read系统调用时被复制到应用进程提供的读缓存，同时BPF设备继续向存储缓存中存放进入分组。图31-20给出了bpfread。

bpf.c

```
344 int
345 bpfread(dev, uio)
346 dev_t dev;
347 struct uio *uio;
348 {
349     struct bpf_d *d = &bpf_dtab[minor(dev)];
350     int error;
351     int s;
352     /*
353      * Restrict application to use a buffer the same size as
354      * as kernel buffers.
355      */
356     if (uio->uio_resid != d->bd_bufsize)
357         return (EINVAL);
358     s = splimp();
359     /*
360      * If the hold buffer is empty, then do a timed sleep, which
361      * ends when the timeout expires or when enough packets
362      * have arrived to fill the store buffer.
363      */
364     while (d->bd_hbuf == 0) {
365         if (d->bd_immediate && d->bd_slen != 0) {
366             /*
367              * A packet(s) either arrived since the previous
368              * read or arrived while we were asleep.
369              * Rotate the buffers and return what's here.
370              */
371             ROTATE_BUFFERS(d);
372             break;
373         }
374         error = tsleep((caddr_t) d, PRINET | PCATCH, "bpf", d->bd_rtout);
375         if (error == EINTR || error == ERESTART) {
376             splx(s);
377             return (error);
378         }
379         if (error == EWOULDBLOCK) {
380             /*
381              * On a timeout, return what's in the buffer,
382              * which may be nothing. If there is something
383              * in the store buffer, we can rotate the buffers.
384              */
385             if (d->bd_hbuf)
```

图31-20 bpfread 函数

```

386          /*
387          * We filled up the buffer in between
388          * getting the timeout and arriving
389          * here, so we don't need to rotate.
390          */
391          break;
392      if (d->bd_slen == 0) {
393          splx(s);
394          return (0);
395      }
396      ROTATE_BUFFERS(d);
397      break;
398  }
399  }
400  /*
401  * At this point, we know we have something in the hold slot.
402  */
403  splx(s);
404  /*
405  * Move data from hold buffer into user space.
406  * We know the entire buffer is transferred since
407  * we checked above that the read buffer is bpf_bufsize bytes.
408  */
409  error = uiomove(d->bd_hbuf, d->bd_hlen, UIO_READ, uio);
410  s = splimp();
411  d->bd_fbuf = d->bd_hbuf;
412  d->bd_hbuf = 0;
413  d->bd_hlen = 0;
414  splx(s);
415  return (error);
416 }

```

bpf.c

图31-20 (续)

344-357 通过最小设备号在bpf\_dtab中寻找相应的BPF设备。如果读缓存不能匹配BPF设备缓存的大小，则返回EINVAL。

#### 1. 等待数据

358-364 因为多个应用进程能够从一个BPF设备中读取数据，如果有某个进程已先读取了数据，while循环将强迫读操作继续。如果暂留缓存中存在数据，循环被跳过。这与两个应用进程通过两个不同的BPF设备过滤同一个网络接口的情况(见习题31.2)是不同的。

#### 2. 立即方式

365-373 如果设备处于立即方式，且存储缓存中有数据，则轮回缓存，while循环被终止。

#### 3. 无可用的分组

374-384 如果设备不处于立即方式，或者存储缓存中没有数据，则应用进程进入休眠状态，直到某个信号到达，读定时器超时，或者有数据到达暂留缓存。如果有信号到达，则返回EINTR或ERESTART。

记住，应用进程不会见到ERESTART，因为syscall函数将处理这一错误，且不会向应用进程返回这一错误。

#### 4. 查看暂留缓存



385-391 如果定时器超时，且暂留缓存中存在数据，则循环终止。

#### 5. 查看存储缓存

392-399 如果定时器超时，且存储缓存中没有数据，则 read 返回 0。应用进程执行限时读取时，必须考虑到这种情况。如果定时器超时，且存储缓存中存在数据，则把存储缓存转给暂留缓存，循环终止。

