# Security Analysis of Encrypted Virtual Machines

Felicitas Hetzelt

Technical University of Berlin
Berlin, Germany
file@sec.t-labs.tu-berlin.de

Robert Buhren

Technical University of Berlin
Berlin, Germany
robert@sec.t-labs.tu-berlin.de

## Abstract

Cloud computing has become indispensable in today's computer landscape. The flexibility it offers for customers as well as for providers has become a crucial factor for large parts of the computer industry. Virtualization is the key technology that allows for sharing of hardware resources among different customers. The controlling software component, called hypervisor, provides a virtualized view of the computer resources and ensures separation of different guest virtual machines. However, this important cornerstone of cloud computing is not necessarily trustworthy or bug-free. To mitigate this threat AMD introduced Secure Encrypted Virtualization, short SEV, which transparently encrypts a virtual machines memory.

In this paper we analyse to what extend the proposed features can resist a malicious hypervisor and discuss the trade-offs imposed by additional protection mechanisms. To do so, we developed a model of SEV's security capabilities based on the available documentation as actual silicon implementations are not yet on the market.

We found that the first proposed version of SEV is not up to the task owing to three design shortcomings. First the virtual machine control block is not encrypted and handled directly by the hypervisor, allowing it to bypass VM memory encryption by executing conveniently chosen gadgets. Secondly, the general purpose registers are not encrypted upon *vmexit*, leaking potentially sensitive data. Finally, the control over the nested pagetables allows a malicious hypervisor to closely monitor the execution state of a VM and attack it with memory replay attacks.

*Keywords*  Secure Encrypted Virtualization, AMD SEV, Cloud Computing

## 1.  Introduction

Cloud computing has been one of the most prevalent trends in the computer industry in the last decade. It offers clear advantages for both customers and providers. Customers can easily deploy multiple servers and dynamically allocate resources according to their immediate needs. Providers can ver-commit their hardware and thus increase the overall utilization of their systems. The key technology that made this possible is virtualization, it allows multiple operating systems to share hardware resources. The hypervisor is responsible for providing temporal and spatial separation of the virtual machines (VMs). However, besides these advantages virtualization also introduced new risks.

Customers who want to utilize the infrastructure of a cloud provider must fully trust the cloud provider. Especially the hypervisor is a critical component provided by the cloud hoster as it has full control over the guest VMs. A malicious or compromised hypervisor can read and write the entire guest memory. This affects the integrity and confidentiality of the customer's secrets and the integrity of the customer's services. Security issues might lead to a full breach of the hypervisor through a hosted VM; bugs of such severity have been reported for all most commonly used hypervisors [10–13, 16]. As a single cloud instance often hosts multiple guest VMs from different customers such security issues allow a malicious tenant to steal confidential data from other customers.

Intel's Software Guard Extensions (SGX) [15] and AMD's Secure Encrypted Virtualization (SEV) [19] are industry's answer to these threats. They extend the features of the processor to reduce the impact of a malicious, higher privileged software in regards to the confidentiality and integrity of lower privileged software. SGX enables the customer to create a secure enclave where special code can be executed in a trusted environment that cannot be tampered with by the hypervisor or the operating system. SGX achieves this by requiring the customer to identify the security sensitive parts of a program and to alter them such that these parts are executed in an SGX enclave. SEV, on the other hand, allows a customer to encrypt the unaltered VM's memory so that the hypervisor is not able to inspect its data. The recent addition of SEV Encrypted State (SEV-ES) extends the cryptographic

protection of the guest VM to its control state and its general purpose registers. As can be seen from the AMD SEV whitepaper [19]:

> *„SEV technology is built around a threat model where an attacker is assumed to have access to not only execute user level privileged code on the target machine, but can potentially execute malware at the higher privileged hypervisor level as well. The attacker may also have physical access to the machine including to the DRAM chips themselves. In all these cases, SEV provides additional assurances to help protect the guest virtual machine code and data from the attacker."*

The advantage of a solution such as AMD's SEV is that it can be easily adopted by customers because no changes to their existing application software are needed.

While the research community has examined Intel's SGX [9, 26, 28], AMD's SEV has not been subject to scientific research so far. It is thus unclear what level of protection SEV can provide. In this paper, we have a first look at the upcoming AMD SEV technology based on publicly available documentation. We identify possible design issues that can be leveraged by a malicious hypervisor to compromise the guest VM. To that end, we implement in total three proof of concept attacks on a currently available system. For the construction of the attacks, we bear in mind not only the restrictions an AMD SEV-enabled system imposes, but also evaluate how the initial SEV design could be hardened without sacrificing further guest transparency or impacting cloud maintenance operation. However we show that even an attacker restricted to basic resource management capabilities, is still able to gain access to the protected guest system. Our contributions are:

- We show how a malicious hypervisor can coerce the guest to leak arbitrary memory content and perform arbitrary write operations on encrypted memory.

- We describe how to completely disable any memory encryption configured by the tenant.

- We implement a replay attack that uses captured login data to gain access to the target system by solely exploiting resource management features of a hypervisor.

For the first two attacks we base our evaluation on the initial design of SEV without the optional SEV-ES extension. The third attack considers the protection of the guest state and general purpose register, as it might be implemented by SEV-ES. As processors featuring SEV are not available yet, it is unclear whether our results will apply to real silicon implementations or future versions of SEV. We therefore emphasize that we did not break AMD SEV itself but rather evaluated the design issues present in the documentation with respect to their capability to protect a guest VM against a malicious or compromised hypervisor.

The rest of the paper is structured as follows: In Section 2 we give an overview on x86 virtualization and AMD SEV. We evaluate the security of the protection mechanisms proposed by SEV in Section 3 and discuss our attack model in Section 4. In Section 5 we present our attack. We discuss possible mitigations to our attack in Section 6. In Section 7 we evaluate alternative approaches to shield execution environments from higher privileged adversaries and present related attacks under similar threat models. Finally, we discuss future work in Section 8 and conclude our work in Section 9.

## 2. Background

In this section we first give a brief introduction to x86 virtualization then we discuss the design of AMD's SEV technology. This information by no means represents a complete overview of these topics. The specification for AMD SVM and AMD SEV are however publicly available. Thus, we refer the interested reader to [2, 4].

### 2.1 x86 Virtualization Technologies

In 2005, both Intel (VT-x) [25] and AMD (SVM) [1] introduced hardware extensions to their x86 processors that added a higher privileged mode to the existing ring 0 to ring 3 privilege levels. This new mode, called host mode, comprises another set of the privilege rings 0 to 3 and is higher privileged than the non-host mode, called guest mode. The host mode is intended to host the hypervisor whereas the guest VM usually executes in the non-host mode. To make use of these extensions, a hypervisor, running in the host mode, uses a special instruction, `vmrun`, to switch the CPU to the guest mode. This instruction takes the address of a control structure as a single argument in the register `rAX`. This control structure, called `vmcb`, contains the guest state, entry controls (pending virtual interrupts) and exit controls. Prior to the initial start of the guest, the hypervisor configures the `vmcb` and initializes the general purpose registers as they are not part of the `vmcb`. Upon issuing the `vmrun` instruction, the CPU copies the values of the `vmcb` fields into the respective hardware registers and starts execution of the guest at the entry point defined in the `vmcb`. An event that is flagged in the `vmcb` as such will lead to a `vmexit` with the exit reason set in the `vmcb`. The hypervisor then handles the exit accordingly.

While the original design of AMD SVM from 2005 allowed a hypervisor to run multiple guests on a single CPU without altering the guest OS, it lacked support to virtualize memory efficiently. In 2008, AMD released a technology called "nested-paging" [3] that enables a hypervisor to virtualize memory in an efficient way. The traditional paging hierarchy was extended with another layer, the nested layer. Instead of just translating from virtual to physical addresses, now the translation involves two steps. The guest pagetable, maintained by the guest operating system, translates from guest virtual to guest physical addresses, whereas

the host pagetable translates from guest physical addresses to physical addresses. This second translation step is fully under control of the hypervisor.

