

# Splitting the Linux Kernel for Fun & Profit

Chris I Dalton [cid@hp.com](mailto:cid@hp.com)\*

HP Labs, Bristol UK



# Splitting the Linux Kernel for Fun & Profit

‘ How to hack a Micro-Kernel interface into Linux’

# Splitting the Linux Kernel for Fun & Profit

‘How to hack a Micro-Kernel interface into Linux’

OR

‘Adding Intra-Kernel protection to Linux using silicon-based Virtualization Extensions’

(Without needing a Hypervisor)

# Outline

- **Part 1: Motivation for the work**
  - If I wanted a secure OS I wouldn't start from Linux...but what if that was all you had?
- Part 2: Background
  - Ways to structure Operating Systems
  - Linux Kernel Weaknesses & Split-Kernel Demo Video
- Part 3: Splitting the Kernel
  - Restructuring Linux using HW Virtualization Support
  - High-level & Some code details
  - Performance & Invasiveness
- Part 4: Current Status, Opportunities & Futures
  - Open Source Links

# Outline

- Part 1: Motivation for the work
  - If I wanted a secure OS I wouldn't start from Linux...but what if that was all you had?
- **Part 2: Background**
  - Ways to structure Operating Systems
  - Linux Kernel Weaknesses & Split-Kernel Demo Video
- Part 3: Splitting the Kernel
  - Restructuring Linux using HW Virtualization Support
  - High-level & Some code details
  - Performance & Invasiveness
- Part 4: Current Status, Opportunities & Futures
  - Open Source Links

# Outline

- Part 1: Motivation for the work
  - If I wanted a secure OS I wouldn't start from Linux...but what if that was all you had?
- Part 2: Background
  - Ways to structure Operating Systems
  - Linux Kernel Weaknesses & Split-Kernel Demo Video
- **Part 3: Splitting the Kernel**
  - Restructuring Linux using HW Virtualization Support
  - High-level & Some code details
  - Performance & Invasiveness
- Part 4: Current Status, Opportunities & Futures
  - Open Source Links

# Outline

- Part 1: Motivation for the work
  - If I wanted a secure OS I wouldn't start from Linux...but what if that was all you had?
- Part 2: Background
  - Ways to structure Operating Systems
  - Linux Kernel Weaknesses & Split-Kernel Demo Video
- Part 3: Splitting the Kernel
  - Restructuring Linux using HW Virtualization Support
  - High-level & Some code details
  - Performance & Invasiveness
- **Part 4: Current Status, Opportunities & Futures**
  - Open Source Links

# Part 1: Original motivation for the work

- Containers offer alternative to using Hypervisors for SW deployments
  - Docker, CoreOS / Rkt, etc



# Original Motivation for the Work

- Containers offer alternative to using Hypervisors for SW deployments
  - Docker, CoreOS / Rkt, etc
- Upside: Lighter-weight than using VMs
  - Only one underlying OS to manage
  - Plus better integration opportunities

# Original Motivation for the Work

- Containers offer alternative to using Hypervisors for SW deployments
  - Docker, CoreOS / Rkt, etc
- Upside: Lighter-weight than using VMs
  - Only one underlying OS to manage
  - Plus better integration opportunities
- Downside: Shared 'host' kernel a significant vulnerability
  - Currently Linux has no **'intra-kernel'** protection
  - All bets are off if you manage to get into the kernel

# Getting into the kernel is not that hard!

E.g. Malware triggering buffer overflows / stack / heap attacks, un-authorized module loading, User-space kernel hijack via code & data redirection (ROP), DMA attacks, etc. even with SMAP / SMEP / PXN support

And of course more recent Spectre / Meltdown / L1TF attacks



What did we want to do?

# What did we want to do?

Reduce the consequences of a Kernel compromise by introducing a degree of **Intra-kernel** Protection into the Linux Kernel

How?

# How?

Restructure kernel into Outer and Inner region based on MMU access

Inner Region can access the MMU directly

Outer Region needs to go through a virtual MMU interface to modify page mappings, etc.

Inner/Outer region separation enforced through CPU HW support for virtualization

# Related Research

- ‘Nested Kernel: An Operating System Architecture for Intra-Kernel Privilege Separation’ (Dautenhahn, et al.)
  - Looks at implementing a vMMU for FreeBSD
  - Relies on privileged instruction removing from kernel binary and code scanning for enforcement not VT-x extensions
- ‘Dune: Safe User-level Access to Privileged CPU Features’ (Belay et al.)
  - Uses Intel Vt-x extensions to safely expose HW to user-space processes (e.g. each process has access to cpu rings 0-3 )
- ‘Address space isolation Inside the Linux Kernel’ (Rapoport, et al 2019)
  - Tries to achieve similar properties and capabilities to our work
  - Uses restricted page table mappings and kernel direct map modifications not VT-x extensions



# What does this buy us?

- Strong control over the integrity of kernel code and core data
  - No 'un-authorized' code gets to run in kernel mode
  - Can protect the integrity of data even against malicious kernel mode code
- Can offer confidentiality guarantees
  - Can protect application secrets even against malicious kernel mode code
- Guard against cross-process code + data compromise
  - Enhanced risk when running multiple 'isolated' containers

# What does this buy us?

- Strong control over the integrity of kernel code and core data
  - No 'un-authorized' code gets to run in kernel mode
  - Can protect the integrity of data even against malicious kernel mode code
- Can offer confidentiality guarantees
  - Can protect application secrets even against malicious kernel mode code
- Guard against cross-process code + data compromise
  - Enhanced risk when running multiple 'isolated' containers

# What does this buy us?

- Strong control over the integrity of kernel code and core data
  - No 'un-authorized' code gets to run in kernel mode
  - Can protect the integrity of data even against malicious kernel mode code
- Can offer confidentiality guarantees
  - Can protect application secrets even against malicious kernel mode code
- Guard against cross-process code + data compromise
  - Enhanced risk when running multiple 'isolated' containers

