# UDP *recvfrom*

As was the case with sending, the socket API provides several mechanisms for receiving a UDP datagram. We begin with a study of *recvfrom()* which, like *sendto(),* does not require the user to pass an address structure and does not support scatter/gather operations via the *iovec* mechanism. The *recvfrom()* function takes the following arguments.

| | |
|---|---|
| fd | File (socket) descriptor. |
| ubuf | Pointer to a buffer to hold the received message. |
| size | The size of the above buffer. |
| addr | Pointer to a structure of type *struct sockaddr_in*. If not NULL and the socket is not connected, the address of the sender of the received message is returned here.. |
| addr_len | A pointer to the size of the structure pointed to by *addr* . If non–zero, its value is changed to the actual size of the address. |
| flags | Can be used to modify the behaviour of the receive operation. Flags supported for UDP sockets include: |

MSG_PEEK: Used to receive data, without dequeuing it. Thus, a subsequent call shall return the same data.

MSG_ERRQUEUE: Used to receive queued errors from socket error queue.

*sys_recvfrom(),* defined in net/socket.c, is the kernel function to which control is eventually passed from *sys_socketcall().* This function is defined in net/socket.c. Its parameters are those passed by the application to *recvfrom().* If an incoming packet is queued for the socket and successfuly copied to the user buffer, *sys_recvfrom* returns its length in bytes. A return value of less than or equal to zero is an indication that an error condition has been encountered and that no data has been returned.

```
1240 asmlinkage long sys_recvfrom(int fd, void * ubuf,
                    size_t size, unsigned flags,
                    struct sockaddr *addr, int *addr_len)
1242 {
1243     struct socket *sock;
1244     struct iovec iov;
1245     struct msghdr msg;
1246     char address[MAX_SOCK_ADDR];
1247     int err,err2;
```

The operation of the sockfd_lookup() function was described in the discussion of UDP sendto. It returns a pointer to struct socket structure corresponding to the fd passed in by the user. If the fd does not index a valid socket NULL is returned and the call fails.

```
1249      sock = sockfd_lookup(fd, &err);
1250      if (!sock)
1251          goto out;
```

Message structures

These structures, introduced in the previous section, are used in both send and receive operations. As with sendto(), the recvfrom() API does not support scatter−gather. Thus it is the responsibility of sys_recvfrom to construct the msghdr and iov.

struct msghdr is defined in include/linux/socket.h.

```
33 struct msghdr {
34     void            *msg_name;
35     int              msg_namelen;
36     struct iovec    *msg_iov;
37     __kernel_size_t  msg_iovlen;
38     void            *msg_control;
39     __kernel_size_t  msg_controllen;
40     unsigned         msg_flags;
41 };
```

Functions of structure elements:

| | |
|---|---|
| msg_name | A pointer to the appropriate struct sockaddr. For TCP/IP sockets this will always be struct sockaddr_in. |
| msg_namelen | The length of the name structure passed in. For TCP/IP sockets this should be sizeof(truct sockaddr_in). |
| msg_iov | A pointer to the IO vector. |
| msg_iovlen | The number of elements in the IO vector which is the number of disjoint fragments of memory comprising the message. |
| msg_control | A pointer to struct cmsghdr. We don't presently under stand the use of cmsgs. |
| msg_controllen | The size of the associated cmsg data. |
| msg_flags | These flags were documented on the first page of this section. |

The struct iovec is defined in include/linux/uio.h.

```
19 struct iovec
20 {
21      void *iov_base;
22      kernel_size_t iov_len;
23 };
```

Functions of structure:

iov_base:               User space pointer to the input data buffer.
iov_len:                Size of buffer pointed to by iov_base. (in bytes)

After recovering the pointer to the struct socket, sys_recvfrom() fills in the struct msghdr and the struct iovec. Since the API supports no mechanism for the receipt of ancillary control data, such data is not collected.

```
1253      msg.msg_control=NULL;
1254      msg.msg_controllen=0;
```

Since the recvfrom() API also permits the caller to pass only a single continous input buffer, a simple I/O vector containing one data block is used.

```
1255      msg.msg_iovlen=1;
1256      msg.msg_iov=&iov;
```

The base pointer for the data block is set to point to the user space address of the data buffer.

```
1257      iov.iov_len=size;
1258      iov.iov_base=ubuf;
```

The name pointer contains the kernel space address of the local buffer in which the address of the sender of the message will be stored temporarily.

```
1259      msg.msg_name=address;
1260      msg.msg_namelen=MAX_SOCK_ADDR;
```

When a socket has O_NONBLOCK flag set, the application will not block(wait) if there is currently no data to receive.

```
1261      if (sock->file->f_flags & O_NONBLOCK)
1262          flags |= MSG_DONTWAIT;
```

The bulk of the work is done by sock_recvmsg(). If the value returned in err is positive, a packet has been successfully received.

```
1263        err = sock_recvmsg(sock, &msg, size, flags);
```

On return to sys_recvmsg() a positive value for err indicates success. If the argument addr, which contains the user space address, is not NULL, and the message was successfully received, the sender's address is returned to user space by the move_addr_to_user() function.

```
1265        if(err >= 0 && addr != NULL && msg.msg_namelen)
1266        {
1267                err2 = move_addr_to_user(address,
                                msg.msg_namelen, addr, addr_len);
1268        if(err2 < 0)
1269            err=err2;
1270        }
1271        sockfd_put(sock);
1272 out:
1273        return err;
1274 }
```

The sock_recvmsg() function, defined in net/socket.c, is the point at which kernel support for all of the recv*() APIs converges.

```
515 int sock_recvmsg(struct socket *sock, struct
        msghdr *msg, int size, int flags)
516 {
517        struct scm_cookie scm;
518
519        memset(&scm, 0, sizeof(scm));
```

