

Poacher turned gamekeeper: Lessons learned from eight years of breaking hypervisors

Rafal Wojtczuk rafal@bromium.com

27 Jul 2014

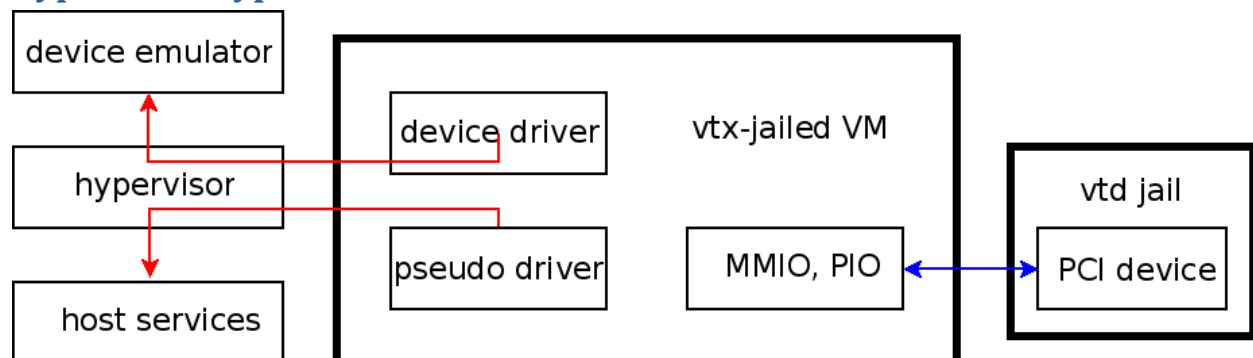
Summary

Hypervisors have become a key element of both cloud and client computing. It is without doubt that hypervisors are going to be commonplace in future devices, and play an important role in the security industry. In this paper, we discuss in detail the various lessons learnt whilst building and breaking various common hypervisors. In particular, we take a trip down memory lane and examine a few vulnerabilities found in popular hypervisors that have led to break-outs, trying to offer a generic mitigation when possible. To add some spice, we will talk about details of four not-yet-discussed vulnerabilities we recently discovered in VirtualBox, and examine DMA attacks against DeepSafe.

Scope

There is a plethora of various hypervisor solutions available nowadays. Some of them are designed from the scratch with security in mind, and e.g. use formal verification to provide assurance about their security. As those solutions are not the mainstream today, in this paper we will focus on the popular commercial virtualization software used commonly nowadays. One of the solution (DeepSafe) is very different from the others, and will be covered in the later part of the paper. We will start with discussion about common Type 1 and Type 2 hypervisors; Xen, VirtualBox, VMWare, ESX, HyperV all belong to these two categories.

Type 1 and Type 2 attack surface



What are the attack vectors that a malicious code running in VTx VM can use to break out of VM ? They can be grouped in the following categories:

- 1) Vulnerabilities in VTx itself. Certainly, if CPU does not isolate properly the code running in non-root mode, then all bets are off. There are certain CPU erratas related to VTx support (eg AAJ66, see [1]). Most of them require “complex micro-architectural conditions”, that are not disclosed in the erratas. We are not aware of any successful exploit based of VTx CPU errata. Also note that CPU erratas have the potential to circumvent any security software, not only hypervisors.
- 2) Vulnerabilities in the core hypervisor (code near vmexit handler). Such vulnerabilities are definitely thinkable; however, the size of the hypervisor core is relatively small, therefore the potential for vulnerabilities is limited.
- 3) Device emulators. There are many cases of such vulnerabilities. If a device emulator runs in a privileged environment (e.g. Type 2 host), such a vulnerability is fatal.
- 4) Other hypervisor-related services. One common example is shared folders service. By their very nature, they usually run in a privileged environment and again compromise of them is fatal.
- 5) Pass-through PCI devices. VM that can command a PCI device can mount certain additional attacks. Notably, hypervisor must protect its memory via VTd from DMA attacks; other types of attacks are also known ([2]). Note that it is possible to create a usable virtualization system without the need to grant any PCI devices to VMs. Also, in certain scenarios, PCI passthrough actually decreases attack surface – see below discussion on service VMs.

Some virtualization solutions are focused on providing as much features and functionality as possible, and inherently that increases the attack surface. Note that in case when users do not run potentially malicious code in the VMs, it is the natural approach. On the other hand, if the goal of the hypervisor is to provide a reliable isolation of malicious code, then care must be taken to design and implement the hypervisor in a way that minimizes exposure.