# Constraints

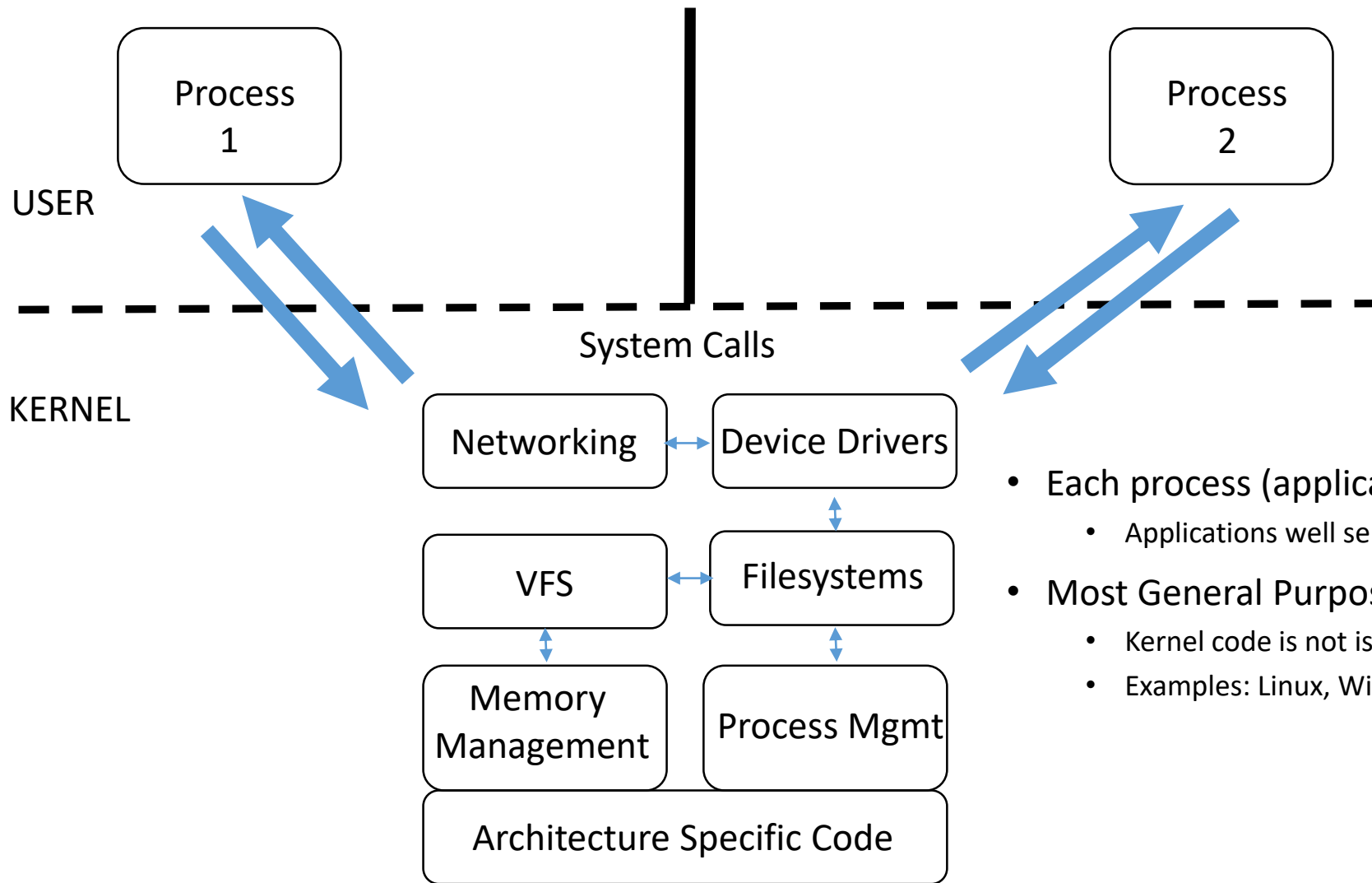
- What no Hypervisor??
  - Strategic control point for others in the Industry (Dell, VMWare, etc)
- Is it still Linux?
  - Can't afford long-term engineering support
  - Needs to be upstream-able
  - Minimize performance overhead c.f. Hypervisors
  - Minimize intrusiveness c.f. L4Linux
  - But still has to offer significant security improvements over 'normal' Linux

# Constraints

- What no Hypervisor??
  - Strategic control point for others in the Industry (Dell, VMWare, etc)
- Is it still Linux?
  - Can't afford long-term engineering support
  - Needs to be upstream-able
  - Minimize performance overhead c.f. Hypervisors
  - Minimize intrusiveness c.f. L4Linux
  - But still has to offer significant security improvements over 'normal' Linux

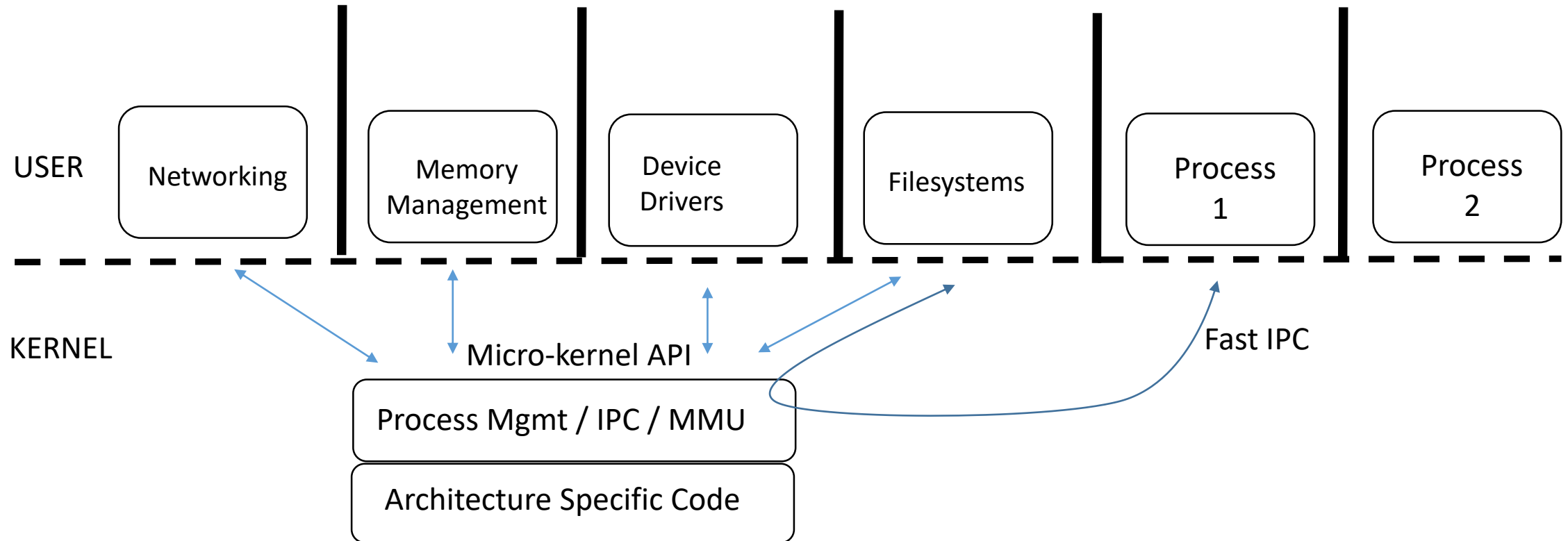
# Part 2: Operating Systems Background

# (Monolithic) Operating Systems



- Each process (application) has its own isolated space
  - Applications well separated
- Most General Purpose Operating Systems are 'Monolithic'
  - Kernel code is not isolated from itself
  - Examples: Linux, Windows