## 2.2 Memory Management

KVM leverages nested page faults as an indicator for the access load of a guest memory page. This information is required to optimize tasks in which guest memory has to be made temporarily unavailable, namely live migration, memory snapshots, and memory overcommitment. Here we give a brief overview of how KVM incorporates nested page fault information in those tasks.

For live migration and memory snapshots, memory is transferred incrementally. First, the current guest memory content is copied without stopping the guest execution. This memory snapshot is extended in later increments by memory content that has been modified during the initial transfer. Only if the estimated remaining transfer time falls below a predefined threshold, the guest is stopped to transfer the remaining pages. To identify modified pages since the last transfer, QEMU instructs KVM to record a list of all pages that have been written to by the guest. To that end, KVM removes the write access permission from all guest memory pages and restores it only after registering the preceding write access fault.

Further, guest memory is subject to Linux standard memory maintenance operations. Based on process page faults, memory is allocated lazily and can be swapped out to disk if unused, thereby enabling memory overcommitment. To do so the nested pagetable entries for the respective pages are removed, therefore causing a nested pagetable violation if the guest tries to access them. KVM handles the nested page fault, restores the page from the swap file and recreates the nested pagetable entries along with the pagetable entries connecting to the QEMU userland process.

A similar method is used to lazily allocate memory for the guest upon initial startup. To associate a memory region with the guest QEMU allocates a chunk of memory and informs KVM of the virtual memory area. Until a page is accessed no host- or nested pagetable entries are created other than those required for the guest kernel image and initial runtime. If the guest accesses any additional page a nested page fault will be triggered and handled by KVM as described in the previous paragraph.

## 2.3 Virtual devices

While the hardware virtualization extensions provide CPU and memory virtualization, handling device virtualization is the obligation of the hypervisor.

On x86, devices are accessed by either IO ports, memory mapped registers or by a combination of both. Accessing IO port based devices requires the use of special instructions (e.g. IN or OUT) whereas memory-mapped devices can be accessed using normal instructions (e.g. mov). If the device itself requires the CPU to handle an event, it raises an interrupt

which diverts the control flow of the CPU to a specific interrupt handling routine. To improve the overall performance, data can also be transferred without the involvement of the CPU. The device reads or writes directly to or from main memory, allowing the CPU to perform other tasks in parallel. The technology is commonly referred to as DMA (Direct Memory Access). To protect against unauthorized accesses from DMA capable devices, an IOMMU can confine devices to only access configured memory regions. Three common approaches to handling devices in a virtualized environment are passthrough, emulation, and para-virtualization.

***Passthrough*** Using this method, one VM has exclusive access to a hardware device. If the device provides only a memory-mapped interface, the corresponding memory pages are mapped into the guest address space via the nested pagetable. In the case of IO ports, the vmcb allows configuring which IO ports are accessible directly by a guest. If a commodity device without special virtualization extensions is passed through, only a single guest can use this device. Some devices can also be configured to provide „virtual functions", through which the same device can be used by multiple guests. An IOMMU is usually required to contain DMA access within configured memory regions.

***Emulation*** The hypervisor can present a virtual device to the guest. It sets up the nested pagetable with a hole in the address space where the guest expects the memory-mapped device. When the guest now accesses these memory ranges to interact with the device, this will trap into the hypervisor. To perform memory access on behalf of the guest the hypervisor must know the value that should be written. The vmcb will contain the fault address, i.e. the location where the data should be written, but not the value itself. The value is usually stored in a general purpose register[1]. The hypervisor must parse the instruction that caused the fault to identify the register holding the value. As the instruction pointer locating this instruction holds a guest-virtual address, the hypervisor must first traverse the guest pagetable to get the guest-physical address of the instruction before it can parse the instruction. As traversing the pagetable imposes a severe bottleneck for device emulation, AMD added decode assists that provide the register location of the value in case of a nested page fault.

***Para-virtualization*** The performance for accessing virtual devices can be enhanced using para-virtualization. Here the hypervisor does not emulate an existing device but provides an interface of an artificial device to the guest that has no corresponding hardware device. This has the advantage that the hypervisor and the guest can agree on an interface that encompasses the peculiarities of the hypervisor and guest communication. For example instead of trapping writes to

---

[1] There are instructions like rep ins movs that take the target and source address as pointers in registers, but those are not commonly used when accessing memory-mapped devices.
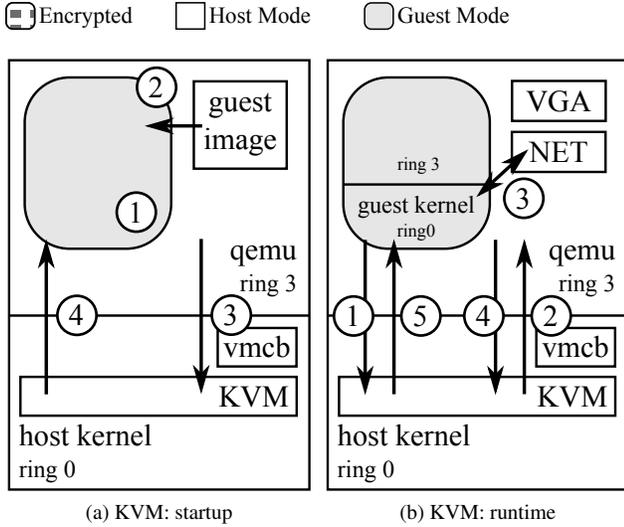
(a) KVM: startup          (b) KVM: runtime

Figure 1: QEMU/KVM architecture



(a) KVM/SEV: startup       (b) KVM/SEV: runtime

Figure 2: SEV-enabled QEMU/KVM architecture

certain memory areas, the guest can use special instructions that cause a trap into the hypervisor. This mechanism is called a hypercall. In contrast to memory accesses, these hypercalls do not cause a pagetable walk by the memory management unit. This can increase the performance of these virtual devices but requires drivers to be adapted to the hypercall interface.

Device emulation is crucial for providing basic VM functionality, like network connectivity, for the guest owner. Besides device emulation, features that are the sole responsibility of the cloud providers, such as host memory management and migration, are mandatory in a cloud environment. Given that substantial modification of the deployed VM's code base goes against customer interest and device passthrough does not scale to larger cloud infrastructures, this draws a lower bound on the limitations imposed on hypervisor control over guest VMs, i.e., host memory management demands that the hypervisor has control over the second level page translations.

## 2.4 Linux KVM

In the previous paragraph, we explained AMD's virtualization extensions. We now lay out how this technology is used by the KVM hypervisor which is integrated with the Linux kernel [20].

Virtualizing CPU and memory is not sufficient because guest operating systems also need devices such as video output, network or block devices. As a guest should not directly interfere with the hardware devices itself, they must be either multiplexed or emulated (see Section 2.3). While the KVM hypervisor is responsible for controlling the execution of guest VMs, QEMU is leveraged to handle the device virtualization. Figure 1a depicts the initial startup of a guest VM in a KVM/QEMU setup.
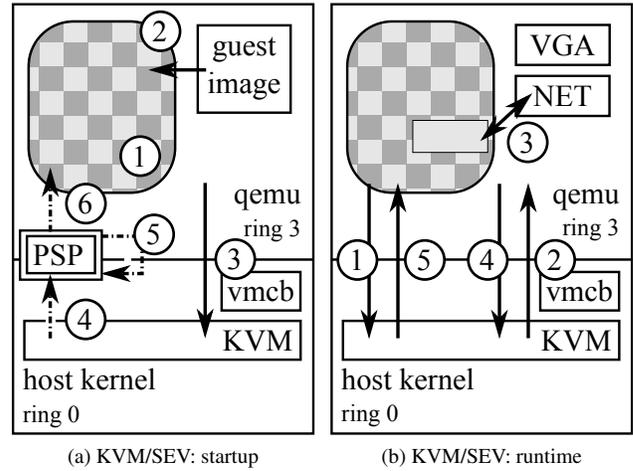
First, QEMU reserves memory for the VM (Figure 1a ①). Then it copies the guest binaries into this reserved memory (Figure 1a ②). By using the /dev/kvm device node, the KVM module of the Linux kernel is instructed to start a new VM (Figure 1a ③). KVM then sets up a vmcb data structure incorporating the information from QEMU and issues the vmrun instruction to start the VM (Figure 1a ④). The processor now enters the guest mode, depicted in grey, and starts execution at the entry point defined in the vmcb.

