### Four Bytes of Power: Exploiting CVE-2021-26708 in the Linux Kernel

#### Alexander Popov

**Positive Technologies** 

April 9, 2021



- Alexander Popov
- Linux kernel developer since 2013

• Security researcher at **POSITIVE TECHNOLOGIES** 

## Agenda

#### • CVE-2021-26708 overview

- Bugs and fixes
- Disclosure procedure
- $\odot$  Exploitation for local privilege escalation on  $\times 86_{64}$ 
  - Hitting the race condition
  - Four-byte memory corruption
  - Long way to arbitrary read/write
- Server by passing SMEP and SMAP
- Possible exploit mitigation

- LPE in the Linux kernel
- Bug type: race condition
- Refers to **5** similar bugs in the virtual socket implementation
- Major Linux distros ship CONFIG\_VSOCKETS and CONFIG\_VIRTIO\_VSOCKETS as a kernel modules

#### Attack Surface

- The vulnerable modules are automatically loaded
- Just create a socket for the AF\_VSOCK domain:

vsock = socket(AF\_VSOCK, SOCK\_STREAM, 0);

- That's available for unprivileged users
- User namespaces are not needed for that

5 / 54

- I used the syzkaller fuzzer with custom modifications
- KASAN got a suspicious kernel crash in virtio\_transport\_notify\_buffer\_size()
- The fuzzer failed to reproduce this crash 🤪
- I inspected the source code and developed the reproducer manually

I found a confusing bug in vsock\_stream\_setsockopt():

```
struct sock *sk;
struct vsock_sock *vsk;
const struct vsock_transport *transport;
```

```
Let me look at it...
```

```
sk = sock->sk;
vsk = vsock_sk(sk);
transport = vsk->transport;
```

lock\_sock(sk);

I found a confusing bug in vsock\_stream\_setsockopt():

```
struct sock *sk;
struct vsock_sock *vsk;
const struct vsock_transport *transport;
```

```
sk = sock->sk;
vsk = vsock_sk(sk);
transport = vsk->transport;
/* vsk->transport value may change here! */
lock_sock(sk);
```





- vsk->transport may change when the socket lock is not acquired
- In that case, the local variable value is out-of-date
- That is an obvious race condition bug
- I found **five** similar bugs in net/vmw\_vsock/af\_vsock.c
- Searching the git history helped to understand the reason

- Initially, the transport for a virtual socket was **not** able to change
- The bugs were implicitly introduced in November 2019 when VSOCK multi-transport support was added
- Fixing this vulnerability is trivial:

```
sk = sock->sk;
vsk = vsock_sk(sk);
- transport = vsk->transport;
lock_sock(sk);
+ transport = vsk->transport;
```

- November 14, 2019 Bugs were introduced
- January 7, 2021 My custom syzkaller got a crash
- January 11, 2021 I started the investigation
- January 30, 2021
  - My PoC exploit and fixing patch were ready
  - I sent the crasher and patch to security@kernel.org
  - Review started

#### Disclosure Procedure (1)

- I got very prompt replies from Linus Torvalds and Greg Kroah-Hartman
- We concluded on this procedure:
  - sending my patch to LKML in public
  - Imaging it to the upstream and backporting to the stable trees
  - (a) informing the distros about the security-relevance via linux-distros ML
  - Isclosing that at oss-security@lists.openwall.com when distros allow me
- The first step is questionable, though

12 / 54

#### Disclosure Procedure (2)

• Linus decided to merge my patch without any disclosure embargo

Linus:

"This patch doesn't look all that different from the kinds of patches we do every day"

• I obeyed and proposed that I should send it to LKML in public

Rationale

Anybody can find kernel vulnerability fixes by filtering kernel commits that didn't

appear on the mailing lists <u>https://arxiv.org/abs/2009.01694</u>

#### Timeline: Part 2

- February 2, 2021 The v2 of my patch was merged into Linus' tree
- February 4, 2021
  - Greg applied it to the affected stable trees
  - ▶ I informed linux-distros ML that the fixed bugs are exploitable
  - ▶ I asked how much time Linux distros need before my public disclosure
  - But I got this reply:

If the patch is committed upstream, then the issue is public. Please send to oss-security immediately.