For UDP and other sockets of family AF_INET, this indirect call is to inet_recvmsg(). As with sendto() the scm cookie is not used.

```
521        size = sock->ops->recvmsg(sock, msg, size, flags, &scm);
```

If there was no error, scm_receive() is called to processes any control message data. Recall that control messages are used in the passing of sockets among unrelated processes.

```
522        if (size >= 0)
523            scm_recv(sock, msg, &scm, flags);
524
525        return size;
526 }
```

4

The function inet_recvmsg() is defined in net/ipv4/af_inet.c.

```
740 int inet_recvmsg(struct socket *sock, struct msghdr
        *msg, int size, int flags, struct scm_cookie *scm)
742 {
743     struct sock *sk = sock->sk;
744     int addr_len = 0;
745     int err;
```

An indirect call is made to udp_recvmsg. Note that  argument the  "scm " is not used by inet_recvmsg regardless of the transport protocol that is in use and that the flags are repartitioned into the no block flag and the flags originally passed in by the user.

```
747     err = sk->prot->recvmsg(sk, msg, size,
748         flags&MSG_DONTWAIT, flags&~MSG_DONTWAIT, &addr_len);
```

If no error occured then the address length which was filled in by udp_recvmsg() is copied back to the msg->msg_namelen field.

```
749     if (err >= 0)
750         msg->msg_namelen = addr_len;
751     return err;
752 }
```

The udp_recvmsg() function is defined in net/ipv4/udp.c.

```
627 int udp_recvmsg(struct sock *sk, struct msghdr *msg,
            int len, int noblock, int flags, int *addr_len)
629 {
```

The variable sin  is declared to be a pointer of struct sockaddr_in  type and set to the  value of msg->msg_name.

```
630     struct sockaddr_in *sin = (struct sockaddr_in *)
                                        msg->msg_name;
631     struct sk_buff *skb;
632     int copied, err;
```

Set length of address to size of struct sockaddr_in.  No wonder they don't use very many comments in this code.

```
634 /*
635  *      Check any passed addresses
636  */
637     if (addr_len)
638         *addr_len=sizeof(*sin);
```

If MSG_ERRQUEUE is specified in flags, data is received from the error queue of the socket.
What are the operations in this case ––– we need to check out ip_recv_error.

```
640        if (flags & MSG_ERRQUEUE)
641            return ip_recv_error(sk, msg, len);
```

A single UDP packet is retreived from the receive queue associated with the struct sock by skb_recv_datagram().

```
643        skb = skb_recv_datagram(sk, flags, noblock, &err);
```

The skb_recv_datagram() function is defined in net/core/datagram.c.

```
135 struct sk_buff *skb_recv_datagram(struct sock *sk,
        unsigned flags, int noblock, int *err)
136 {
137     int error;
138     struct sk_buff *skb;
139     long timeo;
```

This function sock_error returns value of err flag in sk (struct sock) and clears it.
Here we need to understand how and when the err flag might come to be set..

```
141 /*  Caller is allowed not to check sk->err before
        skb_recv_datagram() */
142     error = sock_error(sk);
143     if (error)
144         goto no_packet;
```

The function sock_error simply returns the negative of the last value to have been stored in sk->err.

```
1198 static inline int sock_error(struct sock *sk)
1199 {
1200     int err=xchg(&sk->err,0);
1201     return -err;
1202 }
```

Continuing in skb_recv_datagram, the value returned by sock_rcvtimeo determines the time to wait (in ticks) for data, if the received packet queue is presently empty.

```
146        timeo = sock_rcvtimeo(sk, noblock);
```

sock_rcvtimeo is defined in include/net/sock.h. sk->rcvtimeo is set to default value of MAX_SCHEDULE_TIMEOUT by sys_socket. It is the maximum value of an unsigned long type. The units are specified in 10 msec jiffies, but this is effectively a wait forever.

```
1241 static inline long sock_rcvtimeo(struct sock *sk,
                                         int  noblock)
1242 {
1243      return noblock ? 0 : sk->rcvtimeo;
1244 }
```

On return to skb_recv_datagram() the main receive loop is entered. Exit from the loop will occur when:

    a datagram has been successfully received
    a timeout occurs (either instantly or after a very long wait)
    an error occurs
    a signal is received

```
148       do {
149       /*   Again only user level code calls this
                function, so nothing interrupt level
150            will suddenly eat the receive_queue.
151
152            Look at current nfs client by the way...
153            However, this function was corrent in
                any case. 8)
154        */
```

If MSG_PEEK is specified in flags, skb_peek() is called. It is passed a pointer to the receive queue header. If the receive queue is non−empty it returns a pointer to the first sk buff without dequeuing it from receive queue.

```
155     if (flags & MSG_PEEK)
156     {
157             unsigned long cpu_flags;
158
159             spin_lock_irqsave(
                        &sk->receive_queue.lock, cpu_flags);
160             skb = skb_peek(&sk->receive_queue);

 97 struct sk_buff_head {
 98 /* These two members must be first. */
 99      struct sk_buff  * next;
100      struct sk_buff  * prev;
101
102      __u32           qlen;
103      spinlock_t      lock;
104 };
```
1

The skb_peek() function is defined in include/linux/skbuff.h. Note that sk_buff_head and sk_buff pointers are used interchangably in line 402. This (bad) practice works correctly because the first two elements of the sk_buff_head structure are the same as those of the sk_buff. If the next pointer points back to the header, the list is empty and NULL is returned.

```
400 static inline struct sk_buff *skb_peek(struct
                                sk_buff_head *list_)
401 {
402     struct sk_buff *list =
                        ((struct sk_buff *)list_)->next;
403     if (list == (struct sk_buff *)list_)
404         list = NULL;
405     return list;
406 }
```