The runtime behaviour is shown in Figure 1b. Upon any event that was configured in the vmcb to cause a vmexit, the CPU leaves guest mode and enters host mode again with a specific error code set in the vmcb (Figure 1b ①). The KVM module can then either handle the exit itself or, in the case of, e.g., a memory-mapped IO operation to an emulated device, can return to QEMU which then handles the request (Figure 1b ②). The emulated device can access guest memory directly to mimic DMA memory transfers (Figure 1b ③). After the request was served, QEMU calls KVM again (Figure 1b ④), which resumes execution of the VM in guest mode (Figure 1b ⑤).

## 2.5 AMD SEV

As indicated in Figure 1, the hypervisor has full access to guest memory while the CPU is in host mode. This demands that a cloud customer must trust not only the employees of the cloud provider but also the integrity of the hypervisor. Bugs such as [10–13, 16] can be used by a malicious tenant to attack the hypervisor itself and thereby gain access to assets of other tenants residing on the same physical machine.

SEV protects guest memory via encryption. The guest specific memory encryption key will never be exposed to the hypervisor. It is only accessed by a secure coprocessor and the memory controller that handles the encryption and de-

cryption transparently. The coprocessor which was added to SEV-enabled CPUs (the „Platform Security Processor" [19], indicated as *PSP* in Figure 2a), handles key management and is responsible for configuring the correct guest key within the memory controller.

Figure 2 shows how the classical KVM architecture looks on an SEV-enabled system. Like detailed in the previous Section, QEMU communicates with the KVM module to prepare the VM for launch (Figure 2a ① to ③). To enable SEV for the newly allocated VM it's memory must first be encrypted. The host kernel calls the coprocessor to initiate the encryption of the VM memory using a three-fold command sequence, LAUNCH_START, LAUNCH_UPDATE and LAUNCH_FINISH (Figure 2a ④, ⑤ and ⑥). By using this command sequence, the hypervisor ensures that the firmware generates an encryption key unique to the VM (LAUNCH_START), encrypts the memory and records a launch receipt of the VM used for remote attestation (LAUNCH_UPDATE and LAUNCH_FINISH). After the encryption of guest memory is completed the firmware provides the recipe to the hypervisor to be passed on to the customer. This recipe includes measurements of the guest image and platform authentication data, which allow the customer to verify that the VM memory was encrypted and initialized correctly. If a customer judges the recipe or the contained measurements to be faulty, he can choose to withhold the provisioning of secrets to the VM.

Each VM uses its unique cryptographic key that is loaded by the secure processor when the corresponding VM is scheduled. Once a guest enables paging, it can mark individual data pages as either *shared* or *private* by setting a physical address bit (the enCrypted- or C-bit) in its pagetable. Memory pages marked as *private* are encrypted using AES with the guest specific key and pages marked as *shared* are either not encrypted or encrypted with the hypervisor key and can thus be used to exchange data with the hypervisor. The C-bit of the guest pagetables has precedence over the C-bit of the hypervisor controlled second level pagetables to secure the page protection configured by the guest VM. In addition to the memory protection mechanism, AMD offers tenants the ability to enforce guest policies. Policy configuration includes amongst others the option to disable debug capabilities of the hypervisor towards the guest VM.

Figure 2b shows the system configuration during runtime. The secure coprocessor (*PSP*) is not shown, as it is used mainly during VM startup. The steps composing the runtime behaviour under SEV (Figure 2b ①, ②, ④ and ⑤) do not differ from the non-SEV configuration. This is due to the fact, that cryptographic operations are handled transparently by the memory controller, while the secure processor handles key management without the involvement of the hypervisor or the VM. However to facilitate DMA memory transfers (Figure 2b ③) similar to a classical setup, the guest

is tasked to configure shared memory regions, which are exempt from encryption.

## 2.6   SEV - Encrypted State

SEV - Encrypted State (SEV-ES) is an extension to SEV that additionally encrypts the guest state, including the general purpose register, using the guest specific encryption key. When the CPU leaves the guest mode, all general purpose registers, as well as the guest saved state, are encrypted. A vmexit event is now classified as either an "Automatic Exit" (AE) or a "Non-Automatic-Exit" (NAE) depending on the exit reason. Any asynchronous event, e.g. an interrupt, is classified as an AE. AE events do not require the hypervisor to read the guest state, which can therefore be encrypted by the secure processor. AES events, on the other hand, are events that potentially require the hypervisor to read the guest state. If such an event occurs, the control is not transferred to the hypervisor. Instead, a new exception is raised in the guest, the "VMM Communication Exception" (#VC). The exception handler of these exceptions in the guest can now decide whether to provide the hypervisor with access to the guest state, or not. To exchange data a shared memory region is used, called "Guest-Hypervisor Communication Block" (GHCB). Data to be shared with the hypervisor must be copied to that memory region. After copying the values to the GHCB, the guest executes vmgexit to transfer control to the hypervisor. Nested page faults are usually AE events, i.e., the hypervisor does not need to read guest state to handle these events. Therefore there is no #VC exception triggered in these cases. In order to allow the hypervisor to provide emulated devices for the guest, the hypervisor can enforce certain nested page faults to be NAE events. To do so, the hypervisor sets a reserved bit in the nested pagetable. Faults caused by accesses to these pages are treated as NAE events, and instead of transferring control to the hypervisor, a VC exception in the guest is raised.

## 3.   AMD SEV Security Considerations

While guest memory is protected from direct hypervisor access by encryption, other security-critical components are not protected at all. By examining the AMD SEV documentation [2, 19] and publicly available comments from AMD employees [21], we found that for a system without SEV-ES:

1. The general purpose registers are not encrypted upon a vmexit [21].

2. The vmcb is subject to manipulation by the hypervisor [21].

3. The hypervisor can access encrypted guest memory due to the lack of memory authentication [2, 19].

Under a system, which also implements the SEV-ES extension only the third point remains valid.

***General Purpose Registers*** Whenever the CPU switches from guest mode to host mode, the general purpose registers of the guest are exposed to the hypervisor. As the guest itself cannot control when the CPU transfers to host mode, these registers can contain potentially confidential data. If such an exit occurs e.g. while the guest is generating an RSA key pair, the key components might be exposed to the hypervisor.

***VMCB*** As mentioned in Section 2.1, the `vmcb` is used to control the execution and state of the guest. The `vmcb` is therefore crucial to guest integrity and exposes the content of privileged guest registers. Among these registers is the instruction pointer of the guest which allows the hypervisor to govern guest control flow.

***Memory Authentication*** The memory is encrypted, but it is otherwise not protected from access. This enables the hypervisor to inject faults into the guest or to capture and replay private guest memory.

Later sections will lay out how these design issues can be leveraged by a malicious hypervisor to a) gain shell access to a guest, b) read protected guest memory and c) fully revert any memory protection configured by the tenant. While SEV-ES successfully protects the general purpose registers and the `vmcb`, the guest memory can still be accessed by the hypervisor, though the hypervisor can only access encrypted pages. This mitigates attack vectors b) and c). Still, as we will show in later sections there is no easy way to prevent a) without sacrificing guest transparency or impacting classic cloud functionality, such as migration and the memory management features of the hypervisor.

## 4. Attack Model

In this section, we describe our attack model, which is based on the AMD SEV security properties (detailed in Section 3).

We assume that a customer successfully deployed his VM on an AMD SEV-enabled system. During startup, we also assume that the hypervisor is uncompromised and compliant with the AMD SEV specification [2]. This means the customer was able to attest the correct setup of his VM using the receipt provided by the hypervisor. From this point on the VM is protected by AMD SEV. Neither the hypervisor nor someone with physical access to the cloud infrastructure is able to read the designated private memory regions of the protected guest.

