Linux Kernel Self-Protection Project (and how we're trying to make C less dangerous)

> UC Santa Cruz, Open Source Programming February 26, 2019

> > Kees ("Case") Cook keescook@chromium.org @kees\_cook

https://outflux.net/slides/2019/uc/kspp.pdf

## Agenda

- Background
  - Kernel Self Protection Project
  - C as a fancy assembler
- Towards less dangerous C
  - Variable Length Arrays are bad and slow
  - Explicit switch case fall-through
  - Always-initialized automatic variables
  - Arithmetic overflow detection
  - Hope for bounds checking
  - Control Flow Integrity: forward edges
  - Control Flow Integrity: backward edges
  - Where are we now?
  - How you can help

#### 000-005 Unexplained Phenomena, Software Programming

## Kernel Security for *this* talk is ...

- More than access control (e.g. SELinux)
- More than attack surface reduction (e.g. seccomp)
- More than bug fixing (e.g. CVEs)
- More than protecting userspace (e.g. namespaces)
- More than kernel integrity/roots of trust (e.g. TPM)
- This is about Kernel Self Protection



## What needs securing?

- Servers, laptops, cars, phones, TVs, space stations, ...
- >2,000,000,000 active Android devices in 2017
  - Majority are running v3.10 (with v3.18 slowly catching up)
- Bug lifetimes are even longer than upstream
- "Not our problem"? Even if upstream fixes every bug found, and the fixes are magically sent to devices, bug lifetimes are still huge.



## **Upstream Bug Lifetime**

- In 2010 Jon Corbet researched security flaw fixes with CVEs, and found that the average time between introduction and fix was about 5 years.
- My analysis of Ubuntu CVE tracker for the kernel from 2011 through 2018 crept closer to 6 years for a while, but has now started to diminish:
  - Critical: 3 at 5.3 years average
  - High: 79 at 5.6 years average
  - Medium: 691 at 5.9 years average
  - Low: 349 at 6.2 years average



#### critical & high CVE lifetimes



#### Attackers are watching

- The risk is not theoretical. Attackers are watching commits, and they are better at finding bugs than we are:
  - http://seclists.org/fulldisclosure/2010/Sep/268
- Most attackers are not publicly boasting about when they found their 0-day...



# Bug fighting continues

- We're finding them
  - Static checkers: gcc, Clang, Coccinelle, Smatch, sparse, Coverity
  - Dynamic checkers: kernel, KASan-family, syzkaller, trinity
- We're fixing them
  - Ask Greg KH how many patches land in -stable
- They'll always be around
  - We keep writing them
  - They exist whether we're aware of them or not
  - Whack-a-mole is not a solution



## Analogy: 1960s Car Industry

- Konstantin Ryabitsev's keynote at 2015 Linux Security Summit
  - http://kernsec.org/files/lss2015/giant-bags-of-mostly-water.pdf
- Cars were designed to run, not to fail
- Linux now where the car industry was in 1960s
  - https://www.youtube.com/watch?v=fPF4fBGNK0U
- We must handle failures (attacks) safely
  - Userspace is becoming difficult to attack
  - Containers paint a target on the kernel
  - Lives depend on Linux



## Killing bugs is nice

- Some truth to security bugs being "just normal bugs"
- Your security bug may not be my security bug
- We have little idea which bugs most attackers use
- Bug might be in out-of-tree code
  - Un-upstreamed vendor drivers
  - Not an excuse to claim "not our problem"



## Killing bug classes is better

- If we can stop an entire kind of bug from happening, we absolutely should do so!
- Those bugs never happen again
- Not even out-of-tree code can hit them
- But we'll never kill all bug classes



## Killing exploitation is best

- We will always have bugs
- We must stop their exploitation
- Eliminate exploitation targets and methods
- Eliminate information exposures
- Eliminate anything that assists attackers
- Even if it makes development more difficult



# **Typical Exploit Chains**

- Modern attacks tend to use more than one flaw
- Need to know where targets are
- Need to inject (or build) malicious code
- Need to locate malicious code
- Need to redirect execution to malicious code