- ► I made the public announcement: <u>https://seclists.org/oss-sec/2021/q1/107</u>
- February 5, 2021 CVE-2021-26708 is assigned

#### Pondering over the Disclosure Procedure

The question is rising: Is this "merge ASAP" procedure compatible with the linux-distros mailing list?

Counter-example: how I reported CVE-2017-2636 to security@kernel.org

- Kees Cook and Greg organized a one-week disclosure embargo
- Linux distributions in the linux-distros ML integrated my fix in their security updates in no rush
- Security updates were published synchronously when the embargo ended
- More info in this article: <a href="https://a13xp0p0v.github.io/2017/03/24/CVE-2017-2636.html">https://a13xp0p0v.github.io/2017/03/24/CVE-2017-2636.html</a>



# NOW ABOUT THE MEMORY CORRUPTION

16 / 54

- I exploited the race condition in vsock\_stream\_setsockopt()
- Reproducing it requires two threads
- The first one calls setsockopt()

setsockopt(vsock, PF\_VSOCK, S0\_VM\_SOCKETS\_BUFFER\_SIZE,

&size, sizeof(unsigned long));

• The second thread should change the virtual socket transport

# Changing VSOCK Transport

• It is performed by reconnecting to the virtual socket:

```
struct sockaddr_vm addr = {
    .svm_family = AF_VSOCK,
};
addr.svm_cid = VMADDR_CID_LOCAL;
connect(vsock, (struct sockaddr *)&addr, sizeof(struct sockaddr_vm));
addr.svm_cid = VMADDR_CID_HYPERVISOR;
connect(vsock, (struct sockaddr *)&addr, sizeof(struct sockaddr_vm));
```

• Meanwhile, vsock\_stream\_setsockopt() in a parallel thread is

#### trying to acquire the lock

Alexander Popov

#### Race Condition: Full Picture

```
Thread 1: reconnecting to vsock
```

```
vsock_stream_connect() /* VMADDR_CID_LOCAL */
```

```
vsock_stream_connect() /* VMADDR_CID_HYPERVISOR */
lock_sock() /* locked successfully */
vsock_assign_transport()
vsock_deassign_transport()
virtio_transport_destruct()
kfree(virtio_vsock_sock)
vsk->transport = NULL
release sock()
```

```
Thread 2: setsockopt() for vsock
```

```
vsock_stream_setsockopt()
```

```
transport = vsk->transport
```

```
lock_sock() /* can't lock, waiting */
```

```
/* finally locked successfully, proceed */
vsock_update_buffer_size()
transport->notify_buffer_size()
virtio_transport_notify_buffer_size()
virtio_vsock_sock->buf_alloc = *val /* UAF */
```

#### Using Out-of-date Value From a Local Variable



20 / 54

#### Memory Corruption

- Write-after-free for virtio\_vsock\_sock object
- $\bullet$  The size of this object is  ${\bf 64}$  bytes
- This object lives in kmalloc-64 slab cache
- $\bullet$  The buf\_alloc field has type u32 and resides at offset 40
- The value written buf\_alloc is controlled by the attacker
- Four controlled bytes are written to the freed memory

## Fuzzing Miracle (1)

• syzkaller didn't manage to reproduce this crash

- I had to develop the reproducer manually
- But why did the fuzzer fail to do that?
- Looking at vsock\_update\_buffer\_size() code gives the answer:

```
if (val != vsk->buffer_size &&
```

```
transport && transport->notify_buffer_size)
```

```
transport->notify_buffer_size(vsk, &val);
```

```
vsk->buffer_size = val;
```

# Fuzzing Miracle (2)

- For memory corruption, setsockopt() should be called with different SO\_VM\_SOCKETS\_BUFFER\_SIZE value each time
- A fun hack from my first reproducer:

## Fuzzing Miracle (3)

- Upstream syzkaller doesn't do things like that
- Syscall params are chosen when syzkaller generates fuzzing inputs
- Inputs don't change when the fuzzer executes them on the target
- I still don't completely understand how syzkaller got this crash
- $\overline{\mathbf{.}}$

24 / 54

 syzkaller did some lucky multithreaded magic with vsock buffer size limits but then failed to reproduce it



# NOW ABOUT EXPLOITATION, STEP BY STEP

25 / 54

- I've chosen Fedora 33 Server as the exploitation target
- The kernel version: 5.10.11-200.fc33.x86\_64
- I had a goal to bypass SMEP and SMAP
- Bypassing KASLR is included, of course

26 / 54

#### Four Bytes of Power

#### Write-after-free of a 4-byte controlled value to a 64-byte kernel object at offset 40

- That's quite limited memory corruption
- I had a hard time turning it into a real weapon



Here and further I use images of the artifacts from the State Hermitage Museum in Russia. I love this wonderful museum!