Then, during runtime, an attacker was able to compromise the hypervisor, thereby gaining root access to the host system. The described scenario is likely, as incidents of the past show [10, 12, 13, 16]. We also assume that the attacker has knowledge of the target system with regards to the versions of the kernel and userland processes. As the cloud provider itself often provides the guest images, this is also likely. We assume that the encryption scheme in use produces the same encrypted data if the input, key and host physical address are identical, similar to other symmetric linear memory encryption schemes. Further, we require, that no integrity check is performed on encrypted data, In addition to that, we initially assume access to nested pagetables, *vmcb* and guest register state, which we later restrict only to nested pagetable access for the replay attack.

## 5. Attacks against Encrypted Virtualization

We now present three attacks against VMs under a compromised hypervisor.

The first two attacks presented in Section 5.1 are directed against the proposed design of SEV without the optional feature SEV-ES, which allows the hypervisor to extract and control guest state through the unencrypted guest control block and registers. Amongst other security concerns for the tenant, this flaw can be used to decrypt guest memory including the internal address mapping, as we will show in our first attack in Section 5.1.1. Building upon this capacity we describe in Section 5.1.2 how the memory protection configured by the tenant can be deactivated, without notifying the guest. After deactivation of memory protection further exploitation, like arbitrary code execution is trivial to execute as the hypervisor now has full access to guest memory.

The third attack already takes the presence of SEV-ES into account, which protects the integrity of the guest control block and registers in the face of a malicious hypervisor. We however show in Section 5.2 that protecting those structures alone is not sufficient. If the hypervisor is in control over guest memory allocations through nested paging, it can use this capability to launch a replay attack. We prove this claim by launching an attack against an OpenSSH server running in the protected guest VM to gain access at potentially high privilege levels.

### 5.1 Attacks based on exposed Guest State

In this section, we present two attacks against an encrypted guest, facilitating hypervisor access to the guest control block and registers. First, we describe a method to exploit guest control flow to read and write arbitrary memory areas of a running guest in decrypted form. Based on this primitive, we construct an advanced attack to disable guest memory protection as documented in [2] altogether.

### 5.1.1 Accessing Protected Memory

Given a system which is capable of encrypting guest memory as described in [2], we now describe how a malicious hypervisor can coerce an encrypted guest into leaking arbitrary memory content. The methods for reading and writing protected guest memory are symmetric, therefore we restrict this section to the description of the memory read primitive.

During guest execution, the memory of the active VM is transparently decrypted by the memory controller. Memory content which, in this state, is transferred into unencrypted areas like the `vmcb`, registers or shared memory, will be ex-

```
mov edi, dword ptr [rbx]
hlt
```

Listing 1: Read Instruction Sequence

```
int new_handle_hlt(struct vcpu* vcpu) {
    u64 rip, edi;
    rip = rip_read(vcpu);
    if(rip == DECRYPT_HLT_INS && decrypting) {
        edi = register_read(vcpu, VCPU_EDI);
        // process decrypted data in EDI ...
        register_write(vcpu, VCPU_RBX,
            current_addr);
        current_addr += 0x4;
        if(current_addr < last_addr)
            rip_write(DECRYPT_HLT_INS);
        return 0;
    } else {
    // handle normal hlt exit ...
    }
}
```

Listing 2: HLT Exit Handler

posed to the hypervisor whenever guest execution is interrupted. Our attack induces an interruption of the guest execution, right after protected data has been transferred from an attacker controlled memory location into an unencrypted register. To divert guest control flow, we set the guest instruction pointer before guest re-entry to the guest virtual address of a suitable instruction sequence. Shortly after the read instruction we force a vmexit to read the decrypted data from the register.

The instruction sequence is required to end with a trappable instruction and to contain an indirect memory read. Listing 1 shows the sequence of instructions, which we used to launch the attack. We extracted this sequence statically from the guest kernel binary, for which we used a modified tool for ropchain generation, called ROPGadget [18]. The code snippet reads four bytes from guest memory into the register eDI, before a vmexit is induced by the instruction hlt. The malicious hypervisor can then conveniently take the decrypted word from the general purpose register. Listing 2 shows the respective exit handler, which the hypervisor could use to handle this particular hlt trap condition. To decrypt an arbitrary section of guest memory, the exit handler re-sets the guest instruction pointer to the guest virtual address of the instruction sequence and the guest register rBX to the guest virtual address of the protected memory to be read.

The diversion of guest control flow can be initiated at any point during host execution. To resume normal guest execution after the attack, the guest registers which are clobbered by the decryption are saved in the host environment and later restored after the final memory element has been read.

*Locating the Instruction Sequence*   The recent introduction of kernel address space layout randomization (*KASLR*) complicates our attack. Now the instruction sequence cannot simply be obtained from the guest kernel binary. Instead, we only obtain the offset of the sequence within the kernel text section via static analysis. The offset is then added to the dynamic load address of the kernel text section, which is randomly initialized during the boot process of the VM. To compute the load address of the kernel text section, we compare control registers exposed through the guest control block, pointing to kernel functions inside the guest's virtual address space. Specifically, we subtract the virtual address of the system call entry function entry_SYSCALL_64 of the running guest from the system call entry address of the non randomized kernel image.

### 5.1.2 Disabling Memory Protection

In this section, we describe how encryption can be disabled for individual guest memory pages or even for the complete guest memory space. The attack is based on manipulation of guest internal pagetable entries.

First, we will describe how we access those entries, even though they are assumed to be located in private guest memory and thereby to be encrypted. To access pagetable entries within the protected guest, we first read the physical pagetable address of the currently active process from the cr3 register value stored in the vmcb. We then use the method described in the previous section to read pagetable entries from protected guest memory. As the read primitive can only operate on guest virtual addresses, we access the pagetable data via the direct physical memory map (referred to as *physmap*). The *physmap* is a contiguous mapping of the physical RAM into the virtual address space of the kernel. The virtual base address of the *physmap* is stored in the kernel variable page_offset_base which is located at a constant offset from the dynamic load address of the kernel's text section. We use the read primitive with the adjusted offset of the page_offset_base variable to read its value from guest memory. To access the pagetable entries, we add the physical pagetable base address to the virtual base address of the *physmap*. Using the write primitive, we are now able to overwrite and add pagetable entries arbitrarily, using the adjusted guest virtual address of the pagetable base as the target location.

Even though we are now able to change the pagetable bit controlling page encryption, clearing the C-bit from a guest pagetable entry will only disable the transparent de- or encryption on subsequent memory read or write accesses. Thus the attacker is required to allocate new unprotected (shared) pages of memory and to copy the protected (private) data into the newly allocated areas.

Figure 3 gives an overview of how an attacker can deactivate the protection of guest pages under an arbitrary guest pagetable entry, without notifying the guest. The process can be split up in two phases, first, duplication and then, replace-

Private Page ......... Encrypted

Shared Page ......... Plain

GVA    GPA    HPA

(a) Duplication

Shared
~~Private~~ Page ......... Encrypted

Shared Page ......... Plain
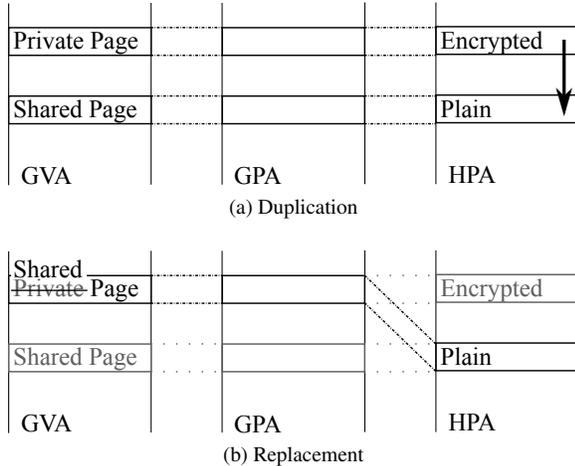
GVA    GPA    HPA

(b) Replacement

Figure 3: Pagetable Modification

ment. In the duplication phase, the protected data is transferred into newly allocated memory as seen in Figure 3a. During the replacement phase, the guest pagetable entry is modified to deactivate the protection, while the nested pagetable entry is redirected to the new data, as shown in Figure 3b.

