

Insomni'Hack 2019

Building a
{Flexible}
Hypervisor-Level
Debugger

Mathieu Tarral

Whoami



- ex F-Secure fellow
- Building stealth, hypervisor-based sandboxes
- Virtual Machine Introspection

https://github.com/KVM-VMI/kvm-vmi





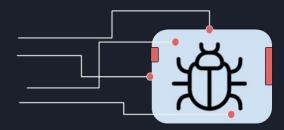
Agenda

- Why?
- History
- Challenges
- Demo
- HyperBreakpoints
- Future

Why?

Debuggers observer effect

- Advanced malware analysis
- No stealth
 - changes in exposed structures
 - altered syscall behavior
 - visible breakpoints
- No robustness
 - how to protect the debugger against code modifications?
 - o runs at the same **privilege** level as the malware (malicious driver)

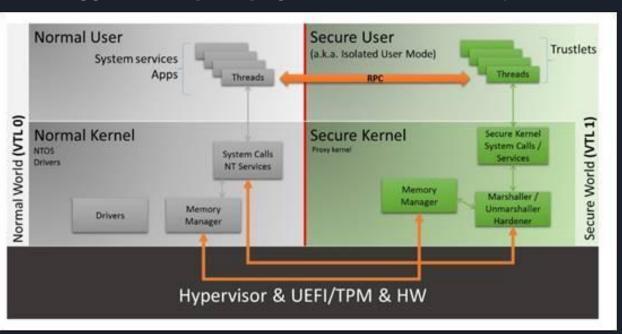


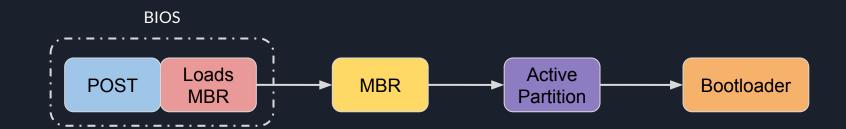
Protected OS features

- The observer effect might sometimes be intentional
 - o protect intellectual property from reverse-engineering
 - Protected Media Path (used to enforce DRM)
- Modern OS security mechanisms are interfering
 - PatchGuard
 - otherwise BSOD!

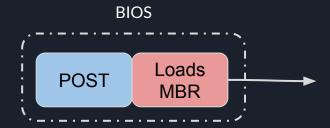
Incomplete system view

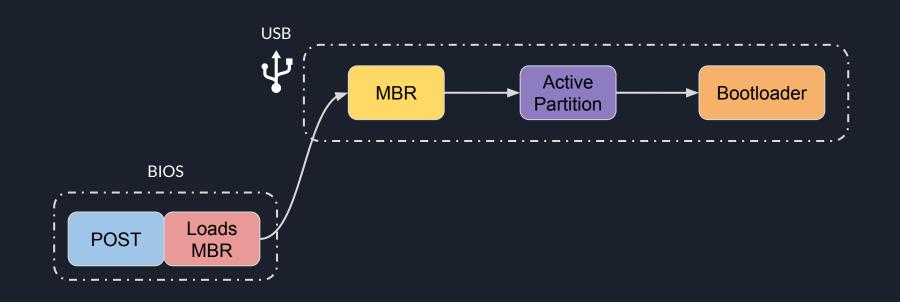
Debuggers are fighting against new OS security features