Note that the user count of the buffer is incremented here. This presumably occurs because when the buffer is eventually returned to the peeker, the count will be decremented, and, since the buffer still resides on the queue, we don't want it deleted! However, this is speculation, not fact, at the moment.

```
161             if(skb!=NULL)
162                 atomic_inc(&skb->users);
163             spin_unlock_irqrestore(
                            &sk->receive_queue.lock, cpu_flags);
```

If the MSG_PEEK flag is not specified, skb_dequeue is called. If the queue is non−empty, skb_dequeue will remove the first sk_buff from the list and return a pointer to it. Otherwise it will return NULL. Queue management is not via standard Linux list structures and the received packet queue is rooted at sk−>receive_queue which is an element of type sk_buff_head.

```
164        } else
165            skb = skb_dequeue(&sk->receive_queue);
```

The skb_dequeue() function is defined in include/linux/skbuff.h. It calls __skb_dequeue after obtaining the list's associated lock.

```
589 static inline struct sk_buff *skb_dequeue(struct
                              sk_buff_head *list)
590 {
591     long flags;
592     struct sk_buff *result;
593
594     spin_lock_irqsave(&list->lock, flags);
595     result = __skb_dequeue(list);
596     spin_unlock_irqrestore(&list->lock, flags);
597     return result;
598 }
```

The __skb_dequeue() function does the work of actually removing an sk_buff from the receive queue. Since the sk_buff_head structure starts with the same link pointers as an actual sk_buff structure, it can masquerade as a list element as is done via the cast in line 564. In line 564 prev is set to point to the sk_buff_head. Then in line 565, the local variable next receives the value of the next pointer in the sk_buff_head. The test in line 567 checks to see if the next pointer still points to the sk_buff_head. If so the list was empty. If not the first element is removed from the list and its link fields are zeroed.

```
560 static inline struct sk_buff *__skb_dequeue(struct
                             sk_buff_head *list)
561 {
562     struct sk_buff *next, *prev, *result;
563
564     prev = (struct sk_buff *) list;
565     next = prev->next;
566     result = NULL;
567     if (next != prev) {
568         result = next;
569         next = next->next;
570         list->qlen--;
571         next->prev = prev;
572         prev->next = next;
573         result->next = NULL;
574         result->prev = NULL;
575         result->list = NULL;
576     }
577     return result;
578 }
```

On return to skb_recv_datagram(), if a pointer to an sk_buff is received, it is returned to udp_recvmsg().

```
167             if (skb)
168                 return skb;
```

If the time to wait is zero, return NULL.

```
170        /* User doesn't want to wait */
171            error = -EAGAIN;
172            if (!timeo)
173                goto no_packet;
```

Otherwise, wait for data arrival.

```
175        } while (wait_for_packet(sk, err, &timeo) ==  0);
176
177        return NULL;
178
179 no_packet:
180        *err = error;
181        return NULL;
182 }
```

The wait_for_packet() function is defined in net/core/datagram.c. For reasons not fully understood at the moment it eschews the use of the kernel service routines designed to provide sleep/wakeup services and implements them internally.

```
 61 static int wait_for_packet(struct sock * sk, int *err,
                                      long *timeo_p)
 62 {
 63     int error;
 64
 65     DECLARE_WAITQUEUE(wait, current);
```

```
144 #define DECLARE_WAITQUEUE(name, tsk)                       \
145     wait_queue_t name = __WAITQUEUE_INITIALIZER(name, tsk)
```

The current process sets its state to TASK_INTERRUPTIBLE and adds itself to the queue of waiting processes. A significant amount of additional processing occurs before the process actually goes to sleep though. sk->sleep is of type wait_queue_head_t.

```
 67     __set_current_state(TASK_INTERRUPTIBLE);
 68     add_wait_queue_exclusive(sk->sleep, &wait);
```

11

The err flag of the struct sock is checked for any errors.

```
70        /* Socket errors? */
71        error = sock_error(sk);
72        if (error)
73             goto out_err;
```

The receive queue is tested for still empty. If not a jump is taken to the end of the function.

```
75        if (!skb_queue_empty(&sk->receive_queue))
76             goto ready;
```

See if the shutdown flag of the sk (struct sock) has been set indicating that some manner of receive close is in progress.

```
78        /* Socket shut down? */
79        if (sk->shutdown & RCV_SHUTDOWN)
80             goto out_noerr;
```

Since, a SOCK_DGRAM type socket is connectionless, we always get past this if–statement. Recall that this function is __skb_dequeue() which can be used by protocols other than UDP.

```
82 /*   Sequenced packets can come disconnected.
         If so we report the problem */
83        error = -ENOTCONN;
84        if (connection_based(sk) &&
                   !(sk->state==TCP_ESTABLISHED
                   ||sk->state==TCP_LISTEN))
85             goto out_err;
```

The connection_based() function is defined in net/core/datagram.c.

```
51 static inline int connection_based(struct sock *sk)
52 {
53     return (sk->type==SOCK_SEQPACKET ||
                   sk->type==SOCK_STREAM);
54 }
```

The signal_pending() function checks the sigpending flag of struct task_struct. If this flag is set, an error is returned.

```
87      /* handle signals */
88      if (signal_pending(current))
89          goto interrupted;
```

signal_pending is defined in include/linux/sched.h.

```
632 static inline int signal_pending(struct task_struct *p)
633 {
634     return (p->sigpending != 0);
635 }
```

Finally schedule_timeout() is called to give up control to the scheduler. Control will not return here until a packet is received, a timeout occurs, or a signal is received.

```
91      *timeo_p = schedule_timeout(*timeo_p);
```

