



[oss-sec](#) mailing list archives



◀ [By Date](#) ▶ ◀ [By Thread](#) ▶



Linux kernel: CVE-2022-1015,CVE-2022-1016 in nf_tables cause privilege escalation, information leak

From: David Bouman <davidbouman35 () gmail com>

Date: Mon, 28 Mar 2022 20:28:21 +0200

Hello list,

I'm reporting two linux kernel vulnerabilities in the nf_tables component of the netfilter subsystem that I found.

CVE-2022-1015 pertains to an out of bounds access in nf_tables expression evaluation due to validation of user register indices. It leads to local privilege escalation, for example by overwriting a stack return address OOB with a crafted nft_expr_payload.

CVE-2022-1015 is exploitable starting from commit 345023b0db3 ("netfilter: nftables: add nft_parse_register_store() and use it"), v5.12 and has been fixed in commit 6e1acfa387b9 ("netfilter: nf_tables: validate registers coming from userspace").

The bug has been present since commit 49499c3e6e18 ("netfilter: nf_tables: switch registers to 32 bit addressing"), but to my knowledge has not been exploitable until v5.12.

CVE-2022-1016 pertains to uninitialized stack data in the nft_do_chain routine. CVE-2022-1016 is exploitable starting from commit 96518518cc41 (original merge of nf_tables), v3.13-rcl, and has been fixed in commit 4c905f6740a3 ("netfilter: nf_tables: initialize registers in nft_do_chain()").

I will be releasing a detailed blog post and exploit code for both vulnerabilities in a few days.

Root cause CVE-2022-1016: (it is the shortest, so I will begin with it)

The nft_do_chain routine in net/netfilter/nf_tables_core.c does not initialize the register data that nf_tables expressions can read from- and write to. These expressions inherently exhibit side effects that can be used to determine the register data, which can contain kernel image pointers, module pointers, and allocation pointers depending on the code path taken to end up at nft_do_chain.

...

```
unsigned int
nft_do_chain(struct nft_pktinfo *pkt, void *priv)
```

```

{
    const struct nft_chain *chain = priv, *basechain = chain;
    const struct net *net = nft_net(pkt);
    struct nft_rule *const *rules;
    const struct nft_rule *rule;
    const struct nft_expr *expr, *last;
    struct nft_regs regs; // <----- VULNERABLE! NOT INITIALIZED.
    unsigned int stackptr = 0;
    struct nft_jumpstack jumpstack[NFT_JUMP_STACK_SIZE];
    bool genbit = READ_ONCE(net->nft.gencursor);
    struct nft_traceinfo info;

    info.trace = false;
    if (static_branch_unlikely(&nft_trace_enabled))
        nft_trace_init(&info, pkt, &regs.verdict, basechain);
do_chain:
    if (genbit)
        rules = rcu_dereference(chain->rules_gen_1);
    else
        rules = rcu_dereference(chain->rules_gen_0);

next_rule:
    rule = *rules;
    regs.verdict.code = NFT_CONTINUE;
    for (; *rules ; rules++) {
        rule = *rules;
        nft_rule_for_each_expr(expr, last, rule) {
            if (expr->ops == &nft_cmp_fast_ops)
                nft_cmp_fast_eval(expr, &regs);
            else if (expr->ops == &nft_bitwise_fast_ops)
                nft_bitwise_fast_eval(expr, &regs);
            else if (expr->ops != &nft_payload_fast_ops ||
                    !nft_payload_fast_eval(expr, &regs, pkt))
                expr_call_ops_eval(expr, &regs, pkt);
            ...
        }
    }
    ...
}

```

Root cause CVE-2022-1015:

(below is pasted from my original security () kernel org report)

Hello, I'm mailing to report a vulnerability I found in nf_tables component of the netfilter subsystem. The vulnerability gives an attacker a powerful primitive that can be used to both read from and write to relative stack data. This can lead to arbitrary code execution by an attacker.