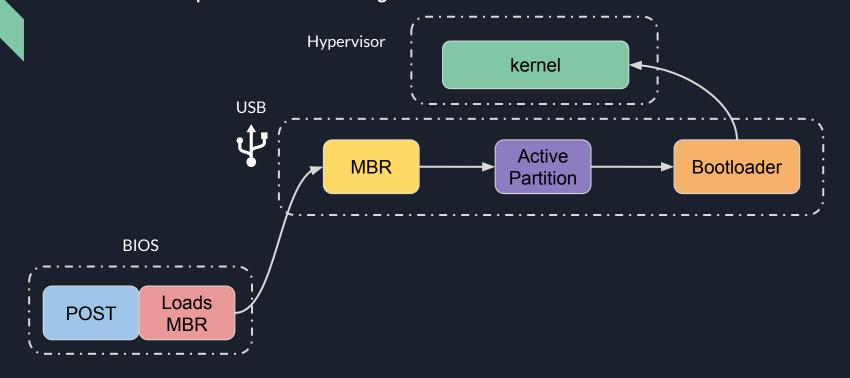


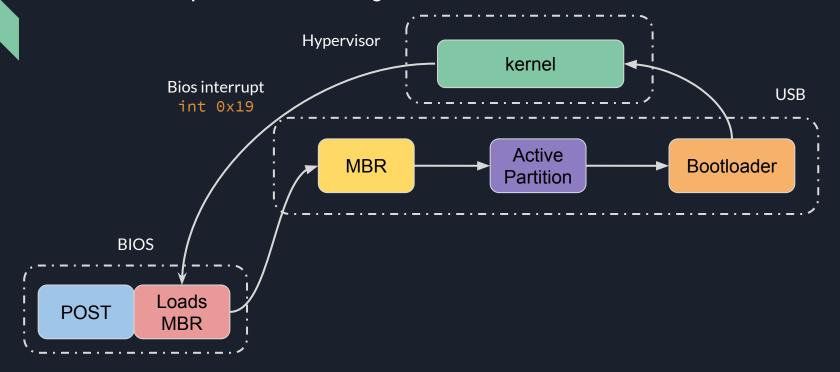


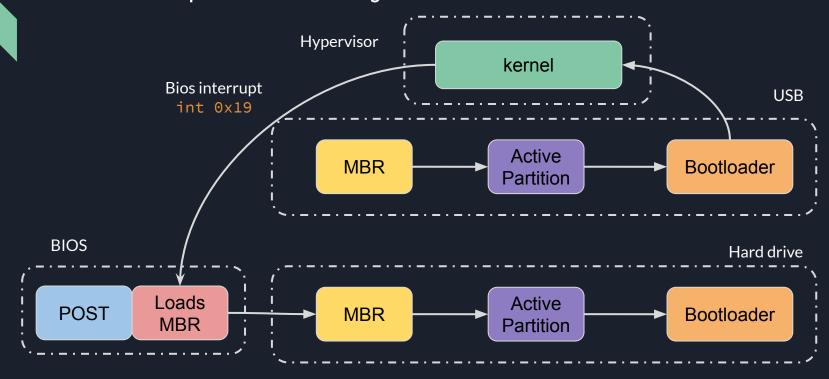


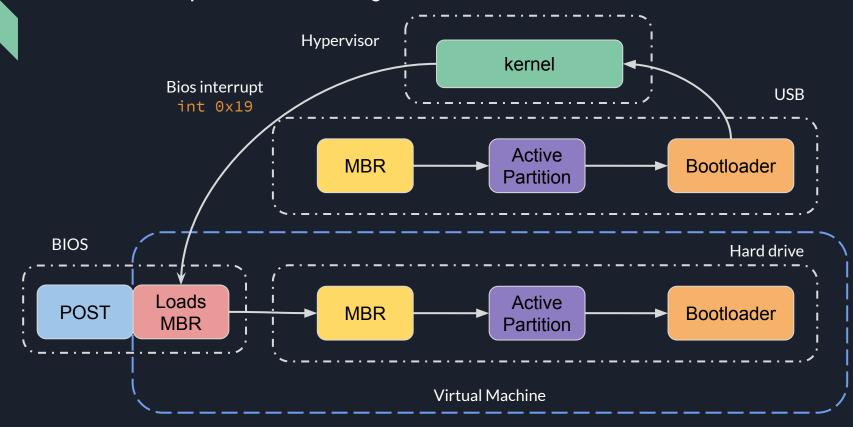






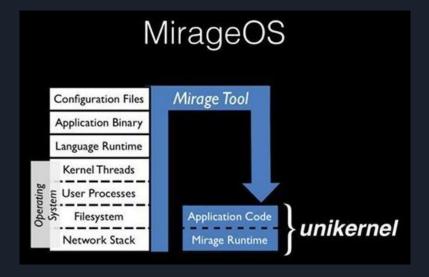






Unikernel debugging

- Specialized embedded system images
- Kernel/app runs in ring0
- One process, one address space
- No debug stub...



Hacker News new | past | comments | ask | show | jobs | submit

cdoxsey 9 months ago | parent | favorite | on: A Rust-Based Unikernel: First Version of a Rust-Ba...

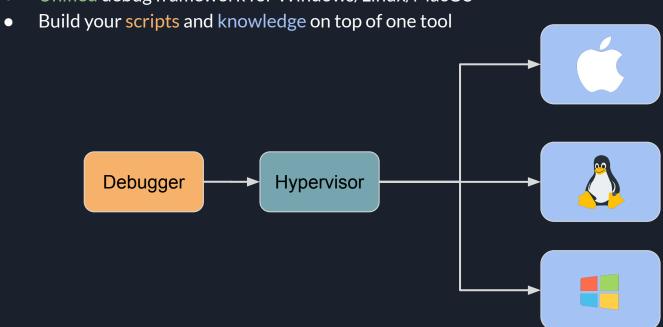
Unikernels can be difficult to debug. All the tools that come with your operating system are missing.

Unmodified guests

- No remote debug agent/stub
- No custom VM setup
 - hardware
 - network card
 - serial cable
 - software
 - install Windows SDK
 - configuration
 - bcdedit /set debug on
 - bcdedit /dbgsettings serial debugport:1 baudrate:115200
- On-the-fly debugging

Cross-platform debugging

Unified debug framework for Windows/Linux/MacOS



Follow malware sandbox trend

- Nowadays, malware sandboxing solutions are hypervisor-based
 - VMRay, Joe Security, DRAKVUF
- Cross-platform
- Agentless
- Stealth breakpoints
- Process hijacking and code injection
 - o start guest binary execution via hypervisor

Why: Recap

- 1. Advanced malware analysis
- 2. OS debug API untrusted
- 3. Boot sequence and unikernel analysis
- 4. Untouched guest VMs
- 5. Unified tool for OS debugging

-> Leverage the hypervisor as a new debugging platform

History

Timeline (2003 - 2019)

- 2003: QEMU GDB stub
- 2007: VMware GDB stub
- 2010: HyperDBG
- 2011:
 - virtdbg
 - o ramooflax
- 2012: VirtualBox GDB stub
- 2014: vmidbg
- 2016: Winbagility

- 2017:
 - PulseDBG
 - PyREBox
 - o rVMI
- 2018:
 - Sandbagility
 - o r2vmi
- 2019:
 - pyvmidbg
 - XenDBG

Categories

- Built-in hypervisor debug stubs
- Bare metal debuggers
 - hyperjacking
 - USB-boot
- Virtual Machine (virtual hardware) debuggers
 - emulated
 - full-virtualization

Built-in hypervisor debug stubs

- QEMU/VirtualBox/VMware GDB stub
- no real stealth
 - VMware: debugStub.hideBreakpoints = "TRUE"
 - Uses debug registers (4 breakpoints...)
- no guest awareness
 - o debug the kernel, that's it
- no flexibility
 - o GDB-only
 - one hypervisor

Bare-metal debuggers (hyperjacking)

- HyperDBG (2010)
 - "I want to take full control of a production system"
 - driver is installed on the host (Windows 7)
- Virtdbg (2011)
 - "I want to debug PatchGuard"
 - o driver is injected via DMA attack

- driver development
- not OS-agnostic

Bare-metal debuggers (USB boot)

- Ramooflax (2011)
 - "Debugging modern operating systems and real BIOS on physical machine"
 - python remote control API
- PulseDBG (2017)
 - "I want a better WinDBG UI"
 - hypervisor is contained in an EFI bootloader (bootx64.efi)
 - SDK to interact with stub

Virtual machine debuggers (emulation)

- PyREBox CISCO Talos (2017)
 - "I want a scriptable dynamic instrumentation system"
 - Instrumentation of QEMU (emulator)
 - Volatility for semantic gap
 - IPython shell
 - Fine grained callbacks

Virtual machine debuggers (full-virt)

- Unmaintained ()
 - vmidbg (2014)
 - GDB stub on top of LibVMI
 - o rVMI FireEye (2017)
 - KVM instrumentation
 - Rekall as introspection layer / debugger interface

Virtual machine debuggers (full-virt)

- Unmaintained ()
 - o vmidbg (2014)
 - GDB stub on top of LibVMI
 - o rVMI FireEye (2017)
 - KVM instrumentation
 - Rekall as introspection layer / debugger interface
- Winbagility (2016)
 - VirtualBox instrumentation
 - KD stub (WinDBG)
- Sandbagility (2018)
 - malware analysis framework
 - based on Winbagility

Virtual machine debuggers (full-virt)

- Unmaintained ()
 - vmidbg (2014)
 - GDB stub on top of LibVMI
 - o rVMI FireEye (2017)
 - KVM instrumentation
 - Rekall as introspection layer / debugger interface
- Winbagility (2016)
 - VirtualBox instrumentation
 - KD stub (WinDBG)
- Sandbagility (2018)
 - o malware analysis framework
 - based on Winbagility
- Xendbg NCC Group (2019)
 - o GDB stub on top of Xen
 - debug unikernels