This is the standard wakeup action. Restore the task state and remove if from the wait queue.

```
93 ready:
94     current->state = TASK_RUNNING;
95     remove_wait_queue(sk->sleep, &wait);
96     return 0;

98 interrupted:
99     error = sock_intr_errno(*timeo_p);
```

sock_intr_errno is defined in include/net/sock.h.

```
1256 /*  Alas, with timeout socket operations are
         not restartable. Compare this to poll().
     */
1259 static inline int sock_intr_errno(long timeo)
1260 {
1261     return timeo == MAX_SCHEDULE_TIMEOUT ?
                             -ERESTARTSYS : -EINTR;
1262 }

 100 out_err:
 101     *err = error;
 102 out:
 103     current->state = TASK_RUNNING;
 104     remove_wait_queue(sk->sleep, &wait);
 105     return error;
 106 out_noerr:
 107     *err = 0;
 108     error = 1;
 109     goto out;
 110 }
```

schedule_timeout is defined in kernel/sched.c.

```
    /*
            schedule_timeout - sleep until timeout
        @timeout: timeout value in jiffies

    Make the current task sleep until @timeout jiffies
    have elapsed. The routine will return immediately
    unless the current task state has been set (see
    set_current_state()).

    You can set the task state as follows -

    %TASK_UNINTERRUPTIBLE - at least @timeout jiffies
    are guaranteed to pass before the routine returns.
    The routine will return 0

    %TASK_INTERRUPTIBLE - the routine may return early
    if a signal is delivered to the current task. In this
     case the remaining time in jiffies will be  returned, or 0
    if the timer expired in time

    The current task state is guaranteed to be
    TASK_RUNNING when this routine returns.

    Specifying a @timeout value of
    MAX_SCHEDULE_TIMEOUT will schedule the CPU away
    without a bound on the timeout. In this case the
    return value will be %MAX_SCHEDULE_TIMEOUT.

    In all cases the return value is guaranteed to be
    non-negative.
    */

410 signed long schedule_timeout(signed long timeout)
411 {
412     struct timer_list timer;
413     unsigned long expire;
414
```

The above comment clearly says that if timeout value is MAX_SCHEDULE_TIMEOUT, the current process waits for data indefinitely. schedule is called to schedule any processes contending for CPU.

```
415        switch (timeout)
416        {
417        case MAX_SCHEDULE_TIMEOUT:
425            schedule();
426            goto out;
427        default:
428         /*
            Another bit of PARANOID. Note that the retval
            will be 0 since no piece of kernel is
            supposed to do a check for a negative retval
            of schedule_timeout() (since it should never
            happens anyway). You just have the printk()
            that will tell you if something is gone wrong
            and where.
434         */
435            if (timeout < 0)
436            {
437                printk(KERN_ERR "schedule_timeout: wrong
                        timeout "
438                        "value %lx from %p\n", timeout,
439                        builtin_return_address(0));
440                current->state = TASK_RUNNING;
441                goto out;
442            }
443        }
```

If a timeout value less than default max value was specified, a timer is initialized and added to list of timers maintained by kernel and schedule() is invoked.

```
445        expire = timeout + jiffies;
446
447        init_timer(&timer);
448        timer.expires = expire;
449        timer.data = (unsigned long) current;
450        timer.function = process_timeout;
451
452        add_timer(&timer);
453        schedule();
454        del_timer_sync(&timer);
```

process_timeout() is defined in kernel/sched.c. When the above timer expires, this function wakes up current process.

```
377 static void process_timeout(unsigned long __data)
378 {
379     struct task_struct * p = (struct task_struct *)
                              __data;
380
381     wake_up_process(p);
382 }
```

When an event wakes the current process, update time remaining to wait and return to skb_recv_datagram to check if an sk_buff was received. If so, a pointer to it is returned to udp_recvmsg. Otherwise, the process either goes to sleep again or returns a NULL value to udp_recvmsg based on the time remaining to wait.

```
456     timeout = expire - jiffies;
457
458 out:
459     return timeout < 0 ? 0 : timeout;
460 }
```

Back in udp_recvmsg a test is made to see if an sk_buff pointer was returned by skb_recv_datagram() and if not a jump to an error exit pointe is made.

```
644     if (!skb)
645         goto out;
```

If a buffer pointer was returned it is necessary to copy the data back to user space and release the buffer . Since skb->len includes the length of the UDP header at this point (but no longer the MAC or IP header) , copied denotes length of user data.

```
647     copied = skb->len - sizeof(struct udphdr);
```

If data to be returned to user exceeds size of buffer provided, adjust the length downward to fit and set an appropriate flag to indicate that the message was truncated.

```
648     if (copied > len) {
649         copied = len;
650         msg->msg_flags |= MSG_TRUNC;
651     }
```

As with send there are multiple copy mechanisms that depend upon the need to validate the UDP checksum.

```
653         if (skb->ip_summed==CHECKSUM_UNNECESSARY) {
```

If check summing is not required, the skb_copy_datagram_iovec() copies the data to the user space buffer described by the iov.

```
654             err = skb_copy_datagram_iovec(skb,
                        sizeof(struct udphdr),
                        msg->msg_iov, copied);
656         } else if (msg->msg_flags & MSG_TRUNC) {
```

If it was necessary to truncate the message and a checksum is required. It is first necessary to call __udp_checksum_complete() to verify UDP checksum over the entire packet. In case of a checksum error, the sk_buff is freed and an error returned. In case of success only the part of the packet that will fit into the buffer provided is copied.

```
657             if (__udp_checksum_complete(skb))
658                 goto csum_copy_err;

659             err = skb_copy_datagram_iovec(skb,
                sizeof(struct udphdr),msg->msg_iov,
660                         copied);
661         } else {
```