In order for an unprivileged attacker to exploit this issue, unprivileged user- and network namespaces access is required (CLONE_NEWUSER | CLONE_NEWNET). The bug relies on a compiler optimization that introduces behavior that the maintainer did not account for, and most likely only occurs on kernels with `CONFIG_CC_OPTIMIZE_FOR_PERFORMANCE=y`. I successfully exploited the bug on x86_64 kernel version 5.16-rc3, but I believe this vulnerability exists across different kernel versions and architectures (more on this later).

Without further ado:

The bug resides in `linux/net/netfilter/nf_tables_api.c`, in the `nft_validate_register_store` and `nft_validate_register_load` routines. These routines are used to check if nft expression parameters supplied by the user are sound and won't cause OOB stack accesses when evaluating the expression.

From my 5.16-rc3 kernel source (d58071a8a76d779eedab38033ae4c821c30295a5: Linux 5.16-rc3):

```

nft_validate_register_store:
...
static int nft_validate_register_store(const struct nft_ctx *ctx,
    enum nft_registers reg,
    const struct nft_data *data,
    enum nft_data_types type,
    unsigned int len)
{
    int err;

    switch (reg) {
        ...
    default:
        if (reg < NFT_REG_1 * NFT_REG_SIZE / NFT_REG32_SIZE)
            return -EINVAL;
        if (len == 0)
            return -EINVAL;
        if (reg * NFT_REG32_SIZE + len >
            sizeof_field(struct nft_regs, data))
            return -ERANGE;

        if (data != NULL && type != NFT_DATA_VALUE)
            return -EINVAL;
        return 0;
    }
}
...

nft_validate_register_load:
...
static int nft_validate_register_load(enum nft_registers reg, unsigned int len)
{
    if (reg < NFT_REG_1 * NFT_REG_SIZE / NFT_REG32_SIZE)
        return -EINVAL;
    if (len == 0)
        return -EINVAL;
    if (reg * NFT_REG32_SIZE + len > sizeof_field(struct nft_regs, data))
        return -ERANGE;

    return 0;
}
...

```

The problem lies in the fact that ``enum nft_registers reg`` is not guaranteed only be a single byte. As per the C89 specification, 3.1.3.3 Enumeration constants: ``An identifier declared as an enumeration constant has type int.``

Effectively this implies that the compiler is free to emit code that operates on ``reg`` as if it were a 32-bit value. If this is the case (and it is on the kernel I tested), a user can forge an expression register value that will overflow upon multiplication with ``NFT_REG32_SIZE`` (4) and upon addition with ``len``, will be a value smaller than ``sizeof_field(struct nft_regs, data)`` (0x50). Once this check passes, the least significant byte of ``reg`` can still contain a value that will index outside of the bounds of the ``struct nft_regs regs`` that it will later be used with.

Take for example a ``reg`` value of ``0xffffffff8`` and a ``len`` value of ``0x40``. The expression ``reg * 4 + len`` will then result in ``0xffffffffe0 + 0x40 = 0x20``, which is lower than ``0x50``. This makes that a value of ``0xf8`` is recognized as a valid index, and is subsequently assigned to a register value in the expression info structs.

Finally, I used a nft payload expression to write my arbitrary data supplied in a packet to the stack in order to overwrite a return address and execute a ROP chain.

An alternative exploitation strategy would be to overwrite to verdict register (including its chain pointer) to arbitrary values, as you can now get an register index of 0 in the same manner.

David Bouman

[← By Date →](#) [← By Thread →](#)

Current thread:

Linux kernel: CVE-2022-1015,CVE-2022-1016 in nf_tables cause privilege escalation, information leak *David Bouman (Mar 28)*



Nmap Security Scanner

Ref Guide

Install Guide

Docs

Download

Nmap OEM

Npcap packet capture

User's Guide

API docs

Download

Npcap OEM

Security Lists

Nmap Announce

Nmap Dev

Full Disclosure

Open Source Security

BreachExchange

Security Tools

Vuln scanners

Password audit

Web scanners

Wireless

Exploitation

About

About/Contact

Privacy

Advertising

Nmap Public Source License