Now to actually decrypt an amount of guest pages, the according amount of pages has to be reserved in host memory. Then using the read primitive, the guest pagetables are browsed for an unallocated slot matching the original entry level. Similarly, an empty slot in the second level pagetables is located. New pagetable entries are created from the host physical, guest physical and guest virtual addresses. The write primitive has to be used to add the entry connecting the guest virtual to the guest physical address in the guest pagetables. Using the read primitive again, the protected guest memory is read and directly written into the newly allocated area by the host. Finally, the pagetable entry of the original protected mapping is modified to clear the C-bit, while the nested pagetable entry is redirected to point to the newly written data.

## 5.2    Attack based on Nested Pagetable Control

Systems that implement the optional SEV-ES feature will be protected against the previously mentioned attacks. Based on these revised security properties we construct a third attack, which relies only on control over nested pagetable structures and interrupt injection capabilities of the hypervisor. We leave the detailed discussion about the necessity of the latter two capabilities for guest transparent VM encryption in a cloud environment to Section 6.

SEV-ES limits the hypervisor's control over a guest. Namely, the feature restricts access to the VM control block and forces the encryption of guest general purpose registers upon a `vmexit`. Here we explain the protection achieved

through deploying these mitigations to motivate the next attack.

Limiting access to the VM control block prevents the execution of the previous attacks on several levels. The leakage of kernel function pointers is prevented; therefore the guest internal address mapping is not revealed. Whether an instruction like `hlt` traps into host mode is controlled via a bitmap contained in the `vmcb`. Therefore the number of instruction sequences suitable for misuse as read and write primitives can be limited by controlling the configuration of this bitmap. Further, the capability of the hypervisor to manipulate guest control flow is restricted, as the instruction pointer, which is also part of the `vmcb`, can be protected from malicious modification. The encryption of guest control registers will handicap the application of read and write primitives by impairing the control over the address of accessed memory as well as the exposure of the decrypted data. We argue that limiting hypervisor control over physical memory assignment would prevent memory overcommitment as well as any dynamic load balancing or migration efforts. In fact, we assume this capability to be crucial in a cloud environment. Comparably critical is the ability to inject virtual interrupts for device virtualization.

We now describe how a malicious hypervisor can launch a replay attack against a VM running in a protected environment, which uses the optional SEV-ES [4] feature in addition to the protection mechanisms proposed by SEV [2].

First, we give a brief overview of replay attacks and explain how we can attack an OpenSSH server running in an unprotected guest by replaying login credentials. Next, we describe how we can infer the correct location and time to capture and replay guest memory *without* insight into the guest memory content, by observing memory access and system call patterns of the guest. Finally, we describe the steps necessary to implement the attack against an encrypted VM. We conclude with an evaluation of the presented attack.

### 5.2.1    Replay Attacks

On a high level, replay attacks exploit the lack of data versioning and authentication, which allows an attacker (in our case a malicious hypervisor) to eavesdrop on the exchange of valid authentication tokens and replay them to pose as the original communication partner. For OpenSSH, we identified the function `userauth_passwd`, shown abbreviated in Listing 3, as a suitable target. In line 5 a password string is read from the network buffer via `packet_get_string` and stored on the stack. The password string is then validated at line 9 by `auth_passwd`. After validation, the `password` is removed from memory at line 11.

To launch the attack against the OpenSSH server executing in a *unencrypted* guest, the hypervisor captures the guest page containing the credential data in between lines 5 and 9. The attacker then initiates a new connection. After the server receives credentials from the attacker controlled client, the hypervisor replaces the invalid credentials of the attacker,

with the data captured in the previous step. The validation of the restored password will then succeed and thereby grant access to the attacker controlled client at the privilege level of the connecting user.

```
 1 static int userauth_passwd(Authctxt* authctxt) {
 2        char *password, *newpass;
 3        // ...
 4        change = packet_get_char();
 5        password = packet_get_string(&len);
 6        // ...
 7        if (change)
 8              logit("password change not
                      supported");
 9        else if (PRIVSEP(auth_password(authctxt,
            password)) == 1)
10              authenticated = 1;
11        memset(password, 0, len);
12        xfree(password);
13        return authenticated;
14 }
```

Listing 3: userauth_passwd

### 5.2.2 Inferring Memory Content

In a classical replay scenario, the hypervisor can monitor memory content to identify location and state of the memory region to be captured and later replayed. If the guest memory is encrypted, the main challenge is to infer those parameters indirectly.

In this section, we describe how we identify when and where to capture and later replay a memory page without insight into its content. The key intuition behind our approach is that memory content can be inferred through the access patterns to individual pages, which we express through system call sequences. First, we explain how we extract information about system calls issued by the guest. Next, we describe how the sequence of system calls issued by the guest is combined with the sequence of writes to guest memory to identify the location of selected data structures as well as their state.

***Trapping System Calls into the Hypervisor***   To record a sequence of system calls issued by the guest we need a mechanism to trap into the hypervisor when a guest user process tries to execute a system call. It is important to note that we can extract system call information without access to the guest register or control state. Instead, we remove the execute permissions on the guest memory page containing the system call entry function `entry_SYSCALL_64` as well as the pages containing the system call handler routines. Thereby we enforce an exit to the hypervisor whenever system call execution is initiated. By examining the fault address, the hypervisor can determine which handler caused the fault. Initially, only the system call entry page is protected; if a `vmexit` is induced by guest execution of this page, we restrict access to the handler pages and restore execute permissions on the entry page. Similarly, if a `vmexit` is induced by guest execution of one of the handler pages we restore execute permissions to all handler pages while restricting access



```
cr3:1 ... - bind - listen            accept - recv - ...
cr3:2                  open - read
| | | listen | listen | | | |
| | | |      |        | | | RAM

cr3:1 ... - bind - listen            accept - recv - ...
cr3:2              open - read
| | | listen | listen | | | |
| | | open  |        | | | RAM
                     ...
cr3:1 ... - bind - listen            accept - recv - ...
cr3:2                  open - read
| read | | listen | listen | | accept |
|      | | open   |        | | recv   |
|      | | read   |        | |        | RAM
```
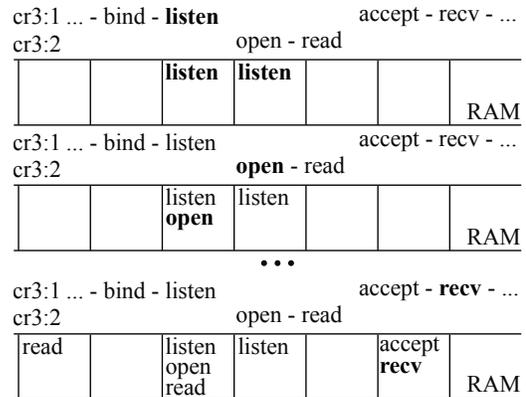
Figure 4: System Call Sequences on Guest Memory

to the entry page. This procedure is necessary, to enable the re-execution of the faulting instruction in the guest.

To remove execute permission from guest pages containing the respective functions, we first need to locate them in guest physical memory. Due to *KASLR*, the physical load offset of the kernel text section is randomly initialized during the VM boot process. We therefore employ a similar method as described in Section 5.1.1 to adjust the guest physical addresses of the system call entry and handler functions accordingly. To obtain a point of reference from which to compute the physical load offset of the kernel text section, the hypervisor can trigger the immediate execution of known functions, like interrupt handler routines, by the guest. By previously marking *all* guest memory pages as non executable through the nested pagetables, the guest will immediately fault, revealing the physical address of the triggered function through the fault metadata provided to the hypervisor. The random physical load offset of the kernel text section is then calculated by subtracting the fault address from the physical address of the function, obtained from a non randomized kernel image.

***Combining System Call and Write Sequences***   Based on the recorded system call sequence, the hypervisor can reason about the state of guest execution. However, we still lack the ability to identify which of the many guest memory pages that are written continuously, contains the data selected for replay. To that end, we cross-reference the sequence of guest memory writes with the system call information by storing a sequence of system calls for each page that preceded a write access to the respective page.