In the final case a checksum is required and the entire packet is to be copied. Here skb_copy_and_csum_datagram_iovec() verifies checksum and copies data from sk buff to I/O vector.

```
662             err = skb_copy_and_csum_datagram_iovec
                        (skb, sizeof(struct udphdr),
                msg->msg_iov);
663
664             if (err == -EINVAL)
665                 goto csum_copy_err;
666         }
667
668     if (err)
669         goto out_free;
```

The sock_recv_timestamp() function records the time stamp, when the sk_buff was received.

```
671        sock_recv_timestamp(msg, sk, skb);
```

sock_recv_timestamp() is defined in include/net/sock.h.

```
1264 static __inline__ void
1265 sock_recv_timestamp(struct msghdr *msg, struct sock
                        *sk, struct sk_buff *skb)
1266 {
```

rcvtstamp is a flag of struct sock. In the case of recvfrom() there is no provision for control messages and the put_cmsg() will do nothing.

```
1267      if (sk->rcvtstamp)
1268           put_cmsg(msg, SOL_SOCKET, SO_TIMESTAMP,
                         sizeof(skb->stamp), &skb->stamp);
1269      else
1270           sk->stamp = skb->stamp;
1271 }
```

Assemble the sender address from elements of the sk_buff.

```
673       /* Copy the address. */
674       if (sin)
675       {
676            sin->sin_family = AF_INET;
677            sin->sin_port = skb->h.uh->source;
678            sin->sin_addr.s_addr = skb->nh.iph->saddr;
679            memset(sin->sin_zero, 0, sizeof(sin->sin_zero));
680       }
```

Depending on control message flags specified, corresponding ancillary control data is collected by ip_cmsg_recv. However, sys_recvfrom does not get any such data.

```
681       if (sk->protinfo.af_inet.cmsg_flags)
682            ip_cmsg_recv(msg, skb);
683       err = copied;
684
```

Free the sk_buff() and return.

```
685 out_free:
686      skb_free_datagram(sk, skb);
687 out:
688      return err;
689
```

In case of a check sum error, the sk_buff is also freed.

```
690 csum_copy_err:
691         UDP_INC_STATS_BH(UdpInErrors);
692
693         /* Clear queue. */
694         if (flags & MSG_PEEK) {
695                 int clear = 0;
696                 spin_lock_irq(&sk->receive_queue.lock);
697                 if (skb == skb_peek(&sk->receive_queue)) {
698                         __skb_unlink(skb, &sk->receive_queue);
699                         clear = 1;
700                 }
701                 spin_unlock_irq(&sk->receive_queue.lock);
702                 if (clear)
703                         kfree_skb(skb);
704         }
705
706         skb_free_datagram(sk, skb);
707
708         return -EAGAIN;
709 }
```

Transfering data from an sk_buff to an I/O vector

This procedure appears (and is) far more complex in the worst case than is actually the case in practice. The problem lies with the ''unususual'' implementation of the sk_buff. In the worst case, a single packet could consist of the following:

- an instance of the struct skbuff buffer header.
- a kmalloc'd ''data'' area allocated to hold the headers of the packet.
- up to 6 unmapped single page structures called fragments holding packet data
- a pointer to an additional sk_buff which may contain all 4 of these elements!

This possibility leads to a recursive implementation of checksumming and data movement code. Fortunately, in practice a UDP packet always consists of only:

- an instance of the struct skbuff buffer header.
- a kmalloc'd ''data'' area allocated holding both packet headers an data.

The ''worst case'' sk_buff structure is mapped by the struct skb_shared_info defined in include/linux/skbuff.h. When used, this structure resides at the end of the kmalloc'd data area and is pointed to by the end element of the struct sk_buff header.

```
119 /* This data is invariant across clones and lives at
120  * the end of the header data, ie. at skb->end.
121  */
122 struct skb_shared_info {
123     atomic_t        dataref;
124     unsigned int    nr_frags;
125     struct sk_buff  *frag_list;
126     skb_frag_t      frags[MAX_SKB_FRAGS];
127 };
```

Functions of structure elements:

| | |
|---|---|
| frag_list | may contain zero a pointer to the next sk_buff in the chain. When is it used?? |
| frags | is a six element array of skb_frag_t type. These are the unmapped single page entities. |
| nr_frags | denotes number of elements of frags array in use. |

```
108 #define MAX_SKB_FRAGS 6
109
110 typedef struct skb_frag_struct skb_frag_t;
111
112 struct skb_frag_struct
113 {
114     struct page *page;
115     __u16 page_offset;
116     __u16 size;
117 };
```

Functions of structure elements:

| | |
|---|---|
| page | Pointer to a struct page. |
| offset | Offset in page from where data is stored. |
| size | Size of data. |

The buffer header of type struct sk_buff contains two members, namely len and data_len, used to describe the length of the received packet. The skb->len field denotes length of the amount of data in the packet that remains to be processed.  That is, it is initially set to the length of all headers and application data.   As headers are removed as the packet is passed up the stack, the value of skb->len is decremented by the length of each network header removed.   The value of skb->data_len is the amount of data that is held in fragments and in chained sk_buffs. It is used by TCP but appears to have no use in processing UDP packets.