如果 tsleep 返回正常且存在数据，同时 while 循环测试失败，则循环终止。

#### 6. 分组可用

400-416 循环终止时，暂留缓存中已有数据。uiomove 从暂留缓存中移出 bd\_hlen 个字节，交给应用进程。把暂留缓存转给空闲缓存，清除缓存计数器，函数返回。 uiomove 调用前的注释指出，uiomove 通常能向应用进程复制 bd\_hlen 字节的数据，因为前面已检查过读缓存大小，确保它大于 BPF 设备缓存的最大值，即 bd\_bufsize。

## 31.6 BPF 的输出

最后，我们讨论如何向带有 BPF 设备的网络接口输出队列中添加分组。首先，应用进程必须构造完整的数据链路帧。对以太网而言，包括源和目的主机的硬件地址和数据帧类型(图4-8)。内核在把它放入接口的输出队列前不会修改链路帧。

### bpfwrite 函数

内核把应用进程的 write 系统调用转给图 31-21 给出的 bpfwrite 处理，数据帧被传给 BPF 设备。

```

437 int
438 bpfwrite(dev, uio)
439 dev_t dev;
440 struct uio *uio;
441 {
442     struct bpf_d *d = &bpf_dtab[minor(dev)];
443     struct ifnet *ifp;
444     struct mbuf *m;
445     int error, s;
446     static struct sockaddr dst;
447     int datlen;

448     if (d->bd_bif == 0)
449         return (ENXIO);

450     ifp = d->bd_bif->bif_ifp;

451     if (uio->uio_resid == 0)
452         return (0);

453     error = bpf_movein(uio, (int) d->bd_bif->bif_dlt, &m, &dst, &datlen);
454     if (error)
455         return (error);

456     if (datlen > ifp->if_mtu)
457         return (EMSGSIZE);

458     s = splnet();
459     error = (*ifp->if_output) (ifp, m, &dst, (struct rentry *) 0);
460     splx(s);

```

bpf.c

图31-21 bpfwrite 函数

```
461     /*
462      * The driver frees the mbuf.
463      */
464     return (error);
465 }
```

— bpf.c

图31-21 (续)

### 1. 检查设备号

437-449 通过最小的设备号选择 BPF设备，它必须已连接到某个网络接口。如果还没有，则返回 ENXIO。

### 2. 向mbuf链中复制数据

450-457 如果 write 给出的写入数据长度等于 0，则立即返回 0。bpf\_movein 从应用进程复制数据到一个 mbuf 链表，并基于由 bpf\_dlt 传递的接口类型计算去除了链路层首部后的分组长度，并在 datlen 中返回该值。它还在 dst 中返回一个已初始化过的 sockaddr 结构。对以太网而言，这个地址结构的类型应该等于 AF\_UNSPEC，说明 mbuf 链中保存了外出数据帧的数据链路层首部。如果分组大于接口的 MTU，则返回 EMSGSIZE。

### 3. 分组排队

458-465 调用 ifnet 结构中指定的 if\_output 函数，得到的 mbuf 链被提交给网络接口。对于以太网，if\_output 等于 ether\_output。

## 31.7 小结

本章中，我们讨论了如何配置 BPF 设备，如何向 BPF 设备递交进入数据帧，及如何在一个 BPF 设备上传送外出数据帧。

一个网络接口可以有多个 BPF 设备，每个 BPF 设备都有自己的过滤程序。存储缓存和暂留缓存最大限度地减少了应用进程为了处理进入数据帧而调用 read 的次数。

本章中只介绍了 BPF 的一些主要特性。有关 BPF 设备过滤程序代码的详细情况和其他一些特性，感兴趣的读者请参阅源代码和 Net/3 手册。

## 习题

31.1 为什么在分组存入 BPF 缓存之前，就能在 catchpacket 中调用 bpf\_wakeup？

31.2 图 31-20 中，我们提到可能会有两个进程在同一 BPF 设备上等待数据。图 31-11 中，我们指出同一时间只能有一个应用进程可以打开一个特定的 BPF 设备。为什么这两种说法都正确呢？

31.3 如果 BIOCSETIF 命令中指定的设备不支持 BPF，会发生什么现象？