It is obvious that hypervisors present non-negligible attack surface, and there are numerous examples of vulnerabilities allowing the malicious code to break out of VM. Should we stop using hypervisors as a method to isolate malware? The answer is no – because they seem to be the best available mechanism to implement isolation. The other method is OS-based sandboxing, with Chrome sandbox being a primary example. However, this method relies on the security of an underlying OS – a vulnerability in OS kernel can be used to break out of OS-based sandbox [3]. At least in case of Windows, its kernel is a vast attack surface – 400 syscalls, 800 win32k.sys syscalls, drivers ioctls/WDDM escapes, resulting in 76 CVEs for kernelmode issues in Windows in year 2013 alone.

It is very difficult to come up with convincing quantitative comparison of a given hypervisor attack surface vs a given OS kernel (let's focus on Windows). The common way is to compare the size of the code base, because it is easy to come up with hard numbers. For example, xen-4.4.0 is ca 1.7 million lines of code. It can be trimmed to 110K lines of usermode code and 60K lines of ring0 code, still retaining rich hypervisor functionality. Windows7 kernel is believed to have ca 2 million lines of code, win32k.sys is probably larger. It shows that a hypervisor can be much smaller. However, it is extremely

difficult to determine what percentage of the code base actually is responsible for processing untrusted and potentially malicious input, which makes all LOC comparison essentially a futile game.

The attack surface should be measured not by total LOC of the whole product, but by the amount and complexity of the untrusted input that must be processed. Unfortunately, it is difficult to measure it in a way that makes it possible to compare attack surface between different solutions. We need to rely on subjective views, based by experience. Most knowledgeable parties agree that a well-written hypervisor has much smaller attack surface than a Windows kernel, the difference being even more significant than what could be concluded by LOC comparison.

Putting subjective views aside, there are certain properties of a hypervisor that are objectively attractive from the security perspective. Particularly, the vmexit boundary is much stronger than syscall and process boundary, which poses a challenge for exploitation of hypervisor memory corruption vulnerabilities. A few facts:

- 1) In case of broker-vulnerability-based sandbox escapes, on Windows attacker knows libraries bases – no ASLR protection
- 2) In case of kernel exploits, attacker can craft useful data structures in usermode that can be misinterpreted by the kernel, because the address space is the same (unless SMAP – but no SMAP for Windows anytime soon)
- 3) Windows kernel hands out its memory layout for free to attacker (better on Windows8.1) [4]
- 4) In case of browser vulnerabilities, attacker has a lot of control over memory layout, thanks to javascript/other scripting, that allows for very high degree of precision when mounting an attack, which is crucial for the reliable exploitation

No such problems exist in case of hypervisor-enforced separation. Usually, just because of the need to bypass ALSR, attacker needs two separate vulnerabilities, a pointer leak and write primitive (while in e.g. browser case, a single use-after-free vulnerability usually provides both). Most VM-escape exploits known to the author relied on ASLR being not functional (executable not position-independent, or libraries with fixed base). One notable exception is Cloudburst [5] exploit against VMWare, because the underlying integer overflow allowed for full read-write access of the host process memory – again, this seems to be a rare case.

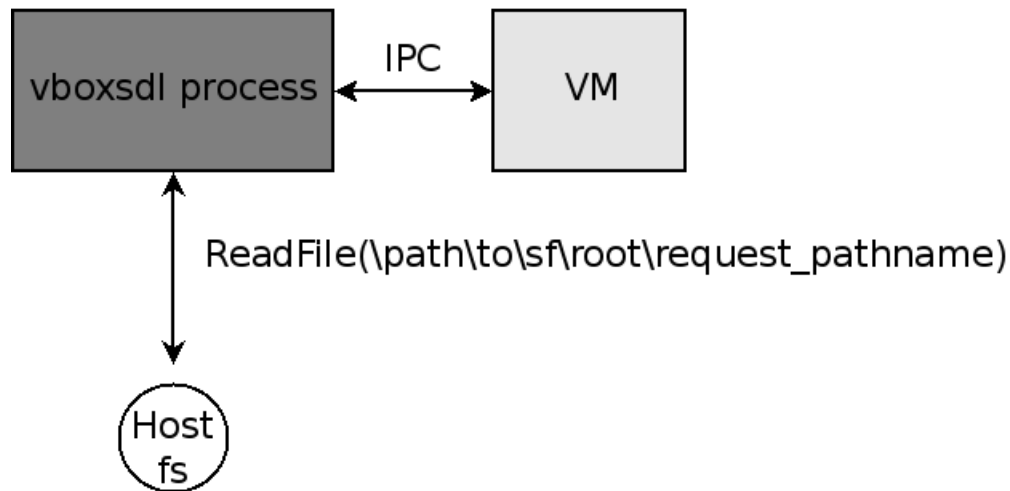
An important point is that often we do not have to choose between OS-based and hypervisor-based isolation – we can combine both. If the following, reasonable conditions are met:

- 1) hypervisor can be attacked only after compromising the VM kernel (note some hypervisors, e.g. VMWare, expose host services to unprivileged VM usermode)
- 2) hypervisor-related drivers in VM do not weaken VM kernel security significantly
- 3) there is nothing valuable in VM

then placing a OS-based sandbox within a VM is a pure gain – attacker would need to pierce both protections.