The function skb_copy_datagram_iovec(), defined in net/core/datagram.c, is used to copy a UDP datagram when checksumming is not required. In this case the value of offset is the size of the UDP header and skb->data_len is 0.

```
204 int skb_copy_datagram_iovec(const struct sk_buff *skb,
              int offset, struct iovec *to, int len)
206 {
207     int i, copy;
```

It is the case that skb->len includes the kmalloc'd stuff but that skb->datalen includes only that which is in the appendicies. Thus start would be set here to the amount of data in the kmalloc'd part which would be just what is needed!

```
208     int start = skb->len - skb->data_len;
```

The comment is misleading and apparently reflects the philosophy that the kmalloc'd part of the sk_buff structure is for storage of network header elements. What is actually happening in the case of UDP is that memcpy_toiovec() is being passed a pointer to the start of the user data along with the length of the user data. In the standard case (no fragments) the value of len will become 0 at line 216 and the function will return.

```
210     /* Copy header. */
211     if ((copy = start-offset) > 0) {
212         if (copy > len)
213             copy = len;
214         if (memcpy_toiovec(to, skb->data + offset, copy))
215             goto fault;
216         if ((len -= copy) == 0)
217             return 0;
218         offset += copy;
219     }
```

The memcpy_to_iovec() function is defined in net/core/iovec.c. It copies kernel data into an I/O vector. Note that as data is copied to the iovec, the len field of the element which is the recipient is decremented and the base pointer is incremented. This strategy makes it possible, albeit slightly inefficient, for callers that are passing multiple fragments of a packet to be copied to always just pass the base address of the iovec. Elements that have been previously filled will just be bypassed in the while loop because the if statement at line 88 will find that such elements have iov_len equal to 0.

```
82 int memcpy_toiovec(struct iovec *iov, unsigned char
                                *kdata, int len)
83 {
84     int err = -EFAULT;
85
86     while(len>0)
87     {
88         if(iov->iov_len)
89         {
```

min_t returns minimum of two arguments.

```
90                 int copy = min_t(unsigned int,
                        iov->iov_len, len);
91                 if (copy_to_user(iov->iov_base, kdata,
                        copy))
92                     goto out;
93                 kdata+=copy;
```

Update available buffer space and base pointer of I/O vector.

```
94                 len-=copy;
95                 iov->iov_len-=copy;
96                 iov->iov_base+=copy;
97             }
98             iov++;
99         }
100     err = 0;
101 out:
102     return err;
103 }
```

If there do exist fragments skb_copy_datagram_iovec() will continue and copy data from page
fragments into the I/O vector.

```
221        /*    Copy paged appendix. Hmm... why does this
                 look so complicated? */
222        for (i=0; i<skb_shinfo(skb)->nr_frags; i++) {
```

skb_shinfo is defined in include/linux/skbuff.h.  It simply returns a pointer to skb_shared_info
structure that is pointed to by skb->end.

```
247        /* Internal */
248        #define skb_shinfo(SKB)
                       ((struct skb_shared_info *)((SKB)->end))

223              int end;
224
225              BUG_TRAP(start <= offset+len);
226
```

In the first iteration of this loop start contains the offset from the start of the packet data (including
UDP header) of the beginning of the paged appendix.   Thus end is set to the offset of the 1st byte
beyond the data in the paged appendix and copy is set to the amount of data in this element of the
paged appendix.

```
227              end = start + skb_shinfo(skb)->frags[i].size;
228              if ((copy = end-offset) > 0) {
229                    int err;
230                    u8  *vaddr;
231                    skb_frag_t *frag =
                            &skb_shinfo(skb)->frags[i];
232                    struct page *page = frag->page;
233
234                    if (copy > len)
235                        copy = len;
```

Get  logical address of page corresponding to page. Copy data  from fragment  into I/O vector
using memcpy_to_iovec.

```
236                    vaddr = kmap(page);
237                    err = memcpy_toiovec(to, vaddr +
                                    frag->page_offset +
                                    offset-start, copy);
239                    kunmap(page);
240                    if (err)
241                        goto fault;
242                    if (!(len -= copy))
243                        return 0;
244                    offset += copy;
245              }
246         start = end;
247     }
```

Finally, if there exist additional sk_buffs in the chain,  the are processed via a recursive call to
skb_copy_datagram_iovec().    This incredible structure is actually a tree of general degree.

```
249        if (skb_shinfo(skb)->frag_list) {
250             struct sk_buff *list;
251
252             for (list = skb_shinfo(skb)->frag_list;
                          list; list=list->next) {
253                 int end;
254
255                 BUG_TRAP(start <= offset+len);
256
257                 end = start + list->len;
258                 if ((copy = end-offset) > 0) {
259                     if (copy > len)
260                         copy = len;
261                     if (skb_copy_datagram_iovec(list,
                                offset-start, to, copy))
262                         goto fault;

263                     if ((len -= copy) == 0)
264                         return 0;
265                     offset += copy;
266                 }
267                 start = end;
268             }
269        }
270     if (len == 0)
271         return 0;
272
273 fault:
274     return -EFAULT;
275 }
```

UDP checksum

The function __udp_checksum_complete() is defined in net/ipv4/udp.c and is used when a UDP datagram must be truncated. Its mission is to make sure that the entire datagram passes the checksum test. If so, then it is safe to return the truncated part to user space. It consists of call to skb_checksum() followed by a call to csum_fold() which converts the 32 bit checksum to a proper 16 bit one.

```
611 static __inline__ int __udp_checksum_complete(struct
                                    sk_buff *skb)
612 {
613     return (unsigned short)csum_fold(skb_checksum(skb,
                    0, skb->len, skb->csum));
614 }
```

skb_checksum is defined in net/core/skbuff.c. It computes checksum on data present in sk buff using function csum_partial.

```
 998 /* Checksum skb data. */
 999
1000 unsigned int skb_checksum(const struct sk_buff *skb,
                int offset, int len, unsigned int csum)
1001 {
1002     int i, copy;
1003     int start = skb->len - skb->data_len;
1004     int pos = 0;
```