# (Micro-Kernel) Operating Systems



- Examples: Fiasco.OC (L4 Micro-kernel), Composite
- Google Fuchsia/Zircon



# Part 3: Monolithic Linux Kernel (In-)Security Demo

```
Terminal
x - □ root@mn: ~/projects/ckernel/test-okernel-mapping/kvc
./kernel_vuln_client 0xPHYSADDR (in hex) [or set KV_PHYSADDR=0xPHYSADDR as env variable]
root@mn:~/projects/ckernel/test-okernel-mapping/kvc# ./kernel_vuln_client 0x42c568000
Physical memory address to access via exploit:=(0x42c568000)
kernel exploit succeeded: (Something Secret) @ (0x42c568000)
root@mn:~/projects/ckernel/test-okernel-mapping/kvc# ../../okernel-usrsrc/okernel_exec2 /bin/bash
61:okernel_exec2.c arg (/bin/bash)
65:okernel_exec2.c started (../../okernel-usrsrc/okernel_exec2) with pid (25680)
72:okernel_exec2.c About to call OKERNEL_ON_CMD ioctl...
Returned successfully from OKERNEL_ON_CMD ioctl.
Current memory segments:

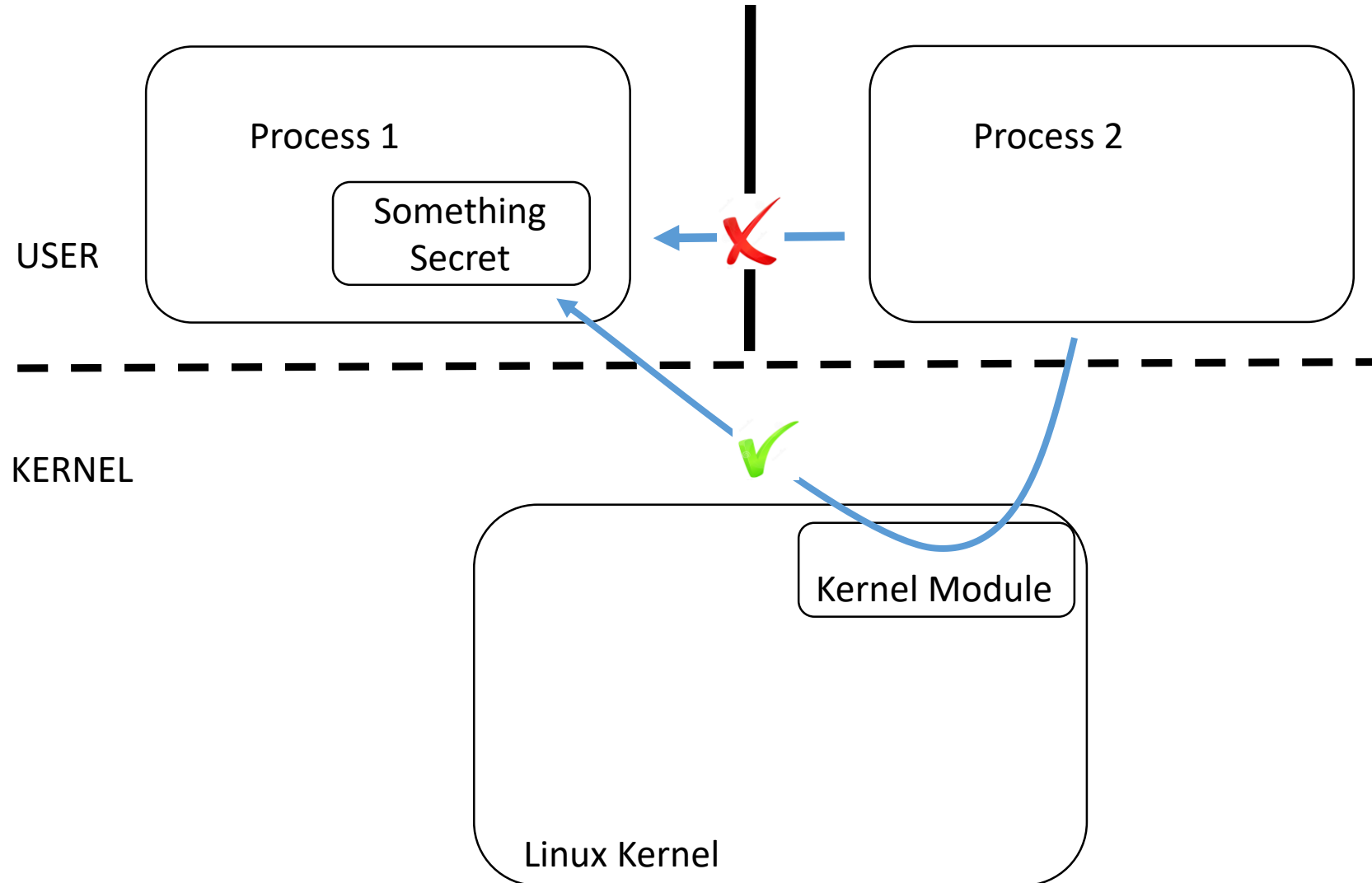
Pid of the process is = 25680
Addresses which fall into:
1) Data segment = 0x602078
2) BSS segment = 0x602080
3) Code segment = 0x400940
4) Stack segment = 0x7ffc36770164
Current memory segments done.
About to exec (/bin/bash)
root@mn:~/projects/ckernel/test-okernel-mapping/kvc#
root@mn:~/projects/ckernel/test-okernel-mapping/kvc#
root@mn:~/projects/ckernel/test-okernel-mapping/kvc#
root@mn:~/projects/ckernel/test-okernel-mapping/kvc#
root@mn:~/projects/ckernel/test-okernel-mapping/kvc#
root@mn:~/projects/ckernel/test-okernel-mapping/kvc# cat /proc/self/okernel
1
root@mn:~/projects/ckernel/test-okernel-mapping/kvc# ./kernel_vuln_client 0x42c568000
Physical memory address to access via exploit:=(0x42c568000)
kernel exploit succeeded: () @ (0x42c568000)
root@mn:~/projects/ckernel/test-okernel-mapping/kvc#

x - □ root@mn: ~/projects/ckernel/test-okernel-mapping/pmc
61:okernel_exec2.c arg (/bin/bash)
65:okernel_exec2.c started (../../okernel-usrsrc/okernel_exec2) with pid (25668)
72:okernel_exec2.c About to call OKERNEL_ON_CMD ioctl...
Returned successfully from OKERNEL_ON_CMD ioctl.
Current memory segments:

Pid of the process is = 25668
Addresses which fall into:
1) Data segment = 0x602078
2) BSS segment = 0x602080
3) Code segment = 0x400940
4) Stack segment = 0x7ffffdccc484
Current memory segments done.
About to exec (/bin/bash)
root@mn:~/projects/ckernel/test-okernel-mapping/pmc#
root@mn:~/projects/ckernel/test-okernel-mapping/pmc#
root@mn:~/projects/ckernel/test-okernel-mapping/pmc#
root@mn:~/projects/ckernel/test-okernel-mapping/pmc#
root@mn:~/projects/ckernel/test-okernel-mapping/pmc#
root@mn:~/projects/ckernel/test-okernel-mapping/pmc# ./protected_memory_client
Starting:=(./protected_memory_client) pid:=(25696)
Phys addr of allocated protected page:=(0x42c568000)
Data to write to protected page:=(Something Secret)
Written data to protected page ok.
Read back protected page data:=(Something Secret)
Just wait now until we are quit...

x - □ root@mn: ~/projects/ckernel/test-okernel-mapping
# entries-in-buffer/entries-written: 57692/57692 #P:4
#
#          -----=> irqsoft
#          -----=> need_resched
#          -----=> hardirq/softirq
#          -----=> preempt-depth
#          delay
#          TASK-PID CPU#  |||||  TIMESTAMP  FUNCTION
#          -----|-----|-----|-----|-----|-----|
# swapper/0-1 [000] .... 3.727002: okernel_init: 03NR: cpu(0) pid(1) okernel_init: St
art initialization...
# swapper/0-1 [000] .... 3.727016: okernel_init: 03NR: cpu(0) pid(1) okernel_init: Cr
eating Device Major_no : 248
# swapper/0-1 [000] .... 3.727208: okernel_init: 03NR: cpu(0) pid(1) okernel_init: De
vice <okernel> initialized in kernel.
# swapper/0-1 [000] .... 3.727211: vmx_init: 03NR: cpu(0) pid(1) e820_end_paddr: last
_pfn = 0x44ee00 max_arch_pfn = 0x400000000
# swapper/0-1 [000] .... 3.727582: okernel_dump_stack info: 03R : cpu(0) pid(1) okern
el_dump_stack_info: thread/stack size (16384) thread_info* (0xffff98caf6e78000) stack in-use (0x1e0
) (480)
# swapper/0-1 [000] .... 3.727583: okernel_dump_stack info: 03R : cpu(0) pid(1) okern
el_dump_stack_info: stack sp0 (0xfffffadca80028000) current sp (0xfffffadca80027e20) end stack (0xffff
fadca80024000)
# swapper/0-1 [000] .... 3.727583: okernel_init: 03R : cpu(0) pid(1) okernel_init: En
abled, initialization done.
# bash-2622 [000] .... 178.339658: ok_device_open: 03R : cpu(0) pid(2622) ok_device_o
pen: Opening device <okernel>
```