We now explain our approach via a simplified example. Figure 4 shows an excerpt of guest execution with two concurrent processes. The processes are identified by their according *cr3* values (either 1 or 2). Each process performs a number of different system calls, whereas the most recent one is highlighted in bold font. Guest pages are subsequently marked with the system call identifier that was last recorded before a write access occurred. To record those we intercept
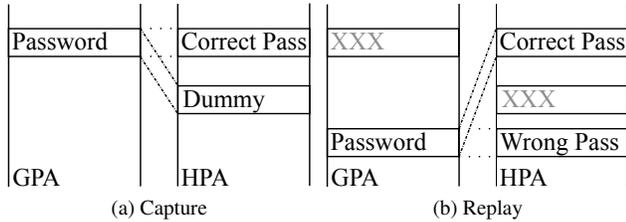
Figure 5: Nested Pagetable Modification



Figure 6: Replay Attack Overview

the guest on memory write access in addition to system call execution. Upon a vmexit induced by execution of a system call handler, we now also remove write permission from all guest pages. Each subsequent write will now trap to the hypervisor, where we firstly restore write permissions for the respective page to allow for the re-execution of the faulting instruction and secondly, mark the accessed page with the last recorded system call identifier. On each memory write we then evaluate these sequences for all guest pages to infer whether a specific page currently contains the data selected for replay.

### 5.2.3 Replay Attacks against Encrypted VMs

In this section we first describe the four phases composing our replay attack, namely offline analysis, tracing, capture, and replay. We then illustrate these steps by describing our procedure to replay OpenSSH login credentials to gain access to an encrypted VM at the privilege level of the connecting user.

***Offline Analysis*** The first stage is an offline analysis of the target application to determine possible replay attack vectors. Currently, we do this manually and on a per-application basis.

***Tracing*** To determine the location of the credential data structure in encrypted guest memory, we first trace system call and memory access patterns of an unencrypted guest running an identical OS and target application. This allows us to scan the unencrypted guest memory for the selected data continuously. If the selected data is detected in a guest memory page, we store the respective access pattern. Due to interrupts, scheduling and input from external sources, execution paths and therefore system call sequences might differ slightly. We account for this by collecting multiple traces and simply extracting the longest trailing sequence occurring in most of the traces. The sequence starts at the syscall after which the memory is replayed. From there we trace backward along all collected syscall sequences until one of them diverges while discarding sequences whose length falls under a predefined threshold.

***Capture and Replay*** In the capture and replay phases, we compare the collected sequence against those generated by the encrypted guest. If the system call sequence of a guest
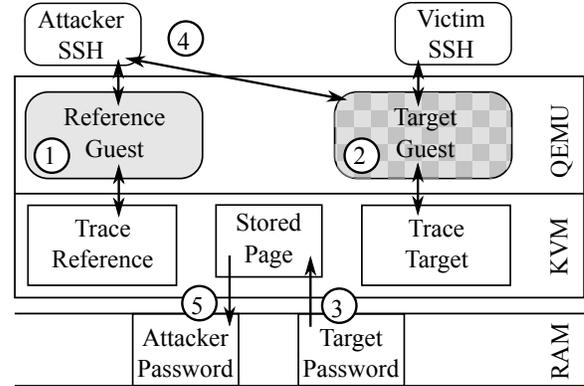
page matches the reference we conclude that the encrypted page contains the selected data and proceed to capture or replay the contained data respectively.

In accordance with the encryption scheme described in Section 4, the cipher text produced by the memory encryption algorithm is influenced not only by the content but also by the host physical address of the memory page. Therefore the replayed data has to be placed at the same host physical address (HPA) as the captured data. This can be achieved by manipulating the nested pagetables, which control the mapping of the guest physical address (GPA) to its HPA as shown in Figure 5. In Figure 5a the guest memory page, containing the valid credentials *"Correct Pass"* has been identified for capture. The nested pagetable entry connecting the GPA to the HPA of this page is then modified to redirect the GPA to a newly allocated page (*Dummy*). This removes the captured data from the guest's address space so that it will not be overwritten. During the replay phase described in Figure 5b the nested pagetable is again modified to redirect the GPA from the HPA containing the *"Wrong Pass"* to the HPA of the previously captured page still containing the *"Correct Pass"*. Using this method the minimum size of data that can be replayed is a single page (usually 4KB), because the address translation can be changed only on page granularity.

Figure 6 gives an overview of our replay attack. Here we assume that offline analysis has already identified data and state for capture and replay. We collect a reference sequence of system calls (*Trace Reference*) for the page containing the identified data, by initiating SSH password logins ① to an OpenSSH server running in an unencrypted guest (*Reference Guest*) with the same software configuration as the target. Next, we wait for an incoming SSH client connection to the protected guest (*Target Guest*) ②, while continuously comparing the access patterns of the protected guest (*Target Trace*) against the reference. If an SSH client, authenticates itself to the server via password, the page containing the credential data structure (*"Target Password"*) is identified and the content of the page is stored ③. We then re-initiate a password login to the protected guest from the attacker

controlled SSH client ④. To grant access to the attacker controlled client, the hypervisor modifies the nested paging structures to redirect the page containing the invalid credentials of the attacker (*"Attacker Password"*) to the stored page ⑤.

### 5.2.4 Impact

To show the effectiveness of the replay attack, we evaluate it by exploiting OpenSSH version 6.7p1-5+de running in a VM. The test was conducted on a AMD Phenom II X4 965 processor with 4GB RAM. As the host operating system, we used Linux with kernel version 4.4.0 and QEMU version 2.7.50 for communication with the KVM driver module. For the evaluation, we disabled symmetric multiprocessing on the host system. The guest was configured with 512MB RAM and ran kernel version 4.9.0-rc5 with the full range of *KASLR* options enabled. As AMD SEV is not available at the time of this writing, we substitute an unencrypted VM as our target. We argue that the results apply to a future real SEV setup, because none of the data structures required for the attack will be obscured even if SEV is enabled, according to the currently available documentation.

The effectiveness of our attack is best classified by the number of successful logins to the target guest that have to be observed on average, before successful execution of the described replay attack. The success rate hinges on the accuracy of page identification via system call and memory access patterns as well as on the structure of the page selected during offline analysis.

***Page Structure*** We found that the offset of the credential data within a memory page varies between four separate values. The distribution of the offset values is shown in Figure 7. To determine this distribution we initiated 387 SSH password logins using a unique password to simplify the identification of the page. By examining the collected traces, we determined that the specific location cannot be extracted from the system call access sequence. Further, we discovered that replaying captured data over a page with mismatching offset will terminate the guest process handling the login. However, termination of the process spawned by the SSH server to handle the connection will not impact the functionality of the guest, since it is immediately respawned by the parent process. Unsuccessful replay attempts will require the re-initiation of a new capture because the captured page was mapped back into the guest's memory space. Overwriting or removing a mapped guest physical page will result in unpredictable behaviour, unless the content of the page is known.

***Trace Accuracy*** To improve the coverage of guest execution paths and thereby the trace accuracy we collect multiple system call sequences. From those, we extract the sequence of system call identifiers that identifies the greatest number of the collected traces correctly. To measure the trace accuracy, we collected 387 traces to compute the reference sequence. We then proceeded to initiate 2155 SSH connections
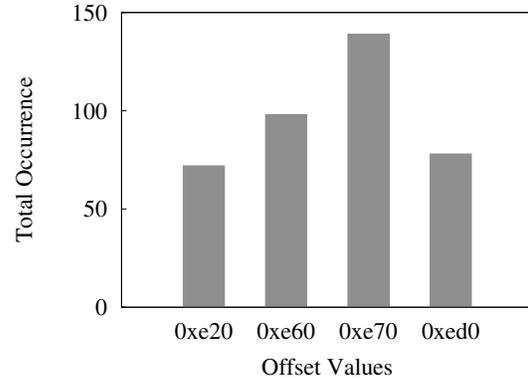


Figure 7: Credential Offsets

to the VM and matched the generated traces against the reference. To verify the correct identification of the page, we choose a unique string as the password and tested whether the guest page identified by the reference sequence contained this string and whether the data structure of the page matched the data structure of page selected for replay. The achieved identification accuracy for the page containing the credential data was **86%**. We encountered no false positives during our test, which is important because the remapping of a falsely identified guest page during capture will most likely have an adverse effect on the guest.