Case studies

Let's have a look at a few new and few old vulnerabilities in hypervisors, trying to pinpoint the design issues, and offer a generic mitigation if possible. We will start with four new VirtualBox vulnerabilities, reported by the author to the vendor in March 2014 and fixed in July 15 Critical Patch Update [6]. Three of them are related to shared folder functionality:



Shared folders need to be enabled in the VM configuration; the path of the root of the shared folder tree (`\path\to\sf\root` in the diagram above) needs to be specified there. VM sends a message containing operation type (read, write, and others) and a pathname (`request_pathname` in the diagram above) to the device model process running on the host (`vboxsd.exe`) over a communication channel (implemented over HGCM; the nature of the channel is irrelevant for our purposes). `vboxsd` performs the requested operation and returns the result to the VM.

VirtualBox shared folders host code is relatively complex:

- 1) Supports utf8 and unicode pathnames
- 2) Supports pathname casing correction (needed when VM has case-insensitive pathnames, and host does not)
- 3) Guest can specify path delimiter; host is supposed to normalize path changing each occurrence to `\`

Likely because of 1), there is no early null-termination check for the received pathname.

VirtualBox issue S0434934

Memory corruption in `vbsfbuildfullpath()`

```
397     /* Correct path delimiters */
398     if (pClient->PathDelimiter != RTPATH_DELIMITER)
```

```

399     {
400         LogFlow(("Correct path delimiter in %ls\n", src));
401         while (*src) // src comes from VM, not null-terminated
402         {
403             if (*src == pClient->PathDelimiter)
404                 *src = RTPATH_DELIMITER;
405             src++;
406         }

```

As we can see, if VM sends a string that is not null-terminated, vbsbuildfullpath() overwrites each occurrence of (attacker-chosen) pClient->PathDelimiter character to RTPATH_DELIMITER (that is defined as '\\'), until a null-terminator is found. If VM sends a string that is not null-terminated, this will corrupt the adjacent memory. The exploitation for code execution is difficult due to limited control the attacker has over the way the memory is corrupted, but not ruled out – particularly, overwriting a location that holds the size of a buffer might pave the way for a subsequent buffer overflow.

This is a textbook example of bad design. As both sides of the channel are VirtualBox code, it should be the guest side that does all the conversions. The host side should accept only a narrow range of inputs, say, only Unicode null-terminated ones – checking for input conformance requires much less code than actually doing the conversions.

VirtualBox issue S0434968

Shared folders directory traversal

Obviously, just concatenating the pathname received from VM to shared folder root leads to directory traversal via „..\..\..\..\request_pathname”, so the service needs to sanitize pathname to prevent this. VirtualBox path sanitize algorithm is:

- 1) Split the path into components (/ or \ is the path separator)
- 2) Start with depth_credit=0
- 3) For each component do: Switch (component)
- 4) Case “.” : do nothing
- 5) Case “..”: depth_credit-- //fail verification if depth_credits becomes negative
- 6) Default: depth_credit++;

So the goal is to check whether the number of “..” components in any path prefix is not greater than the number of “normal” path components. „dirname\.” is ok, „dirname\..\.” is not. If the above algorithm returns success, the request pathname is appended to the shared folder root and passed directly to host file API.

The above algorithm has a fatal flaw: it assumes that “\” is a path separator. It is true on a Windows host, but not on a Posix (say, Linux) host. As a result, assuming the VM is a Linux system with a writable shared folder mounted on /mnt/vboxsf, and the host runs Linux, the following sequence of shared folders operations will allow the VM to access arbitrary file on the host:

- 1) Mkdir /mnt/vboxsf/a\a\a\a\a\a\a
- 2) Access /mnt/vboxsf/a\a\a\a\a\a\a/../../../../../../../../arbitrary/file/on/the/host

Pathnames verification is a tricky subject, but there are well-known methods to deal with possible directory traversal attacks. The best way is to use OS-provided mechanisms – chroot() on Posix, \\?\ prefix on Windows host. If we resort to hand-crafted verification, it should be simple – we can require the guest to take care of “dirname/..” removal, and then host can deny any request with “..” path component in it.