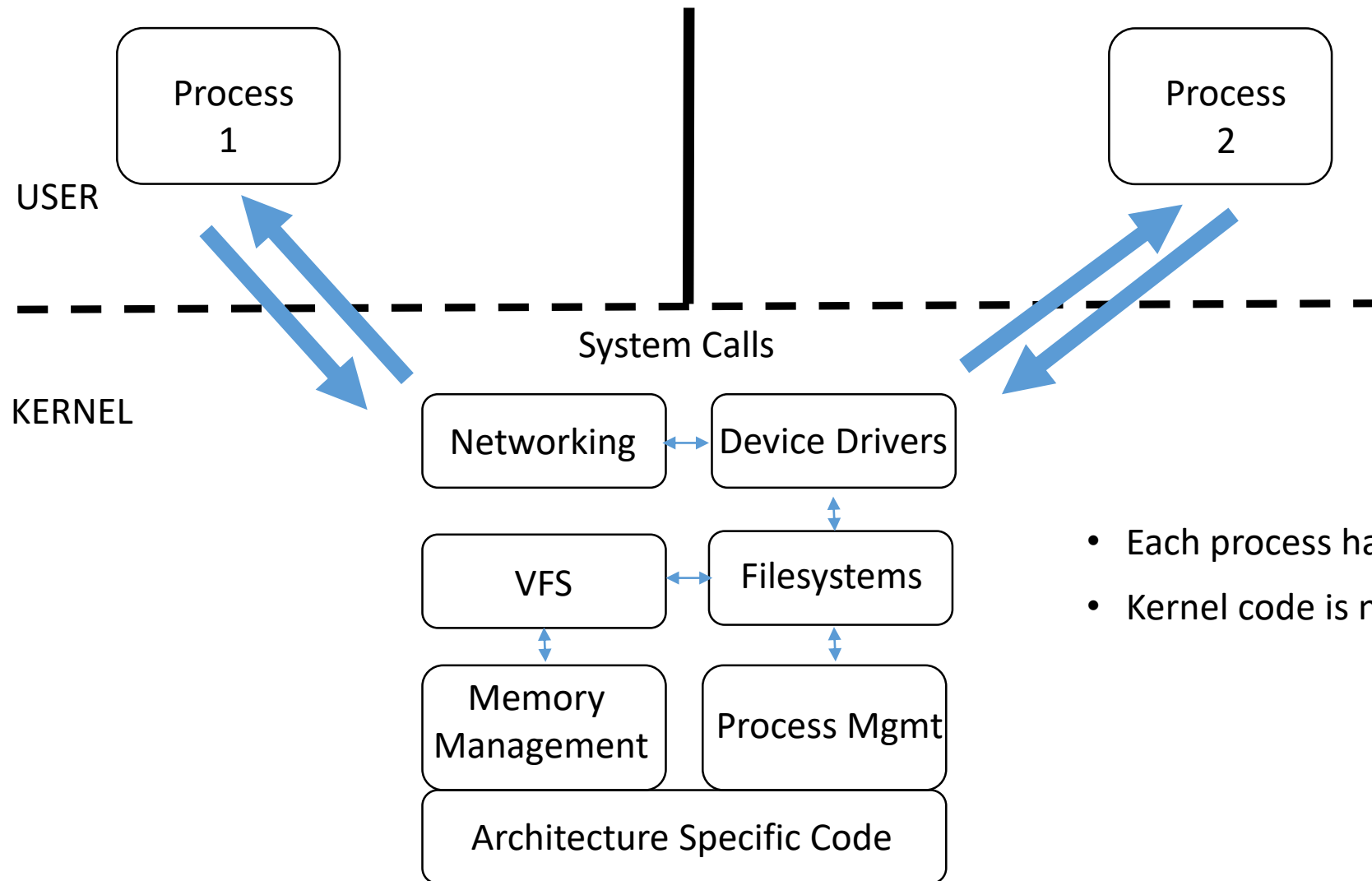
# Demo: Attack Scenario



## Part 4: Splitting the Kernel

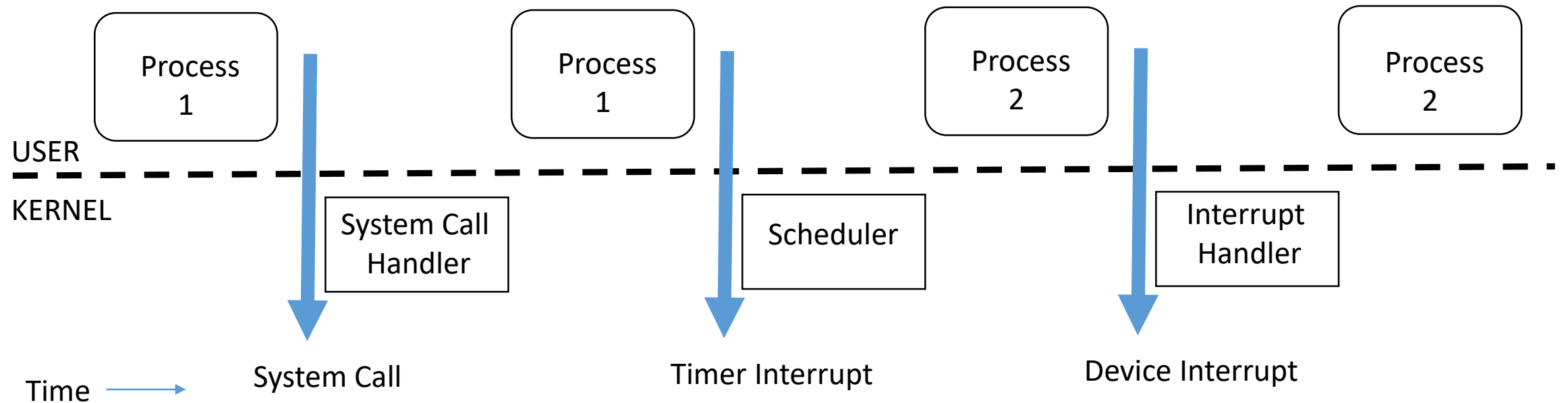
(Hacking a Micro-Kernel interface into Linux)

# Linux Kernel



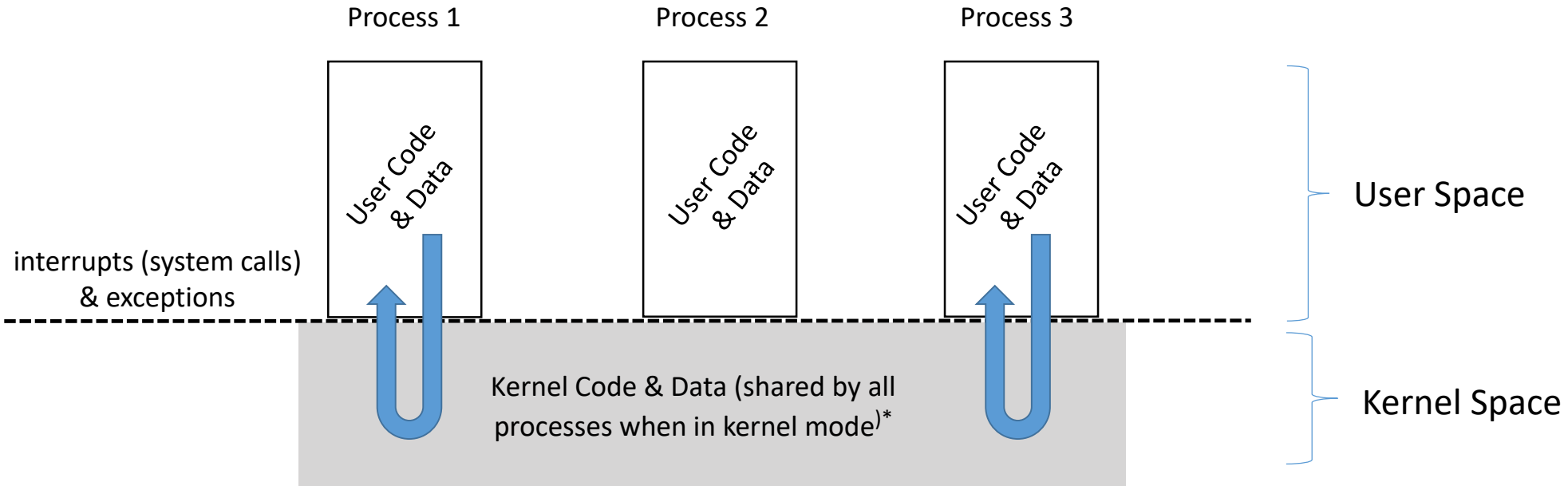
- Each process has its own isolated user space
- Kernel code is not isolated from itself

