

Linux Plumbers Conference 2023

Thara Gopinath, Mickaël Salaün, James Morris



Agenda

Architecture Overview

Hyper-V Implementation

- KVM Implementation (Heki)
- Q&A

Architecture Overview

James Morris

Use virtualization to provide enhanced security for the guest OS, leveraging the hypervisor security boundary.

Protect the integrity of security-critical guest structures.

Prevent bypass of guest security mechanisms and policies.

Support a Trusted Execution Environment (TEE) for running security applications

- Executable code integrity
- Key management
- \cdot Others

Rationale

- \cdot Bring mainline Linux to state of the art.
- Linux is trailing proprietary solutions across Linux and other OSs.
- Attacks continue to evolve, along with motivation levels.

LVBS Architecture

Open-source **architecture** for Linux.

Independent of:

- \cdot ISA
- · Hypervisor
- \cdot VMM
- Security monitor
- TEE implementation

Approach

 Mainline acceptance across ecosystem is critical to success

- Reference implementation (HEKI):
 - · KVM
 - · Linux kernel API
 - Flexible kernel hardening policy
- We are seeking feedback and collaboration.

Hyper-V Implementation

Thara Gopinath

Hyper-V based System





Virtual Secure Mode (VSM)

Separate privileged execution environment within a partition : Virtual Trust Level (VTL)





VSM Features

- \cdot Virtual Processor state isolation
- Memory access hierarchy and protection
- Virtual Interrupt and Intercept handling

<u>https://learn.microsoft.com/en-us/virtualization/hyper-v-on-windows/tlfs/vsm</u>

LVBS and VSM

· Kernel Hardening

- Hardening memory permissions (HVCI)
- $_{\odot}$ Monitoring critical system registers and MSRs
- $_{\odot}$ Monitoring critical kernel data structures
- Offloading policies (control flow integrity, authentication)
- · Offloading secure services (trustlets)

Initial Target : Basic Kernel Hardening.

Threat Model:

Kernel Hardening

- Protect kernel from a user space attacker exploiting a kernel vulnerability
 - Assume that the attacker has arbitrary read write access to guest kernel thanks to exploited vulnerability by malicious
 - User space process
 - Network Packet
 - Block Device
 - · Secure Boot is trusted.
 - Defence in Depth ; but no extra features!!!

H/W Requirements:

- Second Level Address Translation (SLAT, Two-Dimensional Paging, AMD's RVI/NPT)
 - Enable to manage VM memory and add a secondary complementary layer of permissions only controlled by the hypervisor)
- CPU features that allow to differentiate between kernel space and user space memory (MBEC)

(Common Layer)



e.g. Hypervisor Enforced Kernel Integrity (Heki)

Secure Kernel

- \cdot Small TCB
- · Maintainability
- \cdot Ability to support secure interfaces

 Initial choice for secure kernel : Minimal Linux Kernel

Control Interfaces

Synchronous : Explicit VTL Call and Return



Asynchronous : Interrupt based entry and exit Higher VTL gets precedence over lower VTL

Boot

- We trust secure boot !
- $\cdot\,$ VTL0 guest kernel boots up VTL1 secure kernel
- Establish kernel hardening and other policies with secure kernel prior to init process.



Boot Sequence



The Big Picture (Boot)

GUEST PARTITION				
VTL 1	VTLO			
Trustlets	User Applications			
SECURE LINUX KERNEL HVCI VSM-MEM- PROT HV-VSM secure driver HV-VSM secure driver Secure interrupt / intercept handler Optional secure kernel loader HV-VSM-RES- MONITOR	LINUX KERNEL HV-VSM driver Common hypervisor agnostic layer Kernel frameworks (boot, mm, module etc.) HV-VSM-BOOT driver HV-VSM-BOOT			
HW				
	Indicates Hypervisor agnostic layers			

The Big Picture (Late Boot)



Indicates Hypervisor agnostic layers

(Access / Policy Violation)



Indicates Hypervisor agnostic layers

Code

 <u>https://github.com/heki-linux/lvbs-</u> <u>linux/tree/secure-kernel-lvbs</u>

 <u>https://github.com/heki-linux/lvbs-</u> <u>linux/tree/ubuntu-lvbs</u>

KVM Implementation (Heki)

Mickaël Salaün

The Big Picture



Shared frameworks between LVBS-KVM and LVBS-MSHV

RFC v2 patches

Sent <u>RFC v2</u>:

- Guest kernel implementation of the common API
- Two new KVM hypercalls: CR-pinning and memory permission
- KVM interface with the VMM: dedicated VM exits and related capabilities

CR-pinning hypercall

Enforce a bitmask on **control registers** to guard against locked features (e.g. SMEP)

kvm_hypercall3(**KVM_HC_LOCK_CR_UPDATE**, 0, // control register X86_CR0_WP, // flag to pin flags); // options

Can create a **VM exit** on configuration or policy violation for the VMM to be able to do something.

Generate a **GP fault** on policy violation.

Memory protection hypercall Configure (a subset of) EPT permissions.

kvm_hypercall1(**KVM_HC_PROTECT_MEMORY**, pa); // address of a pagelist

The pagelist atomically maps a set of memory ranges with read, write and execute permissions.

Generate a **synthetic page fault** on policy violation.

Executable permission(s)

<u>Issue</u>: efficiently enforce restriction on kernel executable pages without impacting access to user space pages

<u>Solution</u>: leverage Intel's Mode Based Execution Control (**MBEC**)

Split the execution permission into:

- Kernel mode execution
- User mode execution

Kernel memory permissions without MBEC

and radata	executable	OxFFFF
end_rouata	read-only	
vdso_end	read-execute	
vdso_start	read-only	
start_rodata	executable	
_etext	read-execute	
_text	executable	
		0x0000

Kernel memory permissions with MBEC

evel ve dete	non-executable	0xFFFF
end_rodata		
vdso_end	read-only	
vdso_start		
start rodata		
	non-executable	
_etext	read-execute	
_text	non-executable	
		0x0000

Code

https://github.com/heki-linux/linux

branch heki-v2

Wrap up

KVM and Hyper-V supports:

- defense-in-depth mechanism leveraging hardware virtualization
- common API layer across hypervisors

Any feedback?

https://github.com/heki-linux

Q&A

Thank you



<u>Demo: control-register pinning</u> (SMEP)

```
user@heki-host$
static void heki_test_cr_disable_smep(struct kunit *test)
       unsigned long cr4;
       /* SMEP should be initially enabled. */
       KUNIT_ASSERT_TRUE(test, __read_cr4() & X86_CR4_SMEP);
       kunit_warn(test,
                  "Starting control register pinning tests with SMEP check\n")
        * Trying to disable SMEP, bypassing kernel self-protection by not
        * using cr4_clear_bits(X86_CR4_SMEP).
       cr4 = __read_cr4() & ~X86_CR4_SMEP;
       asm volatile("mov %0,%%cr4" : "+r"(cr4) : : "memory");
```

```
/* SMEP should still be enabled. */
KUNIT_ASSERT_TRUE(test, __read_cr4() & X86_CR4_SMEP);
```

62,2-9 25%