#### What can we do?

- Many exploit mitigation technologies already exist (e.g. Grsecurity/PaX) or have been researched (e.g. academic whitepapers), but are not present in the upstream Linux kernel
- There is demand for kernel self-protection, and there is demand for it to exist in the upstream kernel
- http://www.washingtonpost.com/sf/business/2015/11/05/net-of-in security-the-kernel-of-the-argument/



## Kernel Self Protection Project

- https://kernsec.org/wiki/index.php/Kernel\_Self\_Protection\_Project
- KSPP focuses on the kernel protecting the *kernel* from attack (e.g. refcount overflow) rather than the kernel protecting *userspace* from attack (e.g. execve brute force detection) but any area of related development is welcome
- Currently ~12 organizations and ~10 individuals working on about ~20 technologies
- Slow and steady



## **Probabilistic protections**

- Protections that derive their strength from some system state being unknown to an attacker
- Weaker than "deterministic" protections since information exposures can defeat them, though they still have real-world impact
- Familiar examples:
  - stack protector (canary value can be exposed)
  - Address Space Layout Randomization (offset can be exposed)



### **Deterministic protections**

- Protections that derive their strength from organizational system state that always blocks attackers
- Familiar examples:
  - Read-only memory (writes will fail)
  - Bounds-checking (large accesses fail)
- Using a memory-safe langauge...
  - Rewrite the kernel in Rust? :(
  - Let's take some incremental steps



#### C as a fancy assembler: **almost machine code**

- The kernel wants to be as fast and small as possible
- At the core, kernel wants to do very architecture-specific things for memory management, interrupt handling, scheduling, ...
- No C API for setting up page tables, switching to 64-bit mode ...

```
/* Enable the boot page tables */
leal pgtable(%ebx), %eax
movl %eax, %cr3

/* Enable Long mode in EFER (Extended Feature Enable Register) */
movl $MSR_EFER, %ecx
rdmsr
btsl $_EFER_LME, %eax
wrmsr
```

## C as a fancy assembler: undefined behavior

- The C langauge comes with some operational baggage, and weak "standard" libraries
  - What are the contents of "uninitialized" variables?
    - ... whatever was in memory from before now!
  - void pointers have no type yet we can call typed functions through them?
    - ... assembly doesn't care: everything can be an address to call!
  - Why does memcpy() have no "max destination length" argument?
    - ... just do what I say; memory areas are all the same!
- "With undefined behavior, anything is possible!"
  - https://raphlinus.github.io/programming/rust/2018/08/17/undefined-behavior.html



#### Variable Length Arrays are **bad**

- Exhaust stack, linear overflow: write to things following it
- Jump over guard pages and write to things following it
- But easy to find with compiler flag: -Wvla



### Variable Length Arrays are **slow**

- This seems conceptually sound: more instructions to change stack size, but it seems like it would be hard to measure.
- It is quite measurable ... 13% speed up measured during lib/bch.c VLA removal:

https://git.kernel.org/linus/02361bc77888 (Ivan Djelic)

Buffer allocation		Encoding throughput	(Mbit/s)
on-stack, VLA		3988	
on-stack, fixed	I	4494	
kmalloc		1967	



## Switch case fall-through: did I mean it?

- CWE-484 "Omitted Break Statement in Switch"
- Semantic weakness in C ("switch" is just assembly test/jump...)
- Commit logs with "missing break statement": 67



#### Switch case fall-through: new "statement"

- Use -Wimplicit-fallthrough to add a new switch "statement"
  - Actually a comment, but is parsed by compilers now, following the lead of static checkers
- Mark all non-breaks with a "fall through" comment, for example https://git.kernel.org/linus/4597b62f7a60 (Gustavo A. R. Silva)

## Always-initialized local variables: just do it

- CWE-200 "Information Exposure", CWE-457 "Use of Uninitialized Variable"
- gcc -finit-local-vars not upstream
- Clang -fsanitize=init-local not upstream
- CONFIG\_GCC\_PLUGIN\_...
  - STRUCTLEAK (for structs with \_\_user pointers)
  - STRUCTLEAK\_BYREF (when passed into funcs)
  - Soon, plugin to mimic
     finit-local-vars too

From: Linus Torvalds <torvalds@linux-foundation.org> Subject: Re: Fully initialized stack usage

On Tue, Feb 27, 2018 at 3:15 AM, P J P <ppandit@redhat.com> wrote: > ...