# Some More Linux Kernel Background



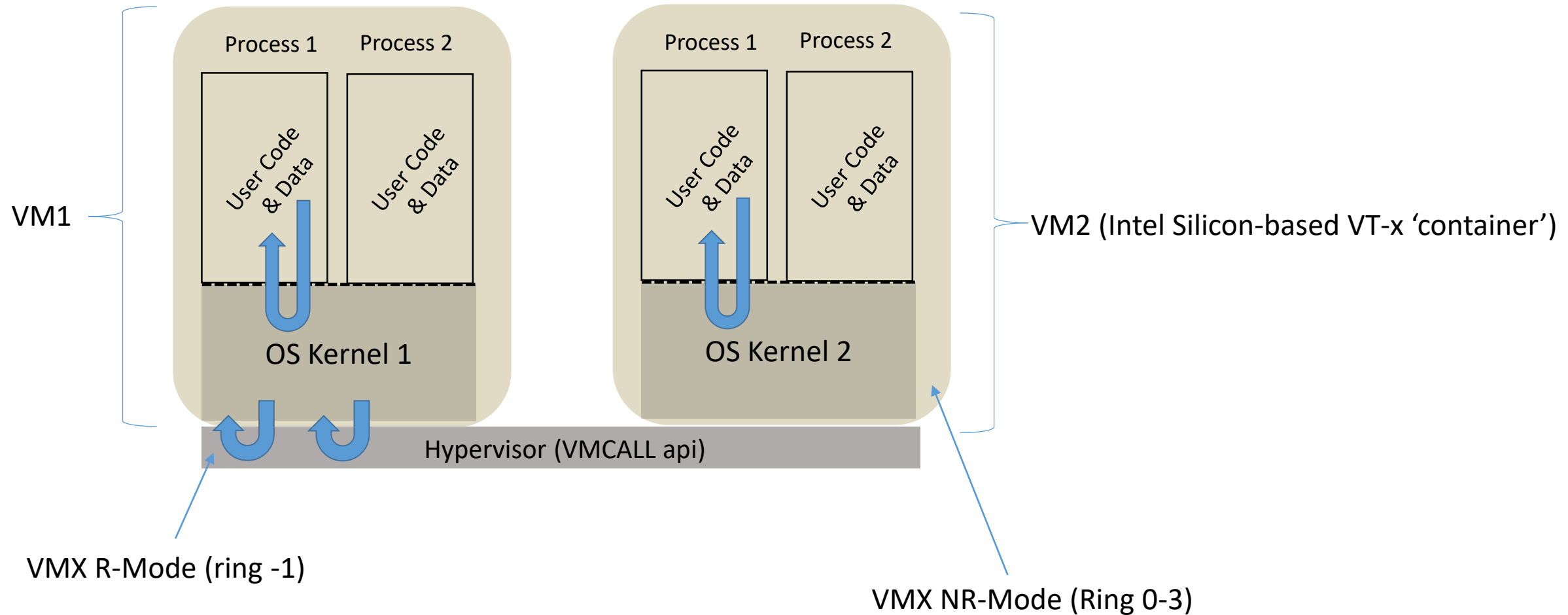
- Each process has its own user space Page table mappings
  - Ignoring process cloning/threads
- Shares kernel mapping (syncd from `init_mm.pgd`)
- Separate kernel stack per process
- Kernel entered through process context
  - It is not a separate 'thing' that runs concurrently
  - Does support the notion of kernel services via kthreads though

