



Compromise-as-a-Service

Our PleAZURE

Felix Wilhelm, Matthias Luft & Enno Rey
{fwilhelm, mluft, erey}@ernw.de





Agenda

- Introduction & Methodology