> This experimental patch by Florian Weimer(CC'd) adds an option

- > '-finit-local-vars' to gcc(1) compiler. When a program(or kernel)
- > is built using this option, its automatic(local) variables are > initialized with zero(0). This could significantly reduce the kernel
- > initialised with zero(0). This could significantly reduce the kernel > information leakage issues.

Oh, I love that patch.

THAT is the kind of thing we should do. It's small, it's trivial, and it's done early in the parsing stage, so later stages will almost certainly end up optimizing things away.

• •

#### Always-initialized local variables: switch gotcha

warning: statement will never be executed [-Wswitch-unreachable]

enum pipe pipe = crtc->pipe; int sprite0\_start, sprite1\_start; uint32\_t dsparb, dsparb2, dsparb3; switch (pipe) { uint32\_t dsparb, dsparb2, dsparb3; case PIPE\_A: dsparb = I915\_READ(DSPARB); dsparb2 = I915\_READ(DSPARB2);

#### Arithmetic overflow detection: gcc?

- gcc's -fsanitize=signed-integer-overflow (CONFIG\_UBSAN)
  - Only signed. Fast: in the noise. Big: warnings grow kernel image by 6% (aborts grow it by 0.1%)
- But we can use explicit single-operation helpers. To quote Rasmus Villemoes:

```
So is it worth it? I think it is, if nothing else for the documentation
value of seeing

if (check_add_overflow(a, b, &d))
   return -EGOAWAY;
   do_stuff_with(d);

instead of the open-coded (and possibly wrong and/or incomplete and/or
UBsan-tickling)

if (a+b < a)
   return -EGOAWAY;
   do_stuff_with(a+b);</pre>
```

## Arithmetic overflow detection: Clang:)

• Clang can do signed and unsigned instrumentation:

-fsanitize=signed-integer-overflow

-fsanitize=unsigned-integer-overflow

```
$ clang overflow.c -fsanitize=signed-integer-overflow && ./a.out
overflow.c:11:12: runtime error: signed integer overflow: 1 + 2147483647 cannot be represented in type 'int'
-2147483648
```

```
-ftrap-function=abort && ./a.out
```

zsh: abort (core dumped) ./a.out

## Bounds checking: explicit checking is slow :(

- Explicit checks for linear overflows of SLAB objects, stack, etc
  - copy\_{to,from}\_user() checking: <~1% performance hit</pre>
  - strcpy()-family checking: ~2% performance hit
  - memcpy()-family checking: ~1% performance hit
- Can we get better APIs?
  - strcpy() is terrible
  - sprintf() is bad
  - memcpy() is weak

## Instead of strcpy(): strscpy()

- strcpy() no bounds checking on destination nor source!
- strncpy() doesn't always NUL terminate (good for non-C-strings, does NUL pad destination) char dest[4]; strncpy(dest, "ohai!", sizeof(dest)); /\* unhelpfully returns dest \*/ dest: "o", "h", "a", "i" ... no trailing NUL byte :(
- strlcpy() reads source beyond max destination size (returns length of source!)
- strscpy() safest (returns bytes copied, not including NUL, or -E2BIG)
   ssize\_t count = strscpy(dest, "ohai!", sizeof(dest)); /\* returns -E2BIG \*/
   dest: "o", "h", "a", NUL
  - Does not NUL-pad destination ... if desired, add explicit memset() (kernel needs a helper for this...)
    if (count > 0 && count + 1 < sizeof(dest))</pre>

```
memset(dest + count + 1, 0, sizeof(dest) - count - 1);
```

## Instead of sprintf(): scnprintf()

- sprintf() no bounds checking on destination!
- snprintf() always NUL-terminates, but returns how much it
  WOULD have written :(

```
int count = snprintf(buf, sizeof(buf), fmt..., ...);
```

```
for (i = 0; i < something; i ++)
```

count += snprintf(buf + count, sizeof(buf) - count, fmt..., ...); copy\_to\_user(user, buf, count);

- scnprintf() always NUL-terminates, returns count of bytes copied
  - Replace in above code!