# Linux Kernel Design



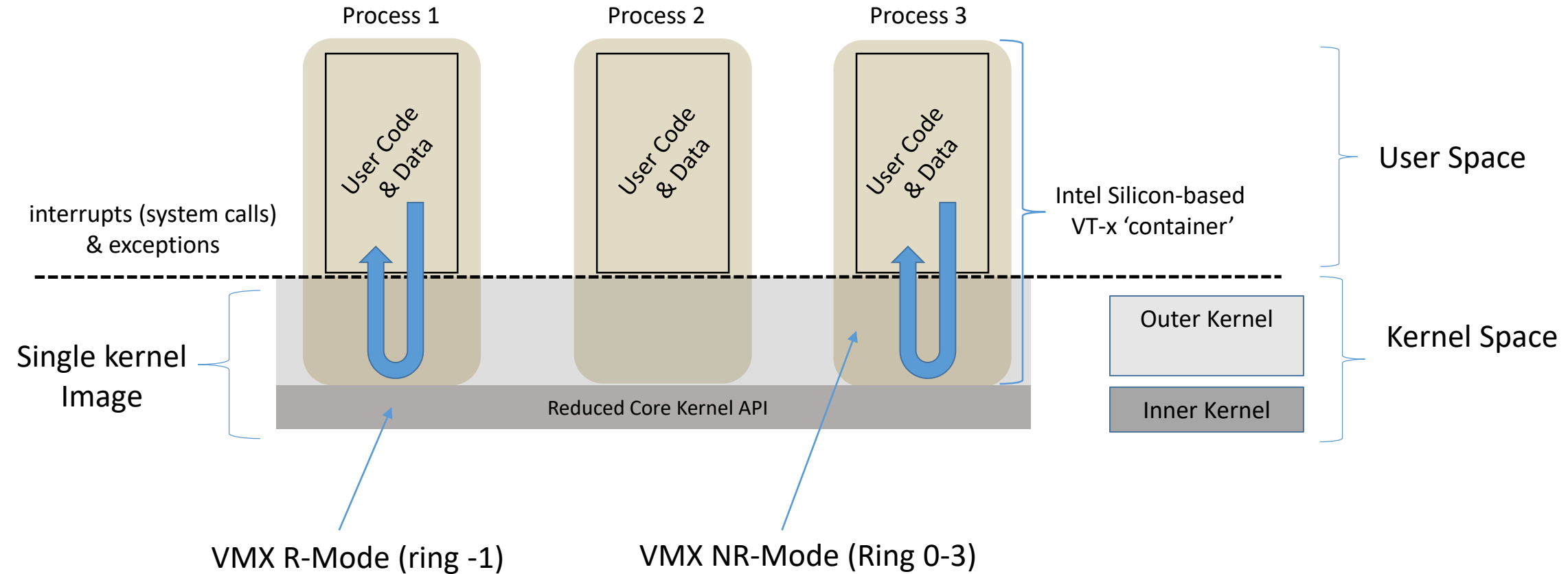
\*separate kernel stack per process

# Intel Virtualization HW support for VMMs



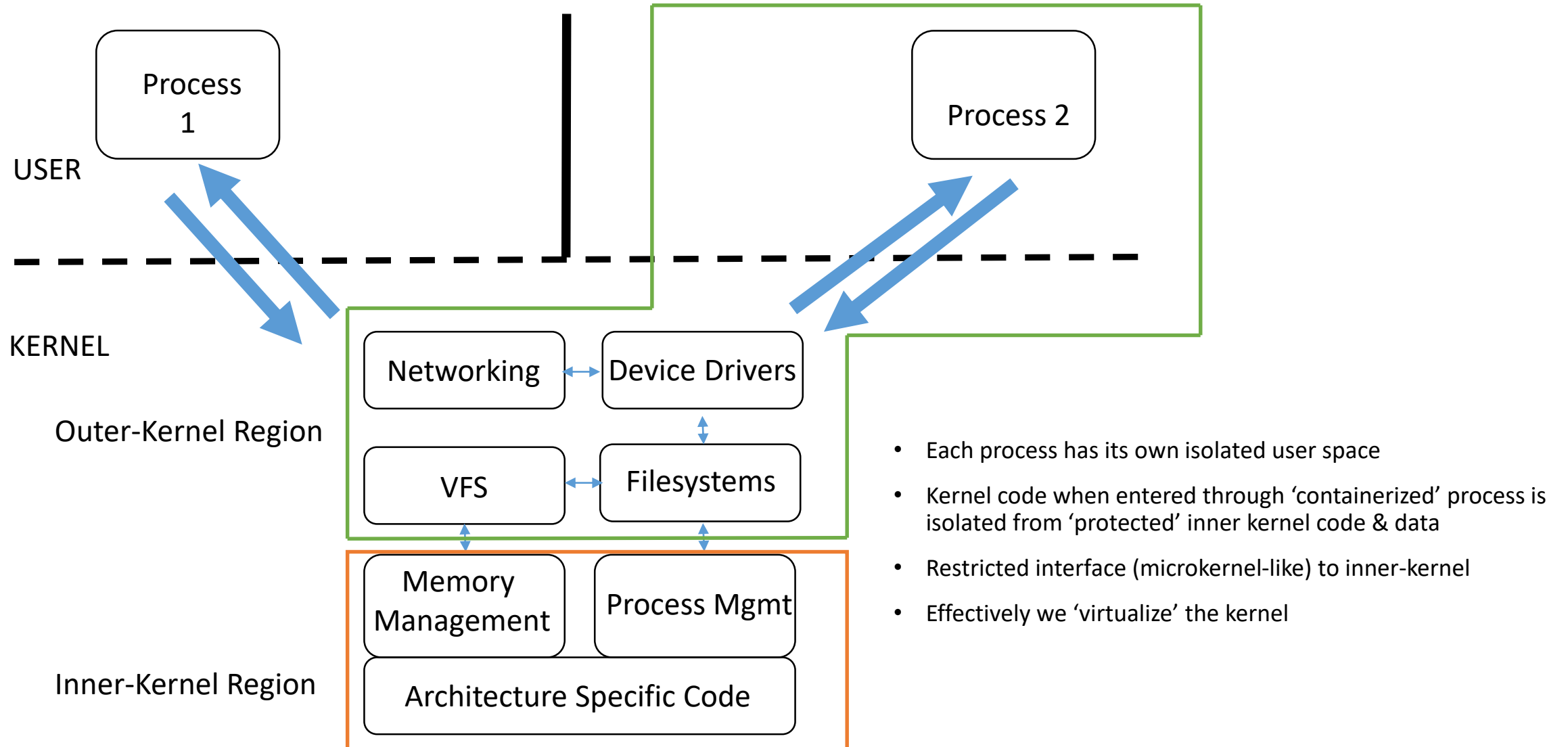
Separate Kernel Code & Data per VM

# Linux Split-Kernel Design using Virtualization HW

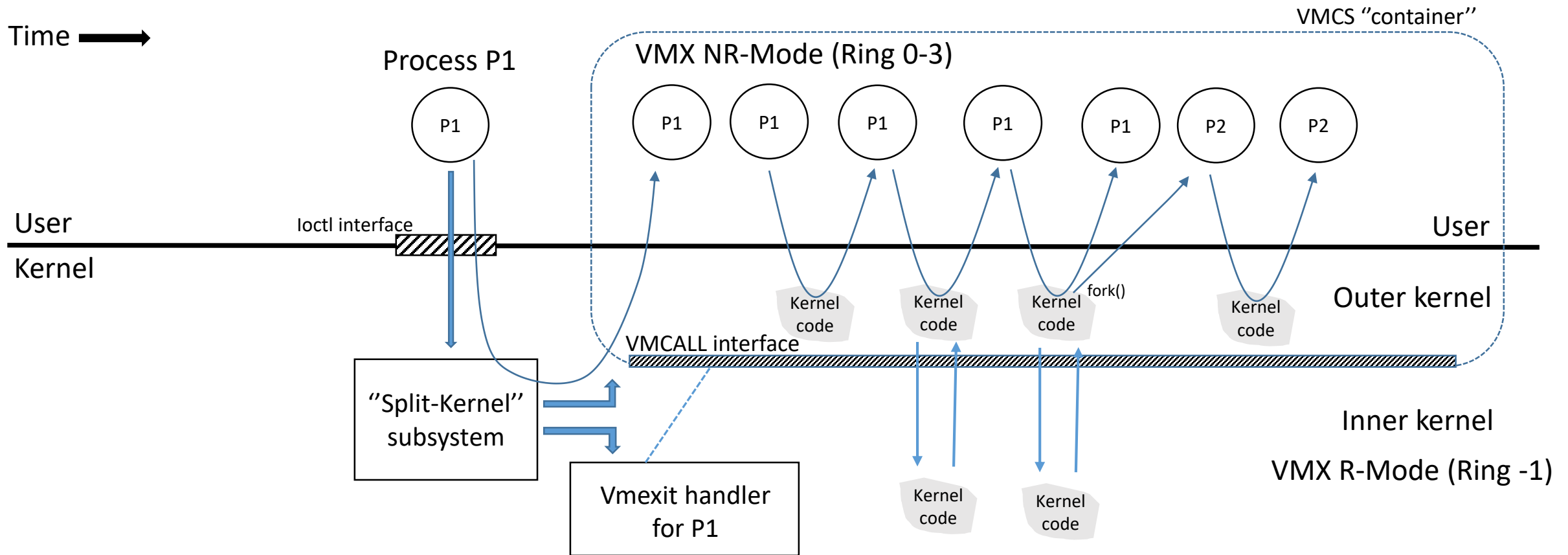




# Split-Kernel (logical view)



# Split-Kernel Process Lifecycle



# Split-Kernel Architecture (x86\_64)

- Use Intel VMX/EPT interposed on the normal linux kernel path
  - Allows vMMU interface for kernel entered through userspace-process to be enforced
  - Split the single shared kernel into an **Outer** region and an **Inner** region
  - Same kernel image but detects during execution whether in outer-kernel or inner-kernel mode
- X86: Each (container) process runs in vmx non-root mode within an EPT 'Container'
  - State defined by its own VMCS record
  - Each VMCS has its EPT pointer set to the same top-level table
  - Direct map 1:1 'host' physical memory except for 'protected' or 'private' memory regions
  - 'Protected' memory regions mapped RO, 'private' regions not mapped unless owned by that process
- do\_schedule() from outer-kernel does VMCALL into inner-kernel
  - R-mode side of the process is scheduled / de-scheduled
- Inner-kernel needs to provide a VMEXIT handler
  - Need to maintain VMCS state across time-slices

# Split-Kernel Architecture (x86\_64)

- Use Intel VMX/EPT interposed on the normal linux kernel path
  - Allows vMMU interface for kernel entered through userspace-process to be enforced
  - Split the single shared kernel into an **Outer** region and an **Inner** region
  - Same kernel image but detects during execution whether in outer-kernel or inner-kernel mode
- X86: Each (container) process runs in vmx non-root mode within an EPT 'Container'
  - State defined by its own VMCS record
  - Each VMCS has its EPT pointer set to the same top-level table
  - Direct map 1:1 'host' physical memory except for 'protected' or 'private' memory regions
  - 'Protected' memory regions mapped RO, 'private' regions not mapped unless owned by that process
- `do_schedule()` from outer-kernel does VMCALL into inner-kernel
  - R-mode side of the process is scheduled / de-scheduled
- Inner-kernel needs to provide a VMEXIT handler
  - Need to maintain VMCS state across time-slices

# Split-Kernel Architecture (x86\_64)

- Use Intel VMX/EPT interposed on the normal linux kernel path
  - Allows vMMU interface for kernel entered through userspace-process to be enforced
  - Split the single shared kernel into an **Outer** region and an **Inner** region
  - Same kernel image but detects during execution whether in outer-kernel or inner-kernel mode
- X86: Each (container) process runs in vmx non-root mode within an EPT 'Container'
  - State defined by its own VMCS record
  - Each VMCS has its EPT pointer set to the same top-level table
  - Direct map 1:1 'host' physical memory except for 'protected' or 'private' memory regions
  - 'Protected' memory regions mapped RO, 'private' regions not mapped unless owned by that process
- do\_schedule() from outer-kernel does VMCALL into inner-kernel
  - R-mode side of the process is scheduled / de-scheduled
- Inner-kernel needs to provide a VMEXIT handler
  - Need to maintain VMCS state across time-slices