VirtualBox issue S0434952

data leak in HGCM, exploitable via shared folders operations

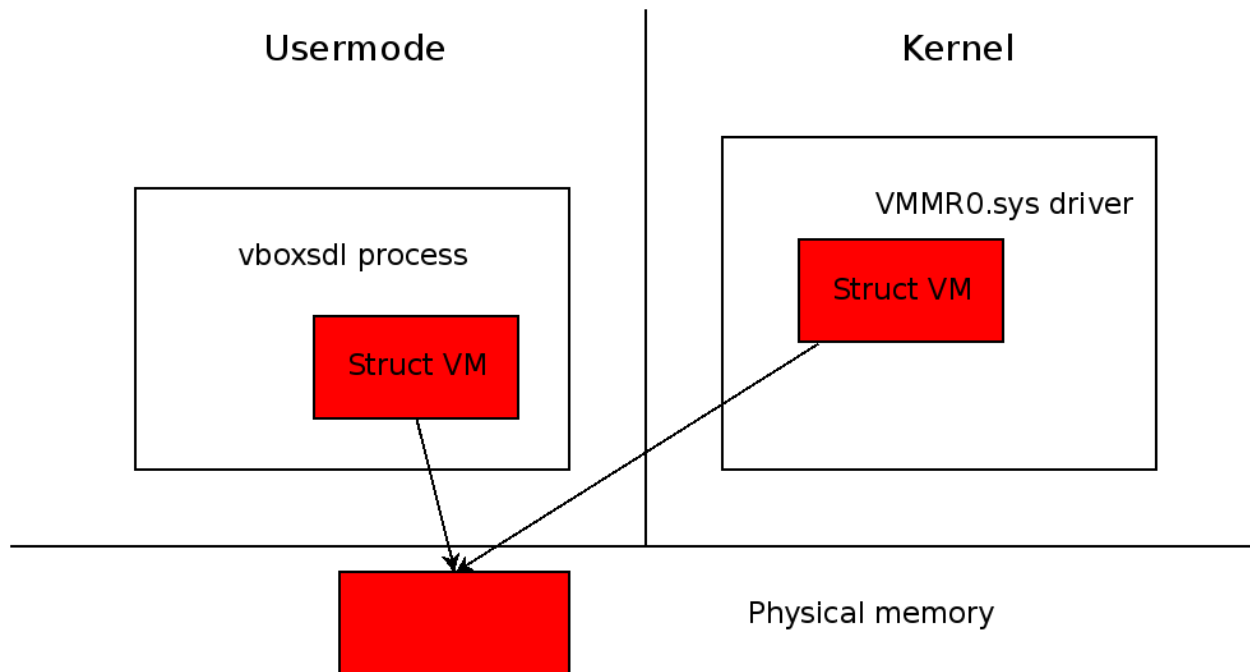
If VM sends SHFL_FN_READ message requesting Nreq bytes to be read, the response from the host contains: a response buffer B and the amount of bytes read Nresp. Due to particular way the underlying HGCM transport mechanism works (and some negligence), B is a character buffer always Nreq bytes long, even if Nresp is actually smaller. So, when VM requests to read 1024 bytes from a zero-length file, host returns 1024 bytes long uninitialized buffer (plus information that 0 bytes have been read), thus leaking contents of uninitialized malloced buffer from the host.

VirtualBox issue S0434947

Frontend to kernel escalation on the host

This vulnerability cannot be exploited by a rogue VM. Instead, it allows an unprivileged user on the Windows host to achieve kernelmode arbitrary code execution.

The problem is that the crucial data structure “struct VM” is mapped read-write both in the kernel component (VMMR0.sys driver) and in usermode frontend process (vboxsd.exe):



It was confirmed that “struct VM” contains kernelmode code pointers, which are frequently dereferenced by VMMR0 driver. A malicious usermode frontend process can overwrite these pointers and gain kernelmode execution.

The above scenario was tested on Windows host – unprivileged user can attach vboxsd.exe with a debugger and perform required malicious actions. On Linux host, the frontend process runs as root, so it cannot be attached with a debugger (and /dev/vboxdrv device default permissions do not allow an unprivileged process to create a VM).

The remaining vulnerabilities referenced below are not new.

CVE-2007-5497

Integer overflow in libext2fs

In order to construct a PV domain on Xen, one must specify the kernel image to use. Xen’s pygrub process runs in privileged domain dom0 and uses libext2fs to retrieve kernel image from VM’s filesystem. The latter is controlled by malicious code running in VM. So, the attacker that has administrative privileges in VM can corrupt VM filesystem in a way that trigger any vulnerability (e.g. CVE-2007-5497) in libraries that parse the filesystem structures, and reboot VM. In order to restart VM, pygrub will be invoked, and it will parse the malicious filesystem - this might result in code execution in dom0.

The problem is elegantly solved by not using `pygrub` and switching to `pvgrub`. The latter is a “stub”, static kernel that one passes to Xen as a VM kernel. `Pvgrub` does all the filesystem parsing, and chainboots the real kernel. Thus, the untrusted filesystem is never parsed in the privileged `dom0`. Lesson – again, offload to VM as much as possible.