Objectives

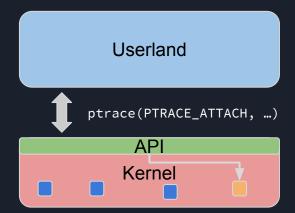
- Emulation, bare-metal: not interested
- I want a {flexible} hypervisor-level debugger to debug VMs
 - o r2vmi (2018)
 - based on LibVMI to be hypervisor-agnostic
 - radare2 for introspection/analysis and shell interactivity
 - pyvmidbg (2019)
 - GDB stub → LibVMI → hypervisor
 - compatibility over pure performance
 - guest aware debug stub

Challenges

Rebuilding the Debugger API

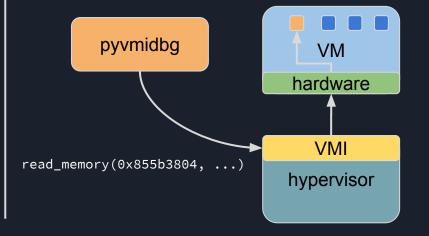
gdb attach <PID>

- Kernel knows how to query process states
- Interface is exposed via defined API



pyvmidbg attach <VM> <PID>

- We don't know the operating system
- Read and interpret raw memory



Filling the Semantic Gap

- We need an intimate knowledge of the guest
- Handle various debug formats (PDBs, DWARFs, etc...)
- Make sense of raw memory
 - forensic tools
- Rekall profiles to the rescue!
 - o JSON format
 - \$CONSTANTS
 - \$ENUMS
 - \$FUNCTIONS
 - \$STRUCTS
- Offsets from kernel base address

```
"$CONSTANTS": {
 "BiosBegin": 2012200,
"str:A driver has added a device obje": 1817744
"SENUMS": {
 "BUS_QUERY_ID_TYPE": {
 "0": "BusQueryDeviceID",
 "1": "BusQueryHardwareIDs",
 "2": "BusQueryCompatibleIDs"
"$FUNCTIONS": {
 "CmDeleteKey": 1434156,
 "IoGetAttachedDevice": 99966,
 "NtCreateProcess": 1024640
},
"$STRUCTS": {
"LIST_ENTRY32": [8, {
 "Blink": [4, ["unsigned long", {}]],
 "Flink": [0, ["unsigned long", {}]]
}],
 "_EPROCESS": [608, {
 "ActiveProcessLinks": [136, ["_LIST_ENTRY", {}]],
 "ActiveThreads": [416, ["unsigned long", {}]],
 "AddressCreationLock": [240, [" FAST MUTEX", {}]]
```

Following OS changes

- pure semantic layer is not enough
- the debugger needs a bit of logic to adapt to operating system



changes

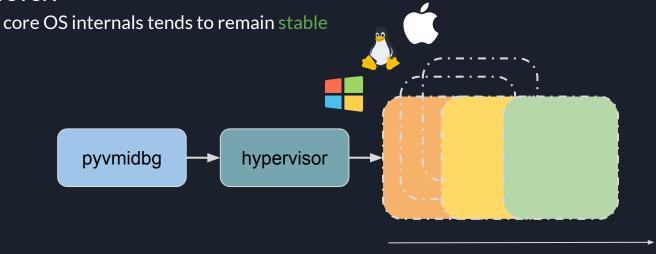
Following OS changes

- pure semantic layer is not enough
- the debugger needs a bit of logic to adapt to operating system
- eg: Windows
 - WinXP: IDLE process outside of main process list
 - Windows Vista: big changes in the kernel
 - Windows 10: VSM, Secure kernel



Following OS changes

- Problem is extended for every OS that you want to support
- Cost of building a hypervisor-based {flexible} debugger
- However:



Hypervisor-agnostic: LibVMI

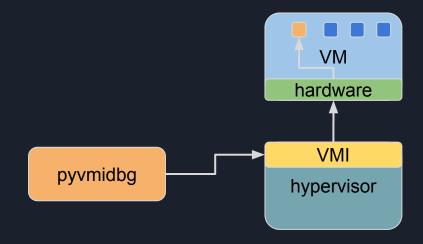
- VMI Abstraction layer
- Offers basic introspection
 - rekall profiles support
- Standard for VMI applications

```
winxp {
    ostype = "Windows";
    rekall_profile = "/etc/libvmi/winxp-profile.json";
}
```

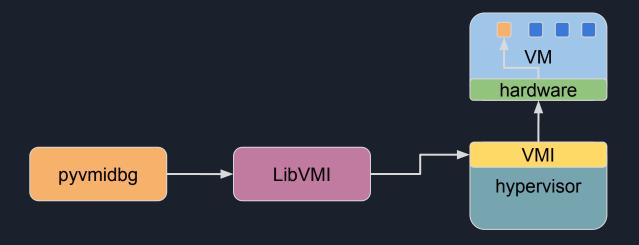
	VCPU Registers	Physical memory	Hardware events
Xen	✓	✓	V
KVM	✓	V	×



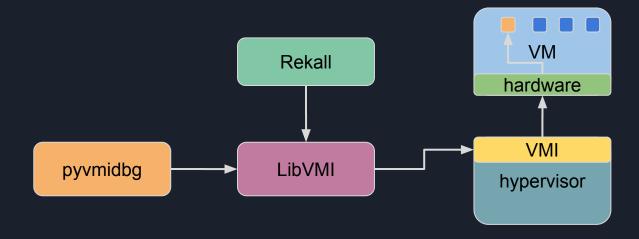
1. Virtual Machine Introspection



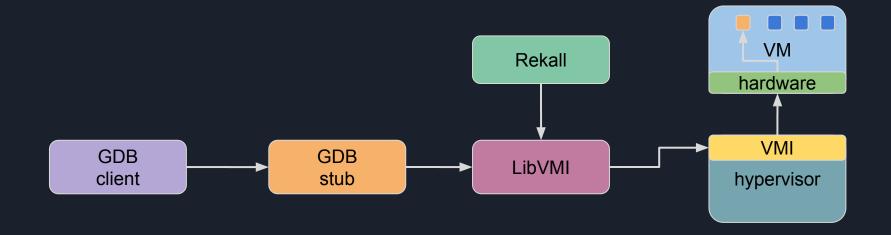
- 1. Virtual Machine Introspection
- 2. LibVMI