#### Heap Spraying Requirements

- I started to work on stable heap spraying
- The exploit should perform some userspace activity that makes the kernel allocate another 64-byte object at the location of freed virtio vsock sock
- 4-byte write-after-free should corrupt the sprayed object instead of unused free kernel memory

#### Experimental Heap Spraying

- I made quick experimental spraying with add\_key syscall
- I called add\_key several times right after the second connect() to vsock while a parallel thread finishes the corrupting setsockopt()
- ftrace allowed to confirm that the freed virtio\_vsock\_sock is overwritten
- I saw that successful heap spraying was possible
- The next step: finding a 64-byte kernel object that can provide a stronger exploit primitive when it has four corrupted bytes at offset 40
- Huh, not so easy!

#### The iovec Technique is Useless Here

• I tried iovec technique from the Bad Binder by Maddie Stone and Jann Horn

A carefully corrupted iovec object can be used

for arbitrary read/write

- No, I got triple fail with this idea:
  - **③** 64-byte iovec is allocated on the kernel stack, not the heap
  - Sour bytes at offset 40 overwrite iovec.iov\_len, not iovec.iov\_base
  - This iovec exploitation trick is dead since the Linux kernel version 4.13, awesome Al Viro killed it with the commit 09fc68dc66f7597b in June 2017

#### Searching for a Special Kernel Object

- I had exhausting experiments with various kernel objects suitable for heap spraying
- I found msgsnd() syscall that creates struct msg\_msg in the kernelspace:

/* message header */					
<pre>struct msg_msg {</pre>					
struct list_head	m_list;	/*	0	16	*/
long int	m_type;	/*	16	8	*/
size_t	m_ts;	/*	24	8	*/
<pre>struct msg_msgseg *</pre>	next;	/*	32	8	*/
void *	security;	/*	40	8	*/
};					
/* message data follows	*/				

• If struct msgbuf in the userspace has 16-byte mtext, the corresponding msg\_msg is

created in kmalloc-64 slab cache, just like virtio\_vsock\_sock!

Alexander Popov

#### Four Bytes of Power

• The 4-byte write-after-free can corrupt the void \*security pointer at offset 40:

/* message header */					
<pre>struct msg_msg {</pre>					
struct list_head	m_list;	/*	0	16	*/
long int	m_type;	/*	16	8	*/
size_t	m_ts;	/*	24	8	*/
<pre>struct msg_msgseg *</pre>	next;	/*	32	8	*/
void *	security;	/*	40	8	*/
};					
/* message data follows	*/				

• Jokingly, I used this security field to break Linux security

 $\overline{\mathbf{\cdot}}$ 

#### Arbitrary Free

- msg\_msg.security points to the kernel data allocated by lsm\_msg\_msg\_alloc()
- It is used by SELinux in the case of Fedora
- It is freed by security\_msg\_msg\_free() when msg\_msg is received
- Corrupting 4 least significant bytes of msg\_msg.security provides arbitrary free!
- That is a much stronger exploit primitive



- After achieving arbitrary free I started to think about where to aim it
- And here I used the trick from my <u>CVE-2019-18683 exploit</u>:
  - > Second connect() to vsock calls vsock\_deassign\_transport()
  - It sets vsk->transport to NULL
  - That makes the vulnerable setsockopt() hit the kernel warning
  - It happens in virtio\_transport\_send\_pkt\_info() just after UAF
  - My exploit can parse this kernel warning and extract useful info!