CVE-2011-1751

Use-after-free in `qemu/KVM`.

This vulnerability and its exploitation was described in detail during BHUSA11. Briefly, a malicious administrative user in KVM VM can request PCI unplug action (via writing to emulated chipset registers) on a device that was not hotplugged in. This results in use-after-free condition in the device model process (running on the host) that leads to code execution on the host.

It is instructive to notice that a feature that was really needed by very few users could be leveraged to compromise all other users. By default, this feature should be disabled.

Generally, it would be good to reduce attack surface by not allowing VM to exercise code paths in the device model related to PCI config space access. The key observation is that this functionality is needed only for VM boot. Thus if we plan to use a given VM as a container for potentially malicious code, we can prepare it in steps:

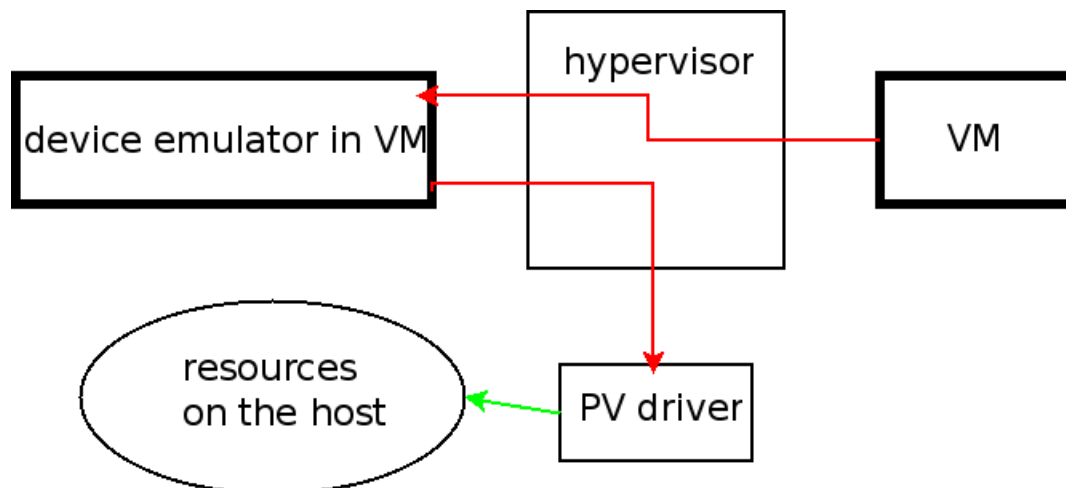
- 1) Start VM with all PCI config space access granted (and other resources needed for boot time only); do not expose VM to malicious input yet
- 2) Save VM
- 3) change the VM configuration so that it denies PCI config space access
- 4) restore VM with this new restrictive configuration
- 5) expose VM to the malicious input

Then even if VM gets compromised at step 5), it runs in a more restricted environment that might mitigate attacks. The term “delusional boot” [7] was coined to describe this technique.

CVE-2012-0029

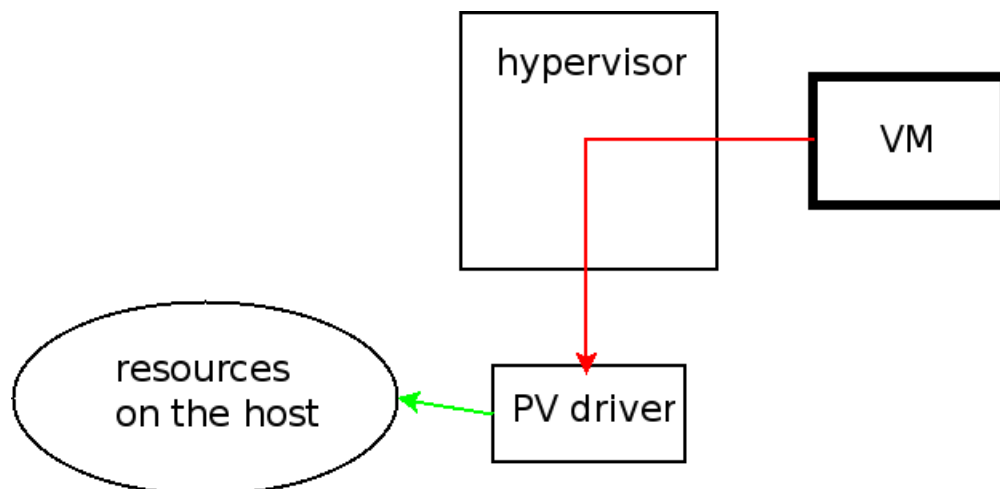
Heap-based buffer overflow in the `process_tx_desc` function in the `e1000` emulation

As mentioned initially, device emulation presents significant attack surface, and this CVE is a good example why. In case of Type 2 hypervisors, the relevant device model process usually runs on the host, and its compromise is fatal. It is thinkable to use OS facilities to restrict privileges of this process – however, if we trusted OS to do sufficiently good job in isolating, we would not need hypervisors for isolation. So a better idea is to sandbox device emulation in a dedicated “stub VM”:



The assumption is that the PV driver interface is much thinner than the interface exposed by stub VM. Therefore, even if a malicious VM is able to achieve code execution in the context of a device emulator, the damage is small – an attacker would need another exploit for the PV driver.