Compute a partial checksum, csum , on UDP header and any data past it. Note that argument offset has the value zero.

```
1006        /* Checksum header. */
1007        if ((copy = start-offset) > 0) {
1008            if (copy > len)
1009                copy = len;
1010            csum = csum_partial(skb->data+offset, copy,
                csum);
```

Prototype of function csum_partial is specified in include/asm−i386/checksum.h. The function is defined in arch/i386/checksum.S. It is written in assembly language.

```
      computes the checksum of a memory block at buff,  length len,
      and adds in "sum" (32-bit) returns a 32-bit number suitable
      for  feeding into itself or csum_tcpudp_magic this function
      must be called with even lengths, except for the last
      fragment, which may be odd it's best to have buff aligned on
      a 32-bit boundary

   17 asmlinkage unsigned int csum_partial(const unsigned
         char * buff, int len, unsigned int sum);


1011            if ((len -= copy) == 0)
1012                  return csum;
1013            offset += copy;
1014            pos = copy;
1015      }
```

On data in each fragment, compute partial checksum, csum2, using csum_partial. Add the two checksums, csum and csum2, using csum_block_add.

```
1017          for (i=0; i<skb_shinfo(skb)->nr_frags; i++) {
1018                  int end;
1019
1020                  BUG_TRAP(start <= offset+len);
1021
1022                  end = start + skb_shinfo(skb)->frags[i].size;
1023                  if ((copy = end-offset) > 0) {
1024                          unsigned int csum2;
1025                          u8 *vaddr;
1026                          skb_frag_t *frag =
                                    &skb_shinfo(skb)->frags[i];
1027
1028                          if (copy > len)
1029                                  copy = len;
1030                          vaddr = kmap_skb_frag(frag);

1031                          csum2 = csum_partial(vaddr +
                                    frag->page_offset +
1032                                    offset-start, copy, 0);
1033                          kunmap_skb_frag(vaddr);
1034                          csum = csum_block_add(csum, csum2, pos);
1035                          if (!(len -= copy))
1036                                  return csum;
1037                          offset += copy;
1038                          pos += copy;
1039                  }
1040                  start = end;
1041          }
```

csum_block_add is defined as an inline function in include/net/checksum.h. It performs an adjustment to csum2, if csum1 is a partial checksum of an odd number of bytes of data. <span style="color:red">How does this work?</span>

```
138 static inline unsigned int
139 csum_block_add(unsigned int csum, unsigned int csum2,
                                    int offset)
140 {
141     if (offset&1)
142         csum2 = ((csum2&0xFF00FF)<<8)
                        +((csum2>>8)&0xFF00FF);
143     return csum_add(csum, csum2);
144 }
```

csum_add is defined as below. It combines two partial checksums.

```
127  static inline unsigned int csum_add(unsigned int csum,
                                          unsigned int addend)
128  {
129      csum += addend;
130      return csum + (csum < addend);
131  }
```

On data present in each sk_ buff on fragment list, compute partial checksum, csum2, using csum_partial. Add the two checksums as done earlier.

```
1043       if (skb_shinfo(skb)->frag_list) {
1044           struct sk_buff *list;
1045
1046           for (list = skb_shinfo(skb)->frag_list; list;
                       list=list->next) {
1047               int end;
1048
1049               BUG_TRAP(start <= offset+len);

1051               end = start + list->len;
1052               if ((copy = end-offset) > 0) {
1053                   unsigned int csum2;
1054                   if (copy > len)
1055                       copy = len;
1056                   csum2 = skb_checksum(list, offset-
                           start, copy, 0);
1057                   csum = csum_block_add(csum, csum2,
                               pos);
1058                   if ((len -= copy) == 0)
1059                       return csum;
1060                   offset += copy;
1061                   pos += copy;
1062               }
1063               start = end;
1064           }
1065       }
1066       if (len == 0)
1067           return csum;
1068
1069       BUG();
1070       return csum;
1071  }
```

csum_fold is defined as below. It folds a 32–bit checksum to a 16–bit value. How does this work?

```
 99 /*
100  *        Fold a partial checksum
101  */
102
103 static inline unsigned int csum_fold(unsigned int sum)
104 {
105            asm__("
106            addl %1, %0
107            adcl $0xffff, %0"
108            "
109            : "=r" (sum)
110            : "r" (sum << 16), "" (sum &  0xffff0000)
111            );
112            return (~sum) >> 16;
113 }
```

skb_copy_and_csum_datagram_iovec is defined in net/core/datagram.c. It gets called when checksum is necessary and message has not been truncated.

```
370 int skb_copy_and_csum_datagram_iovec(const struct
        sk_buff *skb, int hlen, struct iovec *iov)
371 {
372     unsigned int csum;
373     int chunk = skb->len - hlen;

375     /*   Skip filled elements. Pretty silly, look at
            memcpy_toiovec, though 8) */
376     while (iov->iov_len == 0)
377         iov++;

379     if (iov->iov_len < chunk) {
```

We have covered this case, where the size of specified buffer is less than available data.

```
380         if ((unsigned short)csum_fold(skb_checksum
                (skb, 0, chunk+hlen, skb->csum)))
381             goto csum_error;
382         if (skb_copy_datagram_iovec(skb, hlen, iov,
                chunk))
383             goto fault;
384     } else {
```

Obtain checksum value of UDP header. skb_copy_and_csum_datagram gets called which performs both checksumming and copy of data.

```
385         csum = csum_partial(skb->data, hlen,
                            skb->csum);
386         if (skb_copy_and_csum_datagram(skb, hlen,
                iov->iov_base, chunk, &csum))
387             goto fault;
388         if ((unsigned short)csum_fold(csum))
389             goto csum_error;
390         iov->iov_len -= chunk;
391         iov->iov_base += chunk;
392     }
393     return 0;
394
395 csum_error:
396     return -EINVAL;
397
398 fault:
399     return -EFAULT;
400 }
```