#### Instead of memcpy(): uhhh ... be ... careful?

• memcpy() has no sense of destination size :(

```
uint8_t bytes[128];
size t wanted, copied = 0;
```

```
for (i = 0; i < something && copied < sizeof(bytes); i ++) {
   wanted = ...;
   if (wanted > sizeof(bytes) - copied)
      wanted = sizeof(bytes) - copied;
   memcpy(bytes + copied, wanted, source);
   copied += wanted;
}
```

# Bounds checking: memory tagging :)

- Hardware memory tagging/coloring is much faster!
  - SPARC Application Data Integrity (ADI)
  - ARMv8.5 Memory Tagging Extension (MTE)



## Control Flow Integrity: indirect calls

 With memory W^X, gaining execution control needs to change function pointers saved in heap or stack, where all type information was lost!



## CFI, forward edges: just call pointers :(

void call\_one(char \*input)

printf("Printing stuff: %s\n", input);

void call\_two(void)

```
printf("Eek: don't run me\n");
```

```
int main(int argc, char *argv[])
```

void (\*func)(char \*) = call\_ong

```
if (atoi(argv[1]) < 0)
    func = (void *)call_two;</pre>
```

func(argv[0]);

return 0;

Ignore function prototype ...

Normally just a call to a memory address:

```
$ clang demo.c -o demo
$ ./demo 1
Printing stuff: ./demo
$ ./demo -1
Eek: don't run me
$
```

## CFI, forward edges: enforce prototype :)

void call\_one(char \*input)

printf("Printing stuff: %s\n", input);

void call\_two(void)

```
printf("Eek: don't run me\n");
```

```
int main(int argc, char *argv[])
```

```
void (*func)(char *) = call_on
```

```
if (atoi(argv[1]) < 0)
    func = (void *)call_two;</pre>
```

```
func(argv[0]);
```

return 0;

Ignore function prototype ...

#### Clang -fsanitize=cfi will check at runtime:

\$ clang demo.c -o demo -flto -fvisibility=hidden -fsanitize=cfi
\$ ./demo 1
Printing stuff: ./demo
\$ ./demo -1
Illegal instruction (core dumped)

#### CFI, backward edges: two stacks

- Clang's Safe Stack
  - Clang: -fsanitize=safe-stack



## CFI, backward edges: shadow call stack

- Clang's Shadow Call Stack
  - Clang: -fsanitize=shadow-call-stack
  - Results in two stack registers: sp and unspilled x18



## CFI, backward edges: hardware support

- Intel CET: hardware-based read-only shadow call stack
  - Implicit use of otherwise read-only shadow stack during call and ret instructions
- ARM v8.3a Pointer Authentication ("signed return address")
  - New instructions: paciasp and autiasp
  - Clang and gcc: -msign-return-address

+paciasp		
stp	x29, x30,	[sp, #-48]!
mov	x29, sp	
str	w0, [x29,	#28]
mov	w0, #0x0	
ldp	x29, x30,	[sp], #48
+autiasp		
ret		

## Where is the Linux kernel now?

- Variable Length Arrays
  - Finally eradicated from kernel since v4.20 (Dec 2018)!
- Explicit switch case fall-through
  - Steady progress on full markings (232 of 2311 remain)
- Always-initialized automatic variables
  - Various partial options, needs more compiler work
- Arithmetic overflow detection
  - Memory allocations now doing explicit overflow detection
  - Needs better kernel support for Clang and gcc
- Bounds checking
  - Crying about performance hits
  - Waiting (im)patiently for hardware support
- Control Flow Integrity: forward edges
  - Need Clang LTO support in kernel, but works on Android
- Control Flow Integrity: backward edges
  - Shadow Call Stack works on Android
  - Waiting (im)patiently for hardware support



#### **Challenges in Kernel Security Development**

**Cultural**: Conservatism, Responsibility, Sacrifice, Patience **Technical**: Complexity, Innovation, Collaboration **Resource**: Dedicated Developers, Reviewers, Testers, Backporters



#### Thoughts?

Kees ("Case") Cook keescook@chromium.org keescook@google.com kees@outflux.net

https://outflux.net/slides/2019/uc/kspp.pdf

http://www.openwall.com/lists/kernel-hardening/ http://kernsec.org/wiki/index.php/Kernel\_Self\_Protection\_Project