In fact, if we do not require an unmodified VM, we can install proper PV drivers in the VM:



Thus getting rid of a large part of device emulation entirely.

CVE-2007-0069

Windows Kernel TCP/IP/IGMPv3 and MLDv2 Vulnerability, remote code execution

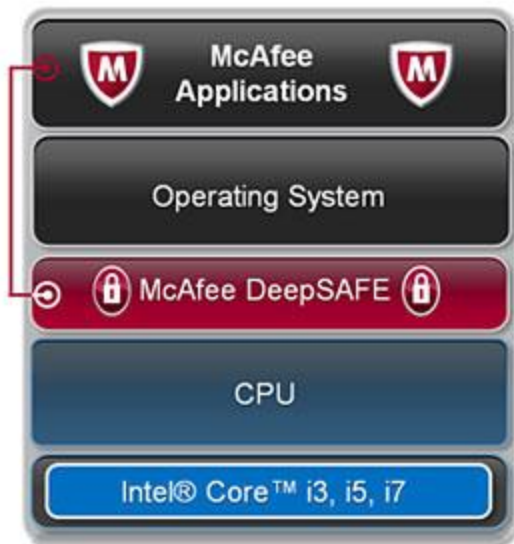
This vulnerability is not related to virtualization, but it is interesting to note how hypervisors can mitigate such problems. Some software components normally run in a privileged environment and are exposed to attacks from the outside (not from VM), networking stack and wireless card drivers being a

primary example. If we move this code to a dedicated “service VM” then the impact of an exploit against this component is greatly reduced. This requires granting the service VM full control over the relevant hardware.

In fact, it is thinkable to turn the whole type 2 host into a giant service VM [8]. If we implement the hypervisor core in a way that protects normal VMs from a compromised host, this provides significant advantage (particularly, solves all the problems with device emulation code). However, this approach needs careful implementation – particularly, hypervisor needs to be protected against hardware-based attacks originating in the host.

Deepsafe and DMA attacks

Deepsafe hypervisor [9] is very different from Type 1 and Type 2 ones:



Deepsafe hypervisor installs itself in a way very similar to Bluepill [10]. Early during Windows OS boot time, it wraps the whole OS in a VM, and resumes OS in VM. There is only one VM in the system – one holding the OS. The benefit is that even if malicious code finds a way into OS and even executes a kernel exploit (at this stage, all security methods implemented in OS are bypassed), then still Deepsafe stays intact and has a chance to detect malicious activity in the OS. Currently Deepsafe focuses on detecting typical rootkit activity, e.g. changes to SSDT, IDT, kernel body and syscall-related MSRs.

Deepsafe hypervisor needs to protect its body from a malicious OS. In order to prevent CPU (running in VM context) from accessing hypervisor memory, EPT mechanism is used, which is sound. Still, as there is no device virtualization, the OS needs full access to hardware. Particularly, it can instruct PCI devices to overwrite any memory via DMA. In order to prevent this attack, VT-d [1] mechanism must be used – by configuring VT-d, hypervisor can disallow certain memory ranges from being accessed via DMA.

It is somewhat unexpected, but Deepsafe (tested the latest available version 1.6.0) does not configure VT-d engine, and thus is vulnerable to DMA-based attacks. They are not something new – such attacks

were demonstrated in BHUSA2008 [12] by the author against Xen. There are also well known discussions [13] about the necessity of it.

One possible explanation for the lack of VT-d protection in DeepSAFE is that the vendor thought they are infeasible. Indeed, if we wanted to program given PCI device registers to initiate custom DMA directly, this would be a tedious device-specific task, and it would require a method to stop OS from interfering with this process, at the same time maintaining system stability. However, there is a simple way to achieve the same – we can use the existing disk drivers for it (a similar approach was used in [12]).