`skb_copy_and_csum_datagram` is defined in net/core/datagram.c.

```
277 int skb_copy_and_csum_datagram(const struct sk_buff   *skb,
        int offset, u8 *to, int len, unsigned int *csump)
278 {
279     int i, copy;
280     int start = skb->len - skb->data_len;
281     int pos = 0;
282
283     /* Copy header.  */
284     if ((copy = start-offset) > 0) {
285         int err = 0;
286         if (copy > len)
287             copy = len;
288         *csump = csum_and_copy_to_user(skb->data
                                offset, to, copy, *csump,&err);
289         if (err)
290                 goto fault;
291         if ((len -= copy) == 0)
292             return 0;
293         offset += copy;
294         to += copy;
295         pos = copy;
296     }
```

`csum_and_copy_to_user` is defined in include/asm−i386/checksum.h.

```
181 /*
182  *        Copy and checksum to user
183  */
184 #define HAVE_CSUM_COPY_USER
185 static __inline__ unsigned int csum_and_copy_to_user
    (const char *src, char *dst, int len, int sum,
        int *err_ptr)
187 {
188     if (access_ok(VERIFY_WRITE, dst, len))
189         return csum_partial_copy_generic(src, dst,
            len, sum, NULL, err_ptr);
190
191     if (len)
192         *err_ptr = -EFAULT;
193
194     return -1; /* invalid checksum */
195 }
```

Prototype of csum_partial_copy_generic is defined in include/asm–i386/checksum.h.

```
    /*
        the same as csum_partial, but copies from src
    while it checksums, and handles user-space pointer
    exceptions correctly, when needed.

        here even more important to align src and dst on a
    32-bit (or even better 64-bit) boundary
    */

    27 asmlinkage unsigned int csum_partial_copy_generic(
    const char *src, char *dst, int len, int sum,
        int *src_err_ptr, int *dst_err_ptr);
```

On return to skb_copy_and_csum_datagram repeat for of each fragment.

```
298        for (i=0; i<skb_shinfo(skb)->nr_frags; i++) {
299            int end;
300
301            BUG_TRAP(start <= offset+len);
302
303            end = start + skb_shinfo(skb)->frags[i].size;
304            if ((copy = end-offset) > 0) {
305                unsigned int csum2;
306                int err = 0;
307                u8  *vaddr;
308                skb_frag_t *frag =
                   &skb_shinfo(skb)->frags[i];
309                struct page *page = frag->page;
310
311                if (copy > len)
312                    copy = len;
313                vaddr = kmap(page);
314                csum2 = csum_and_copy_to_user(vaddr +
                   frag->page_offset +
315                        offset-start, to, copy, 0,
                   &err);
316                kunmap(page);
317                if (err)
318                    goto fault;
319                *csump = csum_block_add(*csump, csum2,
                   pos);
320                if (!(len -= copy))
321                    return 0;
322                offset += copy;
323                to += copy;
324                pos += copy;
325            }
326            start = end;
327        }
```

And then  on data present in each sk buff on fragment list.

```
329        if (skb_shinfo(skb)->frag_list) {
330             struct sk_buff *list;
331
332             for (list = skb_shinfo(skb)->frag_list; list;
                      list=list->next) {
333                  int end;
334
335                  BUG_TRAP(start <= offset+len);
336
337                  end = start + list->len;
338                  if ((copy = end-offset) > 0) {
339                       unsigned int csum2 = 0;
340                       if (copy > len)
341                            copy = len;
342                       if (skb_copy_and_csum_datagram
                                     (list, offset-start, to,
                                 copy, &csum2))
343                            goto fault;
344                       *csump = csum_block_add(*csump,
                          csum2, pos);
345                       if ((len -= copy) == 0)
346                            return 0;
347                       offset += copy;
348                       to += copy;
349                            pos += copy;
350                  }
351                  start = end;
352             }
353        }
354      if (len == 0)
355             return 0;
356
357 fault:
358      return -EFAULT;
359 }
```

Socket control messages

scm_rcv is defined in include/net/scm.h. It is called by sock_recvmsg after the return from the call to inet_recvmsg().

```
45 static __inline__ void scm_recv(struct socket *sock,
      struct msghdr *msg, struct scm_cookie *scm, int
      flags)
47 {
48     if (!msg->msg_control)
49     {
```

Since no control buffer was specified by sys_recvfrom, scm_destroy() is called. As shown below the call is a no–op here because scm–>fpl is NULL. The MSG_CTRUNC flag indicates that control data was discarded.

```
50           if (sock->passcred || scm->fp)
51               msg->msg_flags |= MSG_CTRUNC;
52           scm_destroy(scm);
53           return;
54     }
```

Because of the return on line 53, control can't possibly reach this point when called by recvfrom(). However, other callers might pass a msg_control pointer. In that case, if the passcred flag is set, any socket control data is copied into msg control buffer.

```
56     if (sock->passcred)
57         put_cmsg(msg, SOL_SOCKET, SCM_CREDENTIALS,
                     sizeof(scm->creds), &scm->creds);
59     if (!scm->fp)
60         return;
```

Any passed file descriptors are freed by scm_detach_fds.

```
62     scm_detach_fds(msg, scm);
63 }
```

The scm_destroy function is a wrapper for __scm_destroy.

```
 97 void __scm_destroy(struct scm_cookie *scm)
 98 {
 99       struct scm_fp_list *fpl = scm->fp;
100       int i;
101
102       if (fpl) {
103               scm->fp = NULL;
104               for (i=fpl->count-1; i>=0; i--)
105                       fput(fpl->fp[i]);
106               kfree(fpl);
107       }
108 }
```