# Split-Kernel Architecture (x86\_64)

- Use Intel VMX/EPT interposed on the normal linux kernel path
  - Allows vMMU interface for kernel entered through userspace-process to be enforced
  - Split the single shared kernel into an **Outer** region and an **Inner** region
  - Same kernel image but detects during execution whether in outer-kernel or inner-kernel mode
- X86: Each (container) process runs in vmx non-root mode within an EPT 'Container'
  - State defined by its own VMCS record
  - Each VMCS has its EPT pointer set to the same top-level table
  - Direct map 1:1 'host' physical memory except for 'protected' or 'private' memory regions
  - 'Protected' memory regions mapped RO, 'private' regions not mapped unless owned by that process
- do\_schedule() from outer-kernel does VMCALL into inner-kernel
  - R-mode side of the process is scheduled / de-scheduled
- Inner-kernel needs to provide a VMEXIT handler
  - Need to maintain VMCS state across time-slices

# What does this buy us again?

- Strong control over the integrity of kernel code and core data
  - No 'un-authorized' code gets to run in kernel mode
  - Can protect the integrity of data even against malicious kernel mode code
- Can offer confidentiality guarantees
  - Can protect application secrets even against malicious kernel mode code
- Guard against cross-process code + data compromise
  - Enhanced risk when running multiple 'isolated' containers

# What does this buy us again?

- Strong control over the integrity of kernel code and core data
  - No 'un-authorized' code gets to run in kernel mode
  - Can protect the integrity of data even against malicious kernel mode code
- Can offer confidentiality guarantees
  - Can protect application secrets even against malicious kernel mode code
- Guard against cross-process code + data compromise
  - Enhanced risk when running multiple 'isolated' containers



# What does this buy us again?

- Strong control over the integrity of kernel code and core data
  - No 'un-authorized' code gets to run in kernel mode
  - Can protect the integrity of data even against malicious kernel mode code
- Can offer confidentiality guarantees
  - Can protect application secrets even against malicious kernel mode code
- Guard against cross-process code + data compromise
  - Enhanced risk when running multiple 'isolated' containers

# Some Code: Entering Outer-Kernel Mode

```
09
10 long ok_device_ioctl(struct file *filp, unsigned int cmd, unsigned long arg)
11 {
12     switch(cmd)
13     {
14     case OKERNEL_ON_CMD:
15         /* do okernel_enter() here... */
16         HDEBUB("About to go into okernel on mode via okernel_enter() for pid (%d)\n", current->pid);
17
18         current->okernel_status = OKERNEL_ON;
19
20         if(is_in_vmx_nr_mode()){
21             printk(KERN_CRIT "NR: Calling okernel_enter in NR mode kernel...s\n");
22         }
23
24         ret = okernel_enter(OKERNEL_ENTER_IOCTL);
25
26         if(is_in_vmx_nr_mode()){
27             BXMAGICBREAK;
28             HDEBUB("Returning from okernel_enter (IOCTL_LAUNCH).\n");
29             BXMAGICBREAK;
30             goto nr_exit;
31         }
32         current->okernel_status = OKERNEL_OFF;
33
34         if(ret){
35             printk(KERN_ERR "okernel_enter failed for pid <%d> ret (%d)\n", current->pid, ret);
36         }
37         HDEBUB("outer kernel off for <%d>\n", current->pid);
38         unset_vmx_r_mode();
39         break;
40     default:
41         printk(KERN_ERR "okernel invalid IOCTL cmd.\n");
42         return -ENODEV;
43     }
44 nr_exit:
45     return 0;
46 }
47
```

```
HDEBUB("About to enter vmx handler while loop...\n");
while(1){
    HDEBUB("Entering vmxit handler loop...\n");

    /* Use this instead of local_irq_disable() to save getting tangled in irq on/off
    native_irq_disable();

    /****** GO FOR IT *****/
    ret = vmx_run_vcpu(vcpu);
    /****** GONE FOR IT! *****/

    if((vmcs_readl(GUEST_RFLAGS) & RFLAGS_IF_BIT)){
        native_irq_enable();
    }

    HDEBUB("Returned from vmx_run_vcpu...handling exit condition...\n");

    if(ret==VMX_EXIT_REASONS_FAILED_VMENTRY){
        HDEBUB("vmentry failed (%#x)\n", ret);
        break;
    } else if (ret == EXIT_REASON_EPT_VIOLATION){
        qual = vmcs_readl(EXIT_QUALIFICATION);
        gp_addr = vmcs_readl(GUEST_PHYSICAL_ADDRESS);

        HDEBUB("ept violation exit - qualification=%#lx gpa=%#lx\n",
            qual, gp_addr);
        /* Grant access to protected pages lazily */
        if(__ok_protected_phys_addr(gp_addr)){
            HDEBUB("ok: EPT vmxit on protected address:=(%#lx)\n", gp_addr);
            if(ok_allow_protected_access(gp_addr)){
                (void)add_ept_page_perms(vcpu, gp_addr);
            } else {
                /*
                Map in 'dummy' page for now - need to be careful not
                create another vulnerability.
                */
                (void)remap_ept_page(vcpu, gp_addr, ok_get_protected_dumm
            }
        }
        continue;
    }
    HDEBUB("Done handling exit condition, looping...\n");
}

/* something went wrong - legitimate exit would have been through the
```

# Some Code: Scheduling

```
87
88 asmlinkage __visible void __sched schedule(void)
89 {
90     struct task_struct *tsk = current;
91     #if defined(CONFIG_OKERNEL)
92     #ifdef HPE_DEBUG
93         volatile struct thread_info *ti;
94         int cpu = smp_processor_id();
95         struct tss_struct *tss = &per_cpu(cpu_tss, cpu);
96         unsigned long fs;
97         int orig_cpu = 0;
98         int new_cpu = 0;
99     #endif
100
101     if(is_in_vmx_nr_mode()){
102         /* Return control to the original process running in root-mode VMX */
103         /* shouldn't be holding locks at this point? */
104     #ifdef HPE_DEBUG
105         ti = current_thread_info();
106         rdmsrl(MSR_FS_BASE, fs);
107         HDEBUG("called (pid=%d)  cpu_cur_tos (%#lx) flgs(%#x) MSR_FS_BASE=%#lx\n",
108             current->pid, (unsigned long)tss->x86_tss.sp0, ti->flags, fs);
109         BXMAGICBREAK;
110         HDEBUG("in_atomic(): %d, irqs_disabled(): %d, pid: %d, name: %s\n",
111             in_atomic(), irqs_disabled(), current->pid, current->comm);
112         HDEBUG("preempt_count (%#x) rcu_preempt_depth (%#x)\n",
113             preempt_count(), rcu_preempt_depth());
114     #if defined(CONFIG_TRACE_IRQFLAGS) && defined(CONFIG_PROVE_LOCKING)
115         HDEBUG("current->h_irqs_en (%d)\n",
116             current->hardirqs_enabled);
117     #endif /* CONFIG_TRACE_IRQFLAGS */
118     #endif /* HPE_DEBUG */
119         sched_submit_work(tsk);
120         (void)vmcall(VMCALL_SCHED);
121     } else {
122         sched_submit_work(tsk);
123         do {
124             preempt_disable();
125             okernel_schedule(false);
126             sched_preempt_enable_no_resched();
127         } while (need_resched());
128     }