***Success Rate*** To measure the overall success rate, we compare the number of observed valid logins against the number of times, the attacker was granted access to the system. To that end, we run two identical VMs (reference and target) and extract a reference sequence from 387 observed logins to the reference VM. While using this reference to identify the password data in the target VM's memory, we initiated 2155 SSH password logins with valid login credentials to the target VM. This resulted in 505 successful replays; therefore the success rate of the attack is **23%**. The result is consistent with the measurements of data location distribution and trace accuracy. Especially the variation of credential offsets within a memory page limits the possible success rate to maximally 25%. However we argue that this factor can be alleviated by a more thorough investigation of the target software stack, to identify data structures suited for replay, with less varying location offsets.

## 6. Discussion

In the previous sections we laid out the details how a malicious hypervisor can exploit design issues in AMD's upcoming Secure Encrypted Virtualization technology to a) to gain access to a guest, b) to read encrypted guest memory and c) to fully revert any memory encryption configured by the tenant.

In this chapter, we discuss possible mitigations to these threats and evaluate their projected impact on performance and usability.

## 6.1 Mitigations

Pure software changes can not eliminate the design issues discussed in Section 3. To thwart the attacks presented here we propose the following design changes for future versions of SEV:

- Encrypted general purpose registers
- No access to the vmcb after an initial configuration
- Memory protection against hypervisor access

While the presented mitigations against the first two attacks induce a bearable degradation of guest transparency and performance, the replay attack remains difficult to mitigate given these demands.

***Access to general purpose registers*** The general purpose registers must never be visible to the hypervisor as they leak sensitive guest data on any vmexit. A guest does not have control over exits to the hypervisor, thus the encryption of general purpose registers must be enforced by the hardware. The encryption imposes another difficulty as certain guest operations require the hypervisor to read the general purpose registers. For example, when the guest writes data to a virtual device, this memory access will trap into the hypervisor. If the instruction causing the trap takes the value to write from a register, the hypervisor which is attempting to emulate the access, will not be able to read it when the general purpose registers are encrypted. The vmcb must be extended to contain decode assists for these events to provide the required information. As indicated in [21] decode assists are already in place to allow the hypervisor to read the instruction causing a vmexit. For future versions of SEV, these assists must be extended to contain the register values that hold the arguments to the instruction. The system must only augment vmexit events caused by traps to shared pages with these decode assists, in order to ensure that a malicious hypervisor cannot force a guest to reveal register content through decode assists.

***Access to the vmcb*** Usually, the vmcb is configured only once during the initial setup while at runtime a benign hypervisor does not need to modify the vmcb, with some exceptions. The fact that SEV allows us to alter the vmcb nevertheless imposes a security risk as it permits us to divert the control flow of guest by setting an arbitrary instruction pointer. We propose to change the existing state caching mechanism to enable the creation of a write-once vmcb. Currently, the content of the vmcb is already cached to improve context switch performance. The CPU is permitted to use the cached values of the vmcb unless the hypervisor explicitly clears bits in a special vmcb area called *vmcb clean field* and thereby forces the CPU to reread vmcb data. By prohibiting the hy-

pervisor from altering this field, the CPU is always able to use the cached values. At the first start of a guest the CPU copies the hypervisor provided vmcb into the cache. During runtime, the system always uses the cached vmcb. The initial vmcb can be assumed trustworthy because it is taken into account for the remote attestation. If the hypervisor wants to schedule another guest, hence another vmcb must be loaded, the system must provide a way to store the cached vmcb encrypted in shared memory.

However, there are elements in the vmcb which the hypervisor must be able to modify at runtime. Most importantly, injecting virtual interrupts into a VM requires the modification of several fields, among them V_IRQ, V_INTR_PRIO, V_INTR_MASKING and V_INTR_VECTOR. Efficient injection of multiple pending interrupts also requires access to the EVENTINJ field and the VINTR bit in the generic instruction intercept selection bitmask. These elements would either have to be excluded from our proposed "mandatory caching" scheme, or AMD's interrupt controller virtualization (AVIC) could be declared as a dependency of SEV thus making those vmcb elements obsolete.

***Access to guest memory*** Writes by the untrusted hypervisor to guest memory are dangerous. The fact that no memory authentication is in use opens the door for fault injection and replay attacks as presented in this paper. The most common way to protect memory from unauthorized access are integrity trees. However, they induce a notable performance and memory space overhead [23]. In a more relaxed attack model where physical attacks such as bus intercepts or direct memory accesses are not considered, it is sufficient to prevent the hypervisor from writing encrypted guest memory using mechanisms such as *CIP* as presented in [24]. The exclusion of pages from hypervisor access requires nontrivial changes to the guest operation system as well as the hypervisor. Further, the proposed access restrictions impact or even prohibit major cloud maintenance operations like snapshotting or live migration. Intel's SGX technology uses both encryption and integrity checks to protect the memory of enclaves [15]. However, SGX enclaves are small compared to VMs, and it is thus still an open question whether protecting the memory of complete VMs by integrity trees is feasible.

## 7. Related Work

In the following section, we present some topics which are relevant to this work.

### 7.1 Attacks

While attacks against AMD's SEV have not been published, several attacks against similar systems have been proposed. Checkoway et al. [6] proposed an attack method dubbed Iago whereby a malicious kernel manipulates system call return values to mount arbitrary code execution attacks on a system that protects userland applications from a malicious kernel. This work clearly shows that it is important to secure the

system call interface from an adversary. Linux system calls can be identified by a unique number that is stored in the general purpose registers. As these registers are still subject to manipulation by a hypervisor, this type of attack is also applicable to AMD SEV.

Xu et al. [28] showed how secret data can be extracted by inferring from page faults that specific execution paths inside a protected SGX enclave were executed. Using these execution traces, they were able to reconstruct images that were processed inside this enclave. Their approach of inferring memory content based on pagetable fault information is similar to the approach used in the proposed replay attack. SGX however does not hide the process internal address mapping from the attacker, which allows for a much more direct method of inference. Further, they did not deal with multiple concurrent processes.

Weichbrodt et al. proposed an attack dubbed AsyncShock [26]. They exploit the fact that the operating system is responsible for scheduling SGX enclave threads. By forcing enclave exits during the execution of multithreaded enclave code, they were able to mount use-after-free and TOCTTOU attacks on SGX protected enclaves.

Similar to our replay attack, Branco et al. [5] exploit the lack of memory authentication to compromise an encrypted system. The compromise is accomplished by injecting faults into critical state areas of the login process via the JTAG interface.

### 7.2 Defenses

Protecting applications from higher privileged software has been the subject of research for a long time. Many solutions that target single applications were proposed such as [7, 8, 14, 22]. Many of these solutions assume the existence of a trusted hypervisor to enforce protection of single applications or parts of an application.

A different direction is explored in the publications [17, 24, 24, 29]. The goal of their research is to provide protection mechanisms that ensure the integrity and confidentiality of the guest even in the case of a compromised hypervisor. Zhang et al. proposed CloudVisor [29] where a trusted security manager provides protection of guest VMs by means of nested virtualization. In contrast, Seongwook et al. proposed *H-SVM* [17], a purely hardware-based mechanism to protect guest systems. The guest memory is not mapped into the hypervisor context and a new hardware component, *H-SVM*, is controlling the nested pagetable. The hypervisor cannot access guest memory as it cannot create mappings itself, because the nested pagetables are protected. *H-SVM* protects the guest state by setting aside a dedicated memory area that is also not accessible by the hypervisor. If the hypervisor needs to access guest memory, the corresponding page must be explicitly marked by the guest. Physical attacks are not considered by *H-SVM*.

Similarly, Szefer et al. presented HyperWall [24]. Instead of removing the hypervisor's ability to manage the nested pagetable, an additional protection mechanism is introduced: *Confidentiality and Integrity Protection tables*, short *CIP*. These tables are consulted by the MMU when accessing memory.