The procedure to conduct arbitrary DMA on Windows, in a stable manner, is:

- 1) Achieve kernel privileges
- 2) Allocate a page at virtual address V
- 3) Change PTE of V so that it points to physical address P
- 4) CreateFile(... FILE_FLAG_NO_BUFFERING ...)
- 5) ReadFile/WriteFile(..., V,...) will do DMA to P

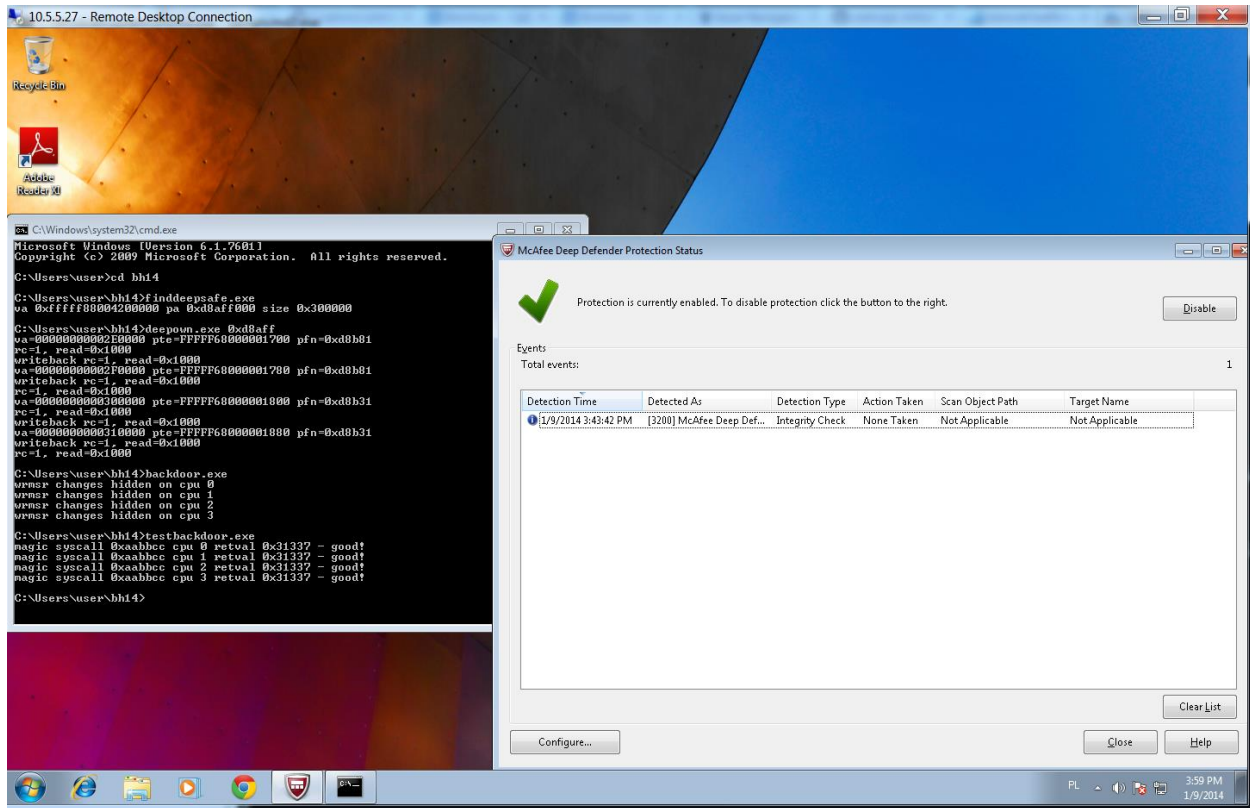
Step 4) instructs the OS to do all file operations on the given file directly via DMA to user buffer. There are restrictions – the transfer size must be aligned with sector size, and BitLocker not activated. Step 3) fools OS into thinking that user buffer is at [arbitrary] physical address P. OS will dutifully program the disk controller so that it will do DMA to the physical address of attacker's choice.

As a result, a compromised OS can overwrite DeepSAFE's body, and take full control over it. One possible way to use this capability is to disable the hypervisor entirely, by returning OS into root mode. While possible, it looks like a lot of work, and may interact with other DeepSAFE components running in the host.

From attacker's perspective, much better idea is to inject rootkit into the DeepSAFE hypervisor. It is the perfect location – secured from the rest of OS, and having much degree of control over it. For the proof-of-concept purposes, a very simple rootkit was implemented, that is installed in the code path handling vmexit for the reasons EXIT_REASON_MSR_READ and EXIT_REASON_MSR_WRITE. Normally, DeepSAFE configures VTx so that reads from LSTAR MSR do not cause vmexit, and writes to it cause vmexit and result in reporting an alert about possible rootkit-like behavior. We change hypervisor body so that

- 1) Writes to LSTAR MSR are allowed and dutifully executed, without any alert being generated
- 2) Reads from LSTAR MSR cause vmexit (change to MSR bitmap is needed for this) and are emulated in a way that return to the OS the original value installed by the OS

Because of the above, attacker in the OS is able to silently alter the LSTAR MSR as he wishes. It allows controlling all the syscall handling by OS, which results in efficient rootkit capabilities. Moreover, Patchguard running in OS is blinded – it still sees the original LSTAR MSR value when running its checks, because the hypervisor lies to it about this MSR value. To sum up, not only we have disabled the DeepSAFE detection, we also have used its power to disable in-OS rootkit checks.