#### Kernel Warning Full of Secrets

WARNING: CPU: 1 PID: 6739 at net/vmw\_vsock/virtio\_transport\_common.c:34 . . . CPU: 1 PID: 6739 Comm: racer Tainted: G W 5.10.11-200.fc33.x86 64 #1 Hardware name: DEMU Standard PC (D35 + TCH9, 2009), BIDS 1,13,0-2,fc32 04/01/2014 RIP: 0010:virtio\_transport\_send\_pkt\_info+0x14d/0x180 [vmw\_vsock\_virtio\_transport\_common] . . . RSP: 0018:ffffc90000d07e10 EFLAGS: 00010246 RAX: 00000000000000 RBX: ffff888103416ac0 RCX: ffff88811e845b80 RDX: 00000000ffffffff BSI: ffffc90000d07e58 RDI: ffff888103416ac0 RBP: 00000000000000 R08: 000000052008af R09: 000000000000000 R10: 00000000000126 R11: 0000000000000 R12: 000000000000000 R13: ffffc90000d07e58 R14: 0000000000000 R15: ffff888103416ac0 00007f2f123d5640(0000) GS:ffff88817bd00000(0000) knlGS:0000000000000000 FS: CS 0010 DS: 0000 ES: 0000 CR0: 000000080050033 CR2: 00007f81ffc2a000 CR3: 000000011db96004 CR4: 000000000370ee0 Call Trace: virtio\_transport\_notify\_buffer\_size+0x60/0x70 [vmw\_vsock\_virtio\_transport\_common] vsock\_update\_buffer\_size+0x5f/0x70 [vsock] vsock stream setsockopt+0x128/0x270 [vsock]

#### Kernel Infoleak

- A quick debugging session with gdb showed that:
  - RCX contains the kernel address of the freed virtio\_vsock\_sock
  - RBX contains the kernel address of vsock\_sock
- On Fedora, unprivileged users can open and parse /dev/kmsg
- If one more warning arrives at the kernel log, the exploit won one more race
- The exploit can parse the kernel log and get the addresses from the registers



My further exploitation plan was to use arbitrary free for use-after-free:

- Free some object at the address that leaked in the kernel warning
- Perform heap spraying to overwrite that object with controlled data
- Get more power using the corrupted object

37 / 54

- Arbitrary free for vsock\_sock address (from RBX) is useless
- It lives in a dedicated slab cache where I can't do heap spraying
- So I invented how to exploit use-after-free on msg\_msg (from RCX)
- For overwriting msg\_msg I used wonderful setxattr() & userfaultfd() heap spraying technique by Vitaly Nikolenko

#### Arbitrary Read with msg\_msg: Part 1

#### Original struct msg\_msg

struct list\_head m\_list = 0xffff8881XXXXXXX;

long int m\_type = 1;

size\_t m\_ts = 16;

struct msg\_msgseg \*next = NULL;

void \*security = 0xffff8881YYYYYYY;

msg\_msg data

#### Overwritten struct msg\_msg

struct list\_head m\_list = 0xa5a5a5a5a5a5a5a5;

long int m\_type = 0x1337;

size\_t m\_ts = 6096;

struct msg\_msgseg \*next = 0xffff8881ZZZZZZ;

void \*security = 0xffff8881YYYYYYY;

msg\_msg data

kernel data for reading

Fake struct msg\_msgseg

kernel data for reading

#### Arbitrary Read with msg\_msg: Part 2

- Receiving this crafted msg\_msg manipulates the System V message queue
- That breaks the kernel because the msg\_msg.m\_list pointer is invalid 😕
- msgrcv() documentation for the win!
- MSG\_COPY flag allows fetching a copy of the message nondestructively in the message nondestructively in the message nondestructively is a set of the message nondestructively in the message nondestructively is a set of the message nondestructively in the message nondestructively is a set of the message nondestructively in the message nondestructively is a set of the message nondestructively is a set of the message nondestructively in the message nondestructively is a set of the message nondest nondestructively is a set of the message nondestructively i



## Exploiting Arbitrary Read (1)

- 1. Get the kernel address of a good msg\_msg
  - win the race on a virtual socket
  - call spraying msgsnd() after the memory corruption
  - parse /dev/kmsg and get the kernel address of this good msg\_msg from RCX
  - $\bullet$  also, save the kernel address of <code>vsock\_sock</code> from <code>RBX</code>

#### Exploiting Arbitrary Read (2)

2. Execute arbitrary free against good msg\_msg using a corrupted msg\_msg



## Exploiting Arbitrary Read (3)

#### 3. Overwrite good msg\_msg with controlled data using setxattr() & userfaultfd()



Alexander Popov

## Exploiting Arbitrary Read (4)

4. Read vsock\_sock to the userspace using msgrcv() for the overwritten msg\_msg



#### Sorting the Loot

That's what I found inside the vsock\_sock kernel object:

- Plenty of pointers to objects from dedicated slab caches 😕
- struct mem\_cgroup \*sk\_memcg pointer at offset 664
  - mem\_cgroup objects live in the kmalloc-4k slab cache
  - ▶ I tried to call kfree() for it and the kernel panicked instantly 😕
- S const struct cred \*owner pointer at offset 840
  - ▶ It points to the credentials that I want to overwrite for privilege escalation
  - It's a pity that cred lives in dedicated cred\_jar slab cache
- void (\*sk\_write\_space)(struct sock \*) function pointer at offset 688
  - It is set to the address of sock\_def\_write\_space() kernel function
  - $\blacktriangleright$  That can be used for calculating the KASLR offset  $\overline{oldsymbol{\odot}}$

#### Good Old Trick with sk\_buff

- I used it in my exploit for CVE-2017-2636 in the Linux kernel
- I turned double free for a kmalloc-8192 object into use-after-free on sk\_buff
- I decided to repeat that trick
  - ► A network-related buffer in the kernel is represented by sk\_buff
  - This object has skb\_shared\_info with destructor\_arg
  - Creating a 2800-byte network packet in the userspace will make skb shared info be allocated in the kmalloc-4k slab cache
  - That's where mem\_cgroup objects live as well!

#### Use-after-free on sk\_buff

- Create one client socket and 32 server sockets (for AF\_INET, SOCK\_DGRAM, IPPROTO\_UDP)
- Send a 2800-byte buffer filled with 0x42 to each server socket using sendto()
- Perform arbitrary read for vsock\_sock (described earlier)
- Calculate possible sk\_buff kernel address as sk\_memcg plus 4096 (the next element in kmalloc-4k)
- Perform arbitrary read for this possible sk\_buff address
- If 0x42 bytes are found, perform arbitrary free against the sk\_buff
- Otherwise, add 4096 to the possible sk\_buff address and go to step 5

#### The Payload for Overwriting skb\_shared\_info



#### Control Flow Hijack

• I didn't manage to find a stack pivoting gadget in vmlinuz-5.10.11-200.fc33.x86\_64

that can work in my restrictions

- So I performed arbitrary write in one shot
- SMEP and SMAP protection is bypassed!

```
/*
 * A single ROP gadget for arbitrary write:
 * mov rdx, qword ptr [rdi + 8] ; mov qword ptr [rdx + rcx*8], rsi ; ret
 * Here rdi stores uinfo_p address, rcx is 0, rsi is 1
 */
uinfo_p->callback = ARBITRARY_WRITE_GADGET + kaslr_offset;
uinfo_p->desc = owner_cred + CRED_EUID_EGID_OFFSET; /* value for "qword ptr [rdi + 8]" */
uinfo_p->desc = uinfo_p->desc - 1; /* rsi value 1 should not get into euid */
```

#### Arbitrary Write Using skb\_shared\_info

This weapon is used twice to get root privileges:

- Write zeros to effective uid and gid
- Write zeros to uid and gid



#### Demo Time



#### Possible Exploit Mitigation

- Exploiting this vulnerability is impossible with the Linux kernel heap quarantine
  - ▶ Because this memory corruption happens very shortly after the race condition
  - See the <u>article</u> about my SLAB\_QUARANTINE prototype
- Against kernel module autoloading by unprivileged users grsecurity MODHARDEN
- Against userfaultfd() abuse setting /proc/sys/vm/unprivileged\_userfaultfd to 0
- $\bullet$  Against infoleak via kernel log setting kernel.dmesg\_restrict sysctl to 1
- Against calling my ROP gadget -
  - Control Flow Integrity (see the technologies on my Linux Kernel Defence Map)
- Against use-after-free (hopefully in the future) -

ARM Memory Tagging Extension (MTE) support for the kernel, on ARM

• [rumors] Against heap spraying -

grsecurity Wunderwaffe called AUTOSLAB (we don't know much about it)

#### Conclusion

 Investigating and fixing CVE-2021-26708, developing the PoC exploit, and preparing this talk was a big deal for me

• I hope you enjoyed it!



- I managed to turn the race condition with a very limited memory corruption into arbitrary read/write for the Linux kernel memory
- I will publish a large and detailed write-up very soon

#### Thanks! Send me your questions!

alex.popov@linux.com @a13xp0p0v

http://blog.ptsecurity.com/ @ptsecurity

POSITIVE TECHNOLOGIES