Xia et al [27] followed this path with *HyperCoffer* and added protection against physical attacks by using encrypted memory with integrity checks. In this later publication they also address the lack of support for common cloud maintenance operations, like live migration or VM snapshotting and restoration.

## 8.  Future Work

After the initial publication of this paper, AMD released a new version of their Programmer's Manual [4]. The updated version details a new set of features called SEV-ES, which encrypts the guest state and enables the guest to finely control which state to share with the hypervisor. Therefore the first two attacks described in 1 are only effective against systems without SEV-ES support. Further research is needed to examine which attack vectors, besides replay, remain despite the proposed mechanisms.

While we clearly show, that SEV and SEV-ES cannot offer protection against a malicious hypervisor, we are confident that those technologies will thwart a substantial amount of attacks that rely on fewer capabilities. The question to what extend less severe hypervisor bugs, which do not lead to a complete compromise, but rather cause data leakage or allow to rewrite the memory of another VM impact the confidentiality and integrity of the tenant's data, has yet to be examined.

## 9.  Conclusion

This paper presents a first security evaluation of the upcoming Secure Encrypted Virtualization technology by AMD. While there are no actual CPUs available yet, the official documents published by AMD give away design issues that can be exploited by a malicious hypervisor.

By implementing three proof-of-concept attacks, we showed that these issues can be exploited to fully circumvent the protection mechanisms introduced by SEV. Furthermore, we showed that even when the hypervisor is not able to control the guest using the *vmcb* and general purpose registers, the control over the nested pagetable combined with the ability to inject interrupts is enough to mount a replay attack. We proposed possible hardware extensions to mitigate our attacks and compared similar solutions presented by the scientific community. Although we discovered serious design issues of AMD's SEV, we still think that the technology is promising considering the mitigations discussed in this paper.

## References

[1] AMD: Secure virtual machine architecture reference manual. Whitepaper, 2005.

[2] Secure Encrypted Virtualization Key Management. `http://support.amd.com/TechDocs/55766_SEV-KM%20API_Spec.pdf`, August 2016.

[3] AMD. Amd-v nested paging. *http://sites.amd.com/us/business/it-solutions/virtualization/Pages/amd-v.aspx*, 2008.

[4] A. AMD. Architecture programmers manual: Volume 2: System programming. *AMD Pub*, (24593), 2016.

[5] R. Branco and S. Gueron. Blinded random corruption attacks. In *2016 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*, pages 85–90. IEEE, 2016.

[6] S. Checkoway and H. Shacham. Iago attacks. *Proceedings of the 18th international conference on Architectural support for programming languages and operating systems (ASPLOS)*, 41(1):253, 2013.

[7] X. Chen, T. Garfinkel, E. C. Lewis, P. Subrahmanyam, C. A. Waldspurger, D. Boneh, J. Dwoskin, and D. R. Ports. Overshadow: a virtualization-based approach to retrofitting protection in commodity operating systems. In *ACM SIGARCH Computer Architecture News*, volume 36, pages 2–13. ACM, 2008.

[8] Y. Cheng, X. Ding, and R. Deng. Appshield: Protecting applications against untrusted operating system. *Singaport Management University Technical Report, SMU-SIS-13*, 101, 2013.

[9] V. Costan and S. Devadas. Intel sgx explained. Technical report, Cryptology ePrint Archive, Report 2016/086, 2016. https://eprint. iacr. org/2016/086, 2016.

[10] CVE Details: The ultimate security vulnerabilty datasource. Microsoft Hyper-V: CVE-2016-0088. `https://www.cvedetails.com/cve/CVE-2016-0088`, September 2016. Accessed: 2016-09-07.

[11] CVE Details: The ultimate security vulnerabilty datasource. VirtualBox CVE-2014-0983. `https://www.cvedetails.com/cve/CVE-2014-0983`, September 2016. Accessed: 2016-09-07.

[12] CVE Details: The ultimate security vulnerabilty datasource. VMWare: CVE-2015-2337. `https://www.cvedetails.com/cve/CVE-2015-2337`, September 2016. Accessed: 2016-09-07.

[13] CVE Details: The ultimate security vulnerabilty datasource. XEN: CVE-2015-5154. `https://www.cvedetails.com/cve/CVE-2015-5154/`, September 2016. Accessed: 2016-09-07.

[14] O. S. Hofmann, S. Kim, A. M. Dunn, M. Z. Lee, and E. Witchel. Inktag: Secure applications on an untrusted operating system. In *ACM SIGARCH Computer Architecture News*, volume 41, pages 265–278. ACM, 2013.

[15] Intel. Intel Software Guard Extensions (Intel SGX). `https://software.intel.com/en-us/sgx`, September 2016. Accessed: 2016-09-07.

[16] Jason Geffner, CrowdStrike. Qemu: VENOM vulnerability. `http://venom.crowdstrike.com/`, September 2016. Accessed: 2016-09-06.

[17] S. Jin, J. Ahn, S. Cha, and J. Huh. Architectural support for secure virtualization under a vulnerable hypervisor. In *Proceedings of the 44th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 272–283. ACM, 2011.

[18] Jonathan Salwan. ROPGadget Tool. `https://github.com/JonathanSalwan/ROPgadget`, September 2016. Accessed: 2016-09-07.

[19] D. Kaplan, J. Powell, and T. Woller. White Paper AMD Memory Encryption. `http://amd-dev.wpengine.netdna-cdn.com/wordpress/media/2013/12/AMD_Memory_Encryption_Whitepaper_v7-Public.pdf`, April 2016.

[20] A. Kivity, Y. Kamay, D. Laor, U. Lublin, and A. Liguori. kvm: the linux virtual machine monitor. In *Proceedings of the Linux symposium*, volume 1, pages 225–230.

[21] Linux kernel mailing list: Thomas Lendacky. Re: [RFC PATCH v1 00/18] x86: Secure Memory Encryption (AMD). `http://www.gossamer-threads.com/lists/linux/kernel/2435682#2435682`, May 2016. Accessed: 2016-09-11.

[22] J. M. McCune, B. J. Parno, A. Perrig, M. K. Reiter, and H. Isozaki. Flicker: An execution infrastructure for tcb minimization. In *ACM SIGOPS Operating Systems Review*, volume 42, pages 315–328. ACM, 2008.

[23] B. Rogers, S. Chhabra, M. Prvulovic, and Y. Solihin. Using address independent seed encryption and bonsai merkle trees to make secure processors os-and performance-friendly. In *Proceedings of the 40th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 183–196. IEEE Computer Society, 2007.

[24] J. Szefer and R. B. Lee. Architectural support for hypervisor-secure virtualization. In *ACM SIGPLAN Notices*, volume 47, pages 437–450. ACM, 2012.

[25] R. Uhlig, G. Neiger, D. Rodgers, A. L. Santoni, F. Martins, A. V. Anderson, S. M. Bennett, A. Kägi, F. H. Leung, and L. Smith. Intel virtualization technology. *Computer*, 38(5): 48–56, 2005.

[26] N. Weichbrodt and P. Pietzuch. AsyncShock: Exploiting Synchronisation Bugs in Intel SGX Enclaves. *ESORICS*, 2016.

[27] Y. Xia, Y. Liu, and H. Chen. Architecture support for guest-transparent vm protection from untrusted hypervisor and physical attacks. In *High Performance Computer Architecture (HPCA2013), 2013 IEEE 19th International Symposium on*, pages 246–257. IEEE, 2013.

[28] Y. Xu, W. Cui, and M. Peinado. Controlled-channel attacks: Deterministic side channels for untrusted operating systems. *Proceedings - IEEE Symposium on Security and Privacy*, 2015-July:640–656, 2015. ISSN 10816011. doi: 10.1109/SP.2015.45.

[29] F. Zhang, J. Chen, H. Chen, and B. Zang. Cloudvisor: retrofitting protection of virtual machines in multi-tenant cloud with nested virtualization. In *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles*, pages 203–216. ACM, 2011.