```

```
void vmx_handle_vmcall(struct vmx_vcpu *vcpu, int nr_irqs_enabled)
{
    struct vmx_vcpu *cpu_ptr;
    /* do_fork_fixup args */
    struct task_struct *p;

    cmd = vcpu->regs[VCPU_REGS_RAX];

    if(cmd == VMCALL_DO_EXEC_FIXUP_HOST){
        /* Next time we take a vmexit we will return using these page tables - should val
        h_cr3 = vcpu->regs[VCPU_REGS_RBX];
        HDEBUG("exec_fixup: Setting saved HOST CR3 to (%#lx) __pa (%#lx)\n",
            (unsigned long)h_cr3, __pa(h_cr3));
        vmx_get_cpu(vcpu);
        vmcs_writel(HOST_CR3, __pa(h_cr3));
        vmx_put_cpu(vcpu);
        BXMAGICBREAK;
        ret = 0;
    } else if (cmd == VMCALL_SCHED){
        unset_vmx_r_mode();
        /* This is the only place we should be swapping CPUs */
        /* There is redundancy here: don't need to do all this flushing */
        vpid_sync_vcpu_global();
        ept_sync_global();
        barrier();

        schedule_r_mode();

        vpid_sync_vcpu_global();
        ept_sync_global();

        /* We may come back here on a different CPU...*/
        set_vmx_r_mode();

        cpu = smp_processor_id();
        tss = &per_cpu(cpu_tss, cpu);

        HDEBUG("ret schedule_r (pid %d) cpu_cur_tos (%#lx) tss.sp0 (%#lx) flgs (%#x)\n",
            current->pid, current_top_of_stack(), (unsigned long)tss->x86_tss.sp0, r_t
        HDEBUG("ret from sched in_atomic(): %d, irqs_disabled(): %d, pid: %d, name: %s\n",
            in_atomic(), irqs_disabled(), current->pid, current->comm);
        HDEBUG("ret from preempt_count (%d) rcu_preempt_depth (%d)\n",
            preempt_count(), rcu_preempt_depth());
        BXMAGICBREAK;
    } else if(cmd == VMCALL_DOEXIT){
        code = (long)vcpu->regs[VCPU_REGS_RBX];
        HDEBUG("calling do_exit(%ld)...\n", code);
        vmx_destroy_vcpu(vcpu);

```

# Performance and Invasiveness

- Still surprised it works at all...
  - Largest test machine 2 physical CPUs with 24 cores & 256 GB memory
  - Focus on functionality not performance or upstream optimization
  - Docker a good testcase for use of obscure Linux kernel features BTW
  - Initial debugging really, really hard
  - Can't use `printk()`, etc but Bochs is really useful 😊
- Benchmarking
  - Linux kernel build, Apache Phoronix, etc.
  - Overhead around 2-5% depending upon number of processor cores
  - Depends on approach to handling process migration across CPU cores
- Invasiveness
  - Core code around 1000-2000 new lines in separate (static) kernel module
  - Plus maybe 100-200 lines of other Linux kernel modifications
  - Does still look like Linux!

# Performance and Invasiveness

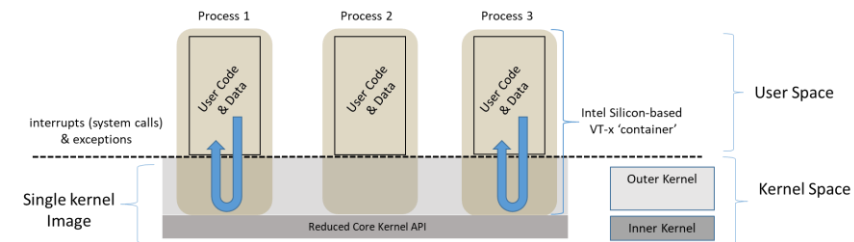
- Still surprised it works at all...
  - Largest test machine 2 physical CPUs with 24 cores & 256 GB memory
  - Focus on functionality not performance or upstream optimization
  - Docker a good testcase for use of obscure Linux kernel features BTW
  - Initial debugging really, really hard
  - Can't use `printk()`, etc but Bochs is really useful 😊
- Benchmarking
  - Linux kernel build, Apache Phoronix, etc.
  - Overhead around 2-5% depending upon number of processor cores
  - Depends on approach to handling process migration across CPU cores
- Invasiveness
  - Core code around 1000-2000 new lines in separate (static) kernel module
  - Plus maybe 100-200 lines of other Linux kernel modifications
  - Does still look like Linux!

# Performance and Invasiveness

- Still surprised it works at all...
  - Largest test machine 2 physical CPUs with 24 cores & 256 GB memory
  - Focus on functionality not performance or upstream optimization
  - Docker a good testcase for use of obscure Linux kernel features BTW
  - Initial debugging really, really hard
  - Can't use `printk()`, etc but Bochs is really useful 😊
- Benchmarking
  - Linux kernel build, Apache Phoronix, etc.
  - Overhead around 2-5% depending upon number of processor cores
  - Processor scaling depends on approach to handling process migration across CPU cores
- Invasiveness
  - Core code around 1000-2000 new lines in separate (static) kernel module
  - Plus maybe 100-200 lines of other Linux kernel modifications
  - Does still look like Linux!

# Current Status / Opportunities / Futures

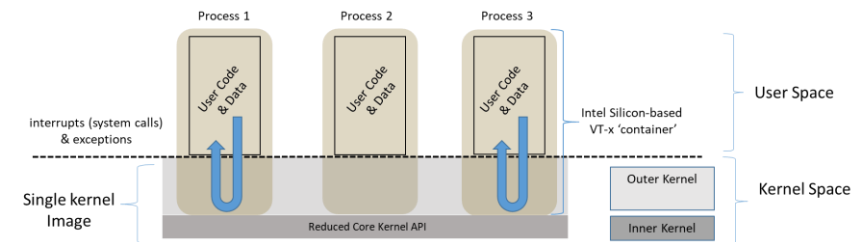
- Available at <https://github.com/linux-okernel>



- We do mention some of the concepts in our HOTOS paper
  - 'Separating Translation from Protection in Address Spaces with Dynamic Remapping', HOTOS 17
- Worth considering for general Linux use
  - Not just container deployments
  - Can run light-dm in outer-kernel mode (e.g. full desktop)
  - Need to fill-out ARM v8.1 implementation
- Good vehicle for kernel malware tracing on top of enhanced security
  - What else?

# Current Status / Opportunities / Futures

- Available at <https://github.com/linux-okernel>

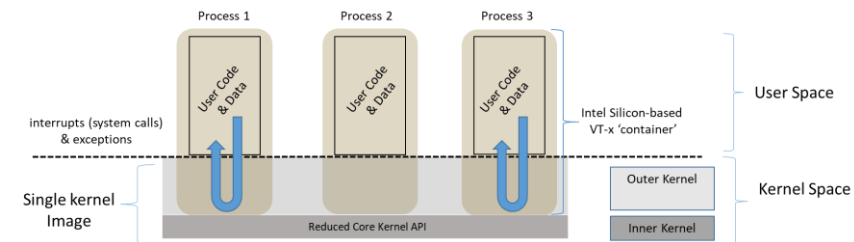


- We do mention some of the concepts in our HOTOS paper
  - 'Separating Translation from Protection in Address Spaces with Dynamic Remapping', HOTOS 17
- Worth considering for general Linux use
  - Not just container deployments
  - Can run light-dm in outer-kernel mode (e.g. full desktop)
  - Need to fill-out ARM v8.1 implementation
- Good vehicle for kernel malware tracing on top of enhanced security
  - What else?



# Current Status / Opportunities / Futures

- Available at <https://github.com/linux-okernel>



- We do mention some of the concepts in our HOTOS paper
  - 'Separating Translation from Protection in Address Spaces with Dynamic Remapping', HOTOS 17
- Worth considering for general Linux use
  - Not just container deployments
  - Can run light-dm in outer-kernel mode (e.g. full desktop)
  - Need to fill-out ARM v8.1 implementation
- Good vehicle for kernel malware tracing on top of enhanced security
  - What else?