There are a few interesting technical details regarding the above hypervisor overwrite. First, malware running in OS needs to know where in physical address space Deepsafe hypervisor is located. Dumping all the physical address space via DMA and doing pattern search in it is possible, but troublesome. A more elegant approach was found – it turns out that when EPT fault occurs because OS tried to read from a physical address belonging to the hypervisor, then Deepsafe does not bother to emulate the instruction, it just skips it. Thus, the following function

```
Mov rax, MAGICVALUE
```

```
Mov rax, [rcx]
```

```
Ret
```

Will return MAGICVALUE if memory at rcx belongs to Deepsafe, and something else (real memory content) if not. Deepsafe allocates a contiguous physical memory region of size 0x300000, so it is easy and fast to find it via scanning all the memory.

Another question is how to trojan Deepsafe effectively. The implemented proof-of-concept code is simple, it works with a given Deepsafe version only, and just places a hook at the first instruction of a vmexit handler; it stores the hook body in the unused end of the last page used by the text section. A more generic approach is possible, based on the fact that it is feasible to find VMCS hypervisor control

structures by pattern matching [14] in the memory area reserved by DeepSAFE. VMCS include the vmexit handler location and CR3 value used by hypervisor, so that we can change both reliably and without the need to hook the beginning of the original vmexit handler in a nondestructive manner.

More concerns about DeepSAFE

Currently DeepSAFE detects only a very specific set of rootkit-like functionality, mostly duplicating PatchGuard functionality (obviously, in a more reliable manner). However, there are legal interfaces exposed by Windows kernel, like filter drivers, which can be abused to gain rootkit-like functionality. As these interfaces are used by non-malicious software (say AV), the mere use of them is not a clear indication of an attack.

Other concerns are, just briefly mentioned for brevity:

- 1) Compromised host kernel can overwrite crucial usermode memory (e.g. disabling usermode DeepSAFE components)
- 2) How secure is mfeib.sys launch, on reboot/S3 resume?
- 3) No trusted UI domain
- 4) Host can mess with PCI config, SMM, BIOS, PCI devices firmware

Bibliography

[1] Intel® Core™ i7-900 Specification Update,

<http://download.intel.com/design/processor/specupdt/320836.pdf>

[2] Rafal Wojtczuk, Joanna Rutkowska,

"Following the White Rabbit: Software attacks against Intel(R) VT-d technology",

<http://invisiblethingslab.com/resources/2011/Software%20Attacks%20on%20Intel%20VT-d.pdf>

[3] "MWR Labs Pwn2Own 2013 Write-up - Kernel Exploit",

<https://labs.mwrinfosecurity.com/blog/2013/09/06/mwr-labs-pwn2own-2013-write-up---kernel-exploit/>

[4] Alex Ionescu, „KASLR Bypass Mitigations in Windows 8.1”, <http://www.alex-ionescu.com/?p=82>

[5] Kostya Kortchinsky, „CLOUDBURST: A VMware Guest to Host Escape Story”, BHUSA09

[6] Oracle Critical Patch Update Advisory - July 2014,

<http://www.oracle.com/technetwork/topics/security/cpujul2014-1972956.html>

- [7] Anh Nguyen, Himanshu Raj, Shravan Rayanchu, Stefan Saroiu, Alec Wolman
"Delusional Boot: Securing Cloud Hypervisors without Massive Re-engineering",
<http://research.microsoft.com/pubs/197590/inception.pdf>
- [8] Ian Pratt, „µXen”, http://www-archive.xenproject.org/xensummit/xs12na_talks/T6.html
- [9] McAfee DeepSAFE, www.mcafee.com/us/solutions/mcafee-deepsafe.aspx
- [10] Joanna Rutkowska, Bluepill, <http://blackhat.com/presentations/bh-usa-06/BH-US-06-Rutkowska.pdf>
- [11] Intel® Virtualization Technology for Directed I/O (VT-d),
<https://software.intel.com/en-us/articles/intel-virtualization-technology-for-directed-io-vt-d-enhancing-intel-platforms-for-efficient-virtualization-of-io-devices>
- [12] Rafal Wojtczuk, "Subverting the Xen hypervisor",
http://www.blackhat.com/presentations/bh-usa-08/Wojtczuk/BH_US_08_Wojtczuk_Subverting_the_Xen_Hypervisor.pdf
- [13] Joanna Rutkowska, "Thoughts on DeepSafe",
<http://theinvisiblethings.blogspot.co.uk/2012/01/thoughts-on-deepsafe.html>
- [14] Gianluca Guida, private conversations