- 1. Virtual Machine Introspection
- 2. LibVMI
- 3. Rekall profiles

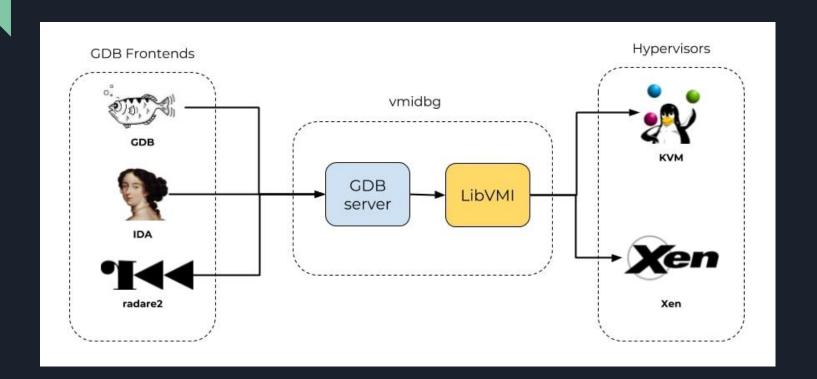


- 1. Virtual Machine Introspection
- 2. LibVMI
- 3. Rekall profiles
- 4. GDB protocol



Demo

pyvmidbg



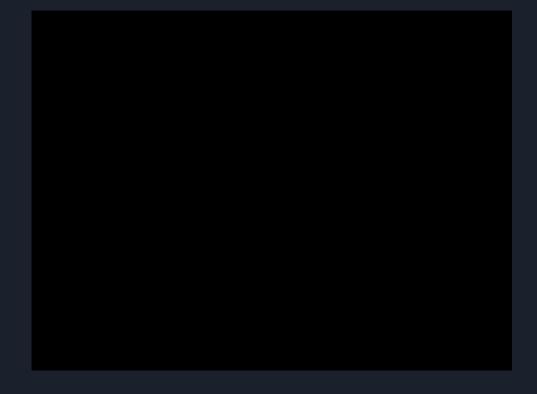
pyvmidbg

- pyvmidbg.py <port> <vm_name> [[cess>]
- 2 modes
 - o raw
 - 1 process -> system
 - thread -> VCPU
 - attach -> pause VM
 - basically your QEMU GDB stub
 - o guest-aware
 - Windows/Linux debug context
 - locate the targeted process descriptor
 - enumerate threads and read execution context
- basic software breakpoints
- singlestep
- breakin

pyvmidbg: raw



pyvmidbg: guest-aware



Demo



Improvements

- info proc
 - o process 1854
 - o cmdline = 'C:\WINDOWS\system32\cmd.exe /C dir'
 - cwd = 'C:\Document and Settings\Vagrant'
 - o exe = 'C:\WINDOWS\system32\cmd.exe'
- info sharedlibrary
 - 0x7c800000 kernel32
 - 0x7c900000 ntdll
 - o 0x7c9c0000 SHELL32
- callstacks

pyvmidbg: dev-friendly

- Strong dev requirements
 - running (modified) hypervisor
 - virtual machines
- Vagrant environment
 - Xen (packaged/from source)
 - o libvmi
 - libvirt
 - o vms
 - ubuntu 16.04/winxp
 - rekall profiles
- Nested virtualization
 - o KVM > Xen





HyperBreakpoints

HyperBreakpoints[®]

- "HyperBreakpoints" Winbagility article (SSTIC 2016)
- Breakpoints using virtualization layer capabilities
- We are not using the best breakpoints available in our debuggers

Basic SoftHyperBreakpoint

- write 0xCC
- listen on int3 events
 - o reinject in guest (if needed)
 - o restore opcode
 - singlestep
 - restore breakpoint



- ✓ fast
- X stealth
- x safe

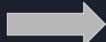
Avoid multi-VCPU race conditions?

--> pause all VCPUs at each event

- V fast-
- X stealth
- **V** safe

HardHyperBreakpoint |

- breakpoints in DRO-DR3
- listen on int1 events
 - o reinject in guest (if needed)
 - disable breakpoint in DR7
 - singlestep
 - enable breakpoint in DR7
- ✓ fastest
- X stealth
- V safe



4 breakpoints limit...

To gain stealth:

- listen on MOV-TO-DR events
 - restore guest's DRx
 - singlestep
 - o restore our DRx

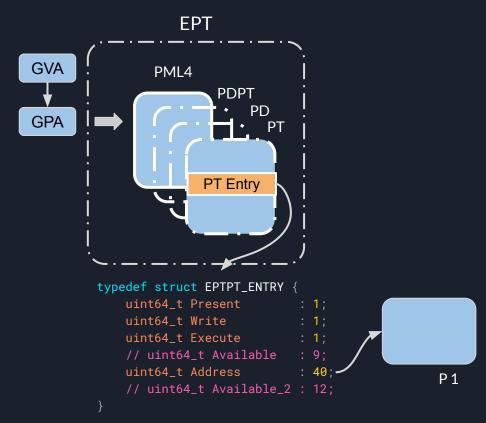
- **V** fastest
- **v** stealth
- ✓ safe

PageHyperBreakpoint

- GVA → GPA → HPA
- Extended Page Tables
- From Guest Virtual Address (GVA)
 - find Guest Physical Page (GPP)
 - translate via EPT
 - find Host Physical Page
 - change permissions
 - o generate #EPTViolation
- fine grained breakpoints (R/W/X)
- X fast
- 🗸 stealth
- X safe

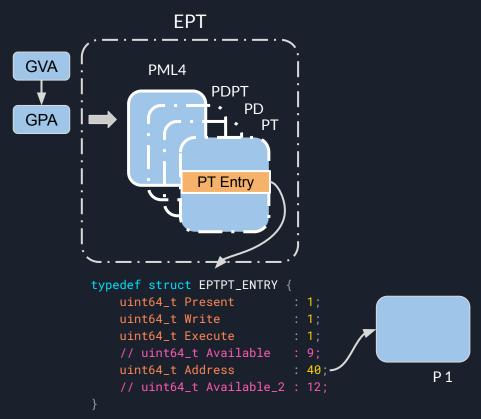
singlestep or emulate?

- emulators are incomplete
- emulation opens vulnerabilities
- heavy emulation can lead to instabilities



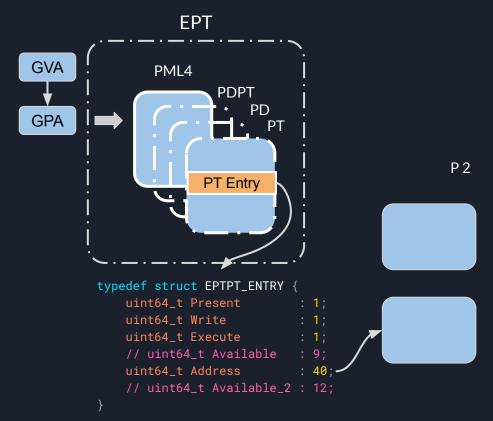
Setup:

1. duplicate Host Physical Page

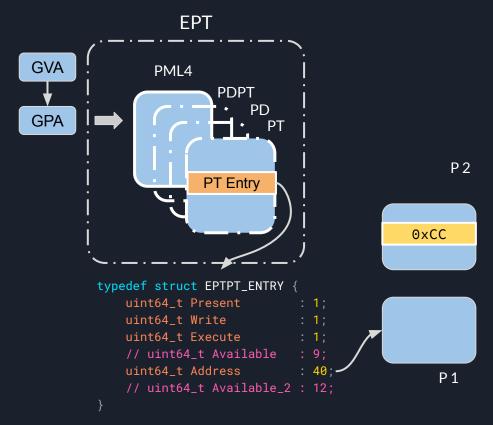


Setup:

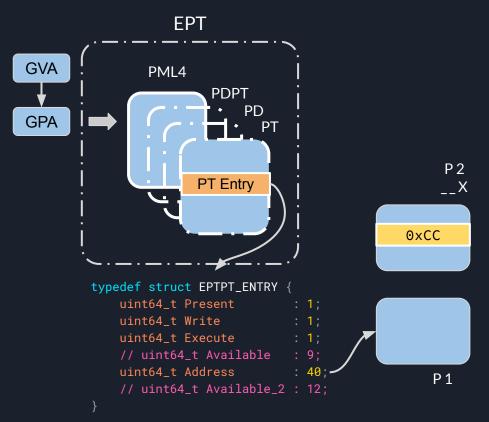
1. duplicate Host Physical Page



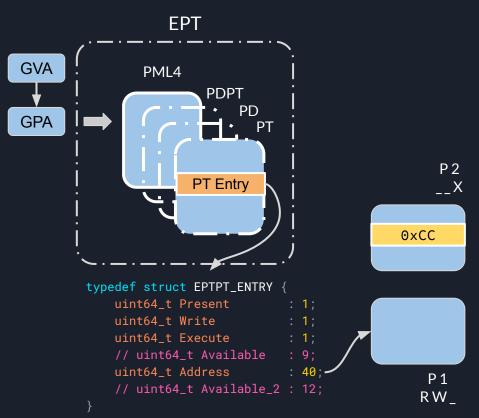
- 1. duplicate Host Physical Page
- 2. insert breakpoint in P2



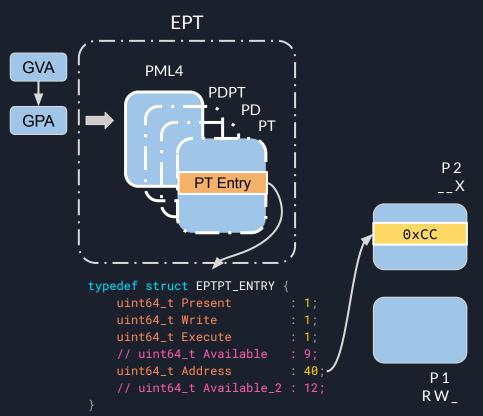
- 1. duplicate Host Physical Page
- 2. insert breakpoint in P2
- 3. set P2 as execute-only



- 1. duplicate Host Physical Page
- 2. insert breakpoint in P2
- 3. set P2 as execute-only
- 4. set P1 as read/write



- 1. duplicate Host Physical Page
- 2. insert breakpoint in P2
- 3. set P2 as execute-only
- 4. set P1 as read/write
- 5. point GPA to P2

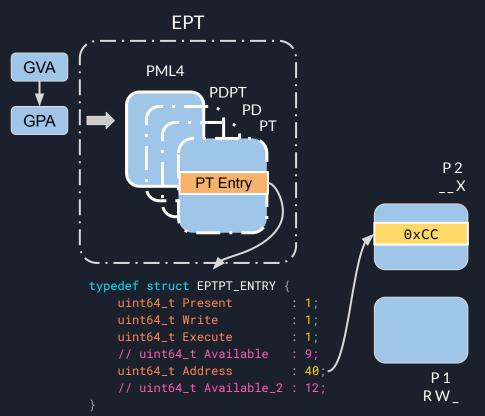


Setup:

- 1. duplicate Host Physical Page
- 2. insert breakpoint in P2
- 3. set P2 as execute-only
- 4. set P1 as read/write
- 5. point GPA to P2

#INT3

handle breakpoint



Setup:

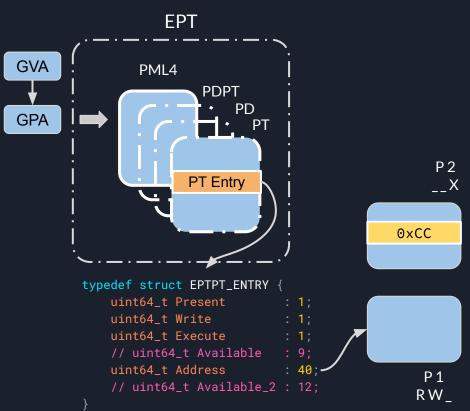
- 1. duplicate Host Physical Page
- 2. insert breakpoint in P2
- 3. set P2 as execute-only
- 4. set P1 as read/write
- 5. point GPA to P2

#INT3

• handle breakpoint

#EPTViolation (Read/Write)

• switch on P1



Setup:

- 1. duplicate Host Physical Page
- 2. insert breakpoint in P2
- 3. set P2 as execute-only
- 4. set P1 as read/write
- 5. point GPA to P2

#INT3

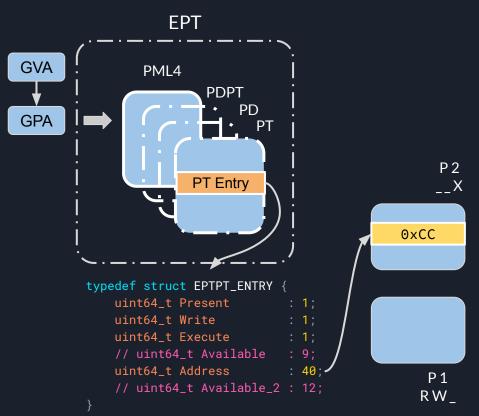
handle breakpoint

#EPTViolation (Read/Write)

• switch on P1

#EPTViolation (Execute)

• switch on P2



#INT3

handle breakpoint

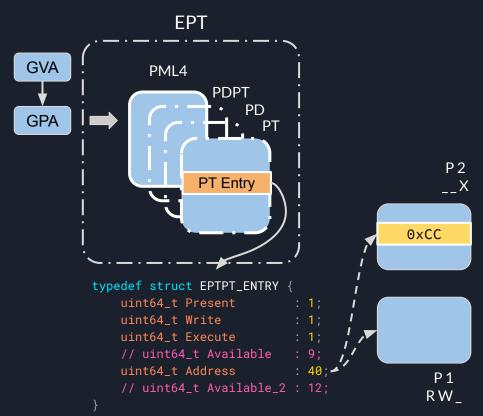
#EPTViolation (Read/Write)

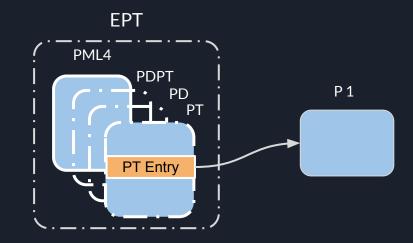
• switch on P1

#EPTViolation (Execute)

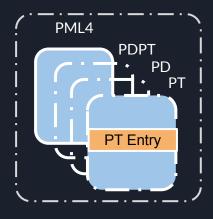
• switch on P2

- fast
- **V** stealth
- X safe

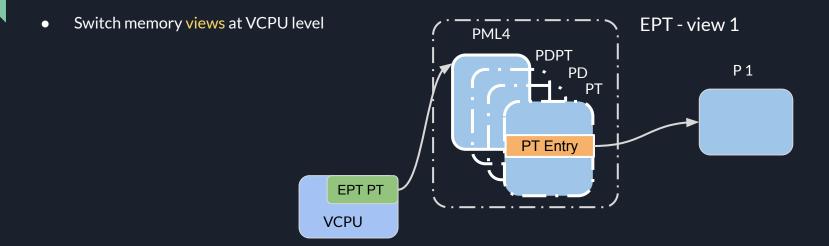


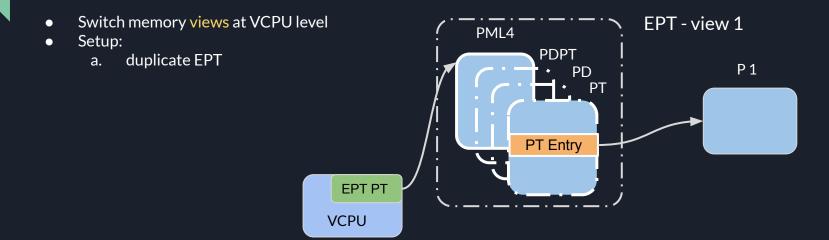


• Switch memory views at VCPU level

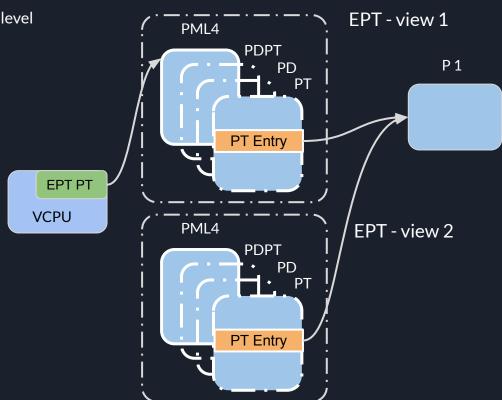


EPT - view 1

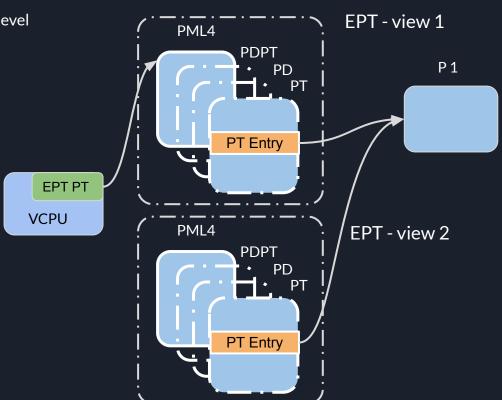




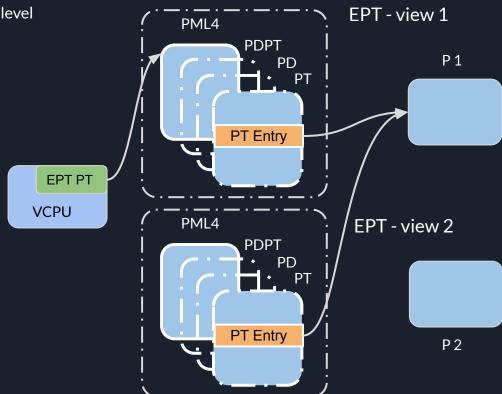
- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT



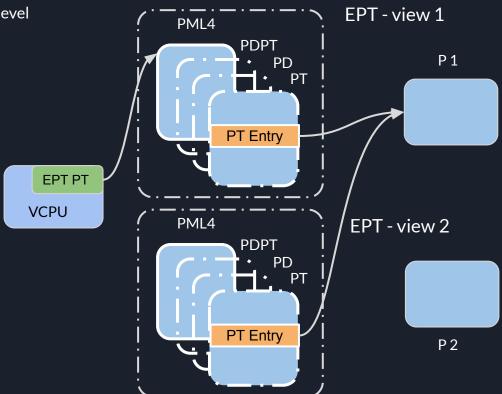
- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1



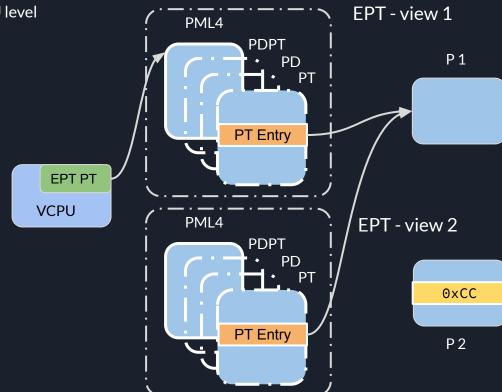
- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1



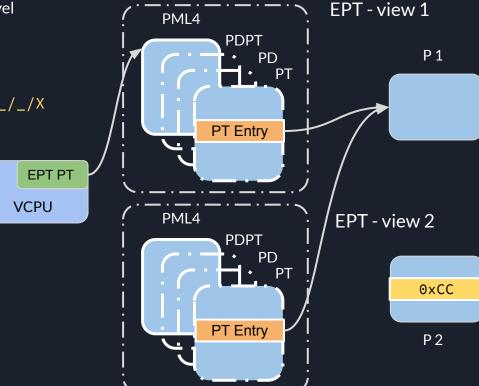
- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2



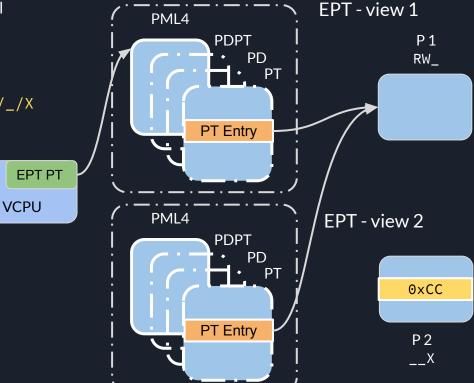
- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2



- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{/}/X$

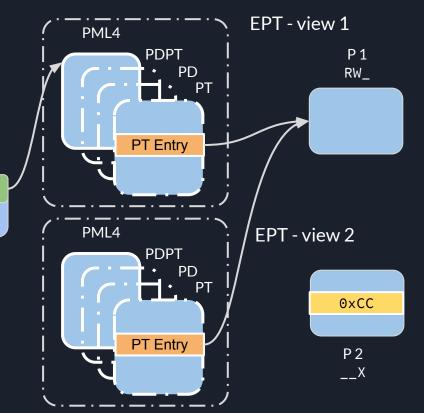


- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{-}/_{X}$



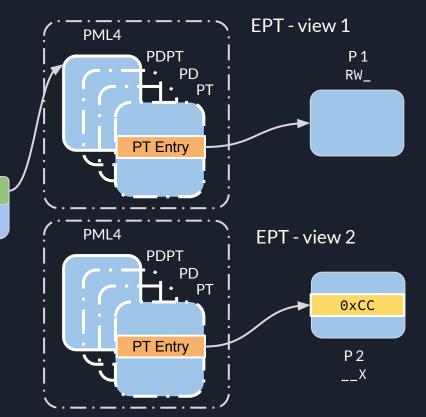
EPT PT

- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{/}/X$
 - e. point view 2 PT Entry to P2



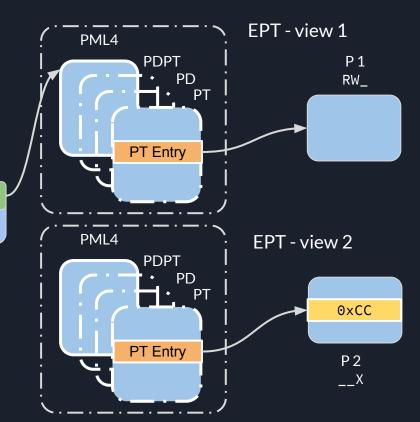
EPT PT

- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{/}/X$
 - e. point view 2 PT Entry to P2



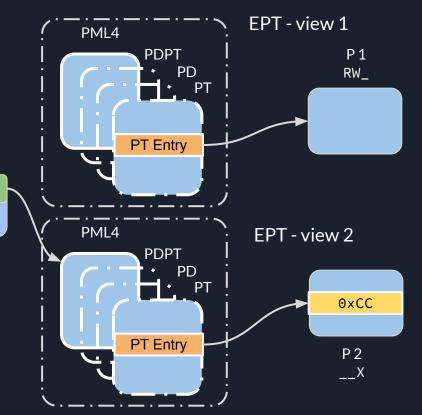
EPT PT

- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{-}/_{X}$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2



EPT PT

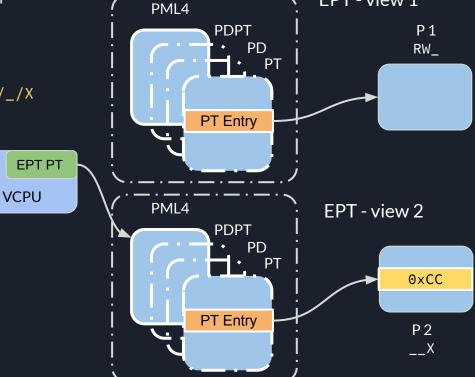
- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{-}/_{X}$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2



- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{-}/_{X}$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2

#int3

• handle breakpoint



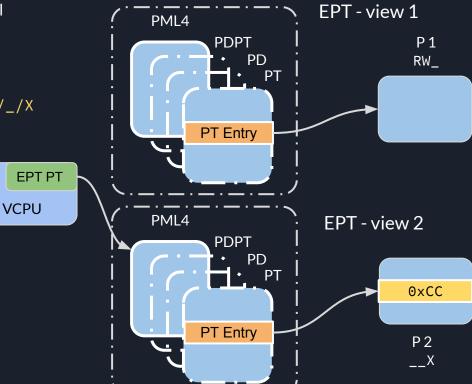
EPT - view 1

- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{/}/X$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2

#int3

handle breakpoint

#EPTViolation (Read/Write)

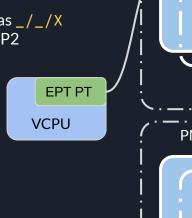


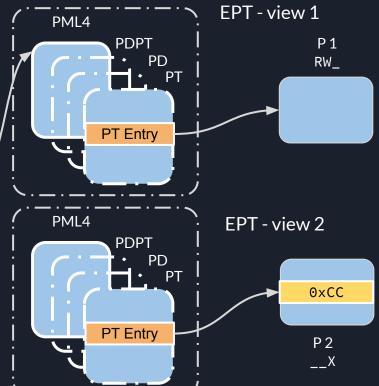
- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{/}/X$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2

#int3

handle breakpoint

#EPTViolation (Read/Write)





- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{/}/X$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2

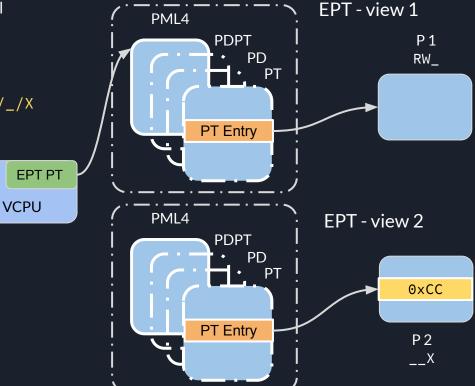
#int3

handle breakpoint

#EPTViolation (Read/Write)

switch VCPU to view 1

#EPTViolation (Execute)



- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_{/}/X$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2

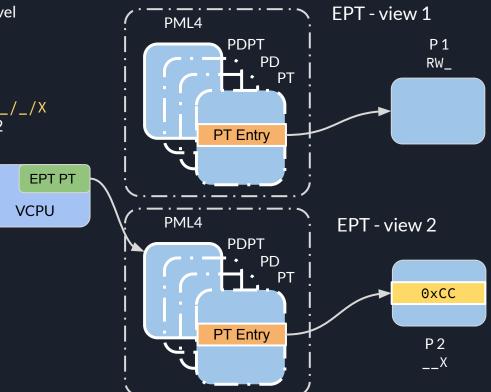
#int3

handle breakpoint

#EPTViolation (Read/Write)

switch VCPU to view 1

#EPTViolation (Execute)



- Switch memory views at VCPU level
- Setup:
 - a. duplicate EPT
 - b. duplicate P1
 - c. insert breakpoint in P2
 - d. set P1 as R/W/ and P2 as $_//X$
 - e. point view 2 PT Entry to P2
 - f. switch VCPU to view 2

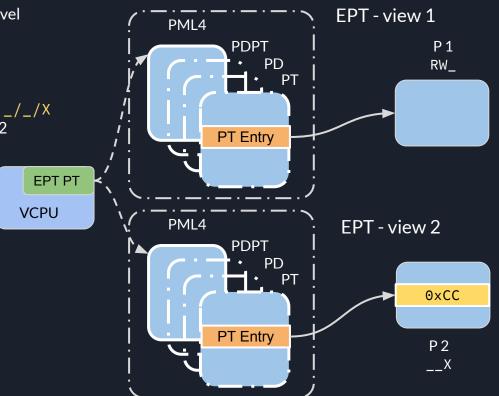
#int3

handle breakpoint

#EPTViolation (Read/Write)

switch VCPU to view 1

#EPTViolation (Execute)



#int3

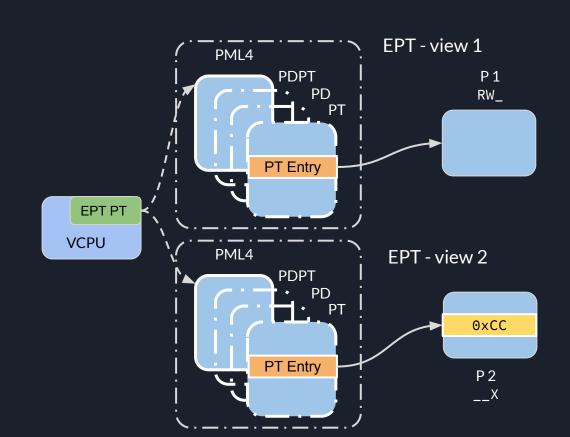
• handle breakpoint

#EPTViolation (Read/Write)

switch VCPU to view 1

#EPTViolation (Execute)

- ✓ fast
- v stealth
- 🗸 safe



#int3

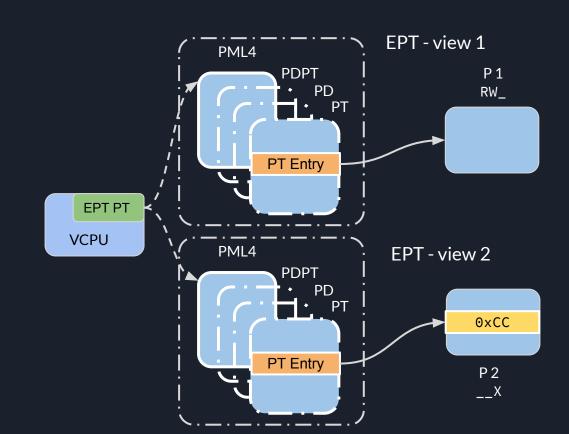
handle breakpoint

#EPTViolation (Read/Write)

• switch VCPU to view 1

#EPTViolation (Execute)

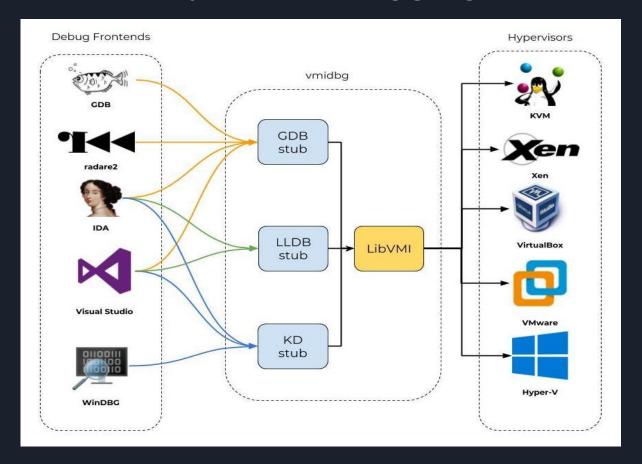
- V fast
- **V** stealth
- ✓ safe



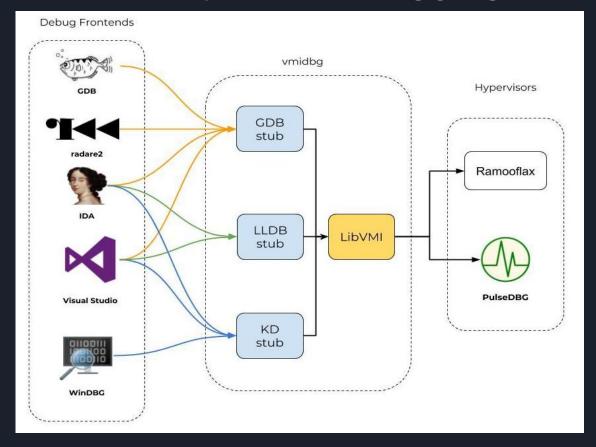
2015 - altp2m implemented in Xen 4.6 EPT - view 1 PML4 **PDPT** P 1 2016 - Available in DRAKVUF sandbox RW_ PT Entry **EPT PT VCPU** PML4 EPT - view 2 **PDPT** fast • PT stealth 0xCC safe PT Entry P 2

Future?

Flexible Full-system Debugging in 202X?



Flexible Boot-sequence Debugging



LibVMI: design and future goals

- Switched to CMake recently
 - o opens Windows drivers
- New ambitions
 - cross-platform
 - adapt to Windows/Linux/MacOS hypervisors
 - o embedded use-case
 - ramooflax/PulseDBG
 - safety concerns
 - runs with very high-privileges (root for Xen APIs)
 - processes untrusted input, likely dangerous (malware)
- Long-term: reimplementation in Rust?
 - to be discussed
 - call to rust developers!



Getting Involved

- Hypervisor developers
 - o open Virtual Machine Introspection API
 - VMware, Hyper-V, VirtualBox

LibVMI

- write a driver
- o improve introspection capabilities
- o performance study for Winbagility integration
- o long-term: rewrite in Rust

Pyvmidbg

- improve GDB stub
- write LLDB/KD stubs
- Windows/Linux internals (scheduling, process state)
- join the Slack! (ask for invites)
 - https://vmidbg.slack.com

Conclusion

- For hypervisor developers
 - o provides a unified debugger interface
 - leverage the full power out of your hypervisor
 - o focus on the right abstractions, together
- For VM users
 - o transform your vision of the Virtual Machines as another process tree
 - keep the same debug environment and tools for multiple OS
- For everyone
 - build hypervisor-level debugging framework as a commodity

Thanks 🙏

- Saar Amar
- @J00ru
- Petr Beneš
- Artem Shishkin
- @aionescu
- @JmpCallPoo
- Tamas K Lengyel

Questions

Building a {flexible} hypervisor-level debugger

https://github.com/Wenzel/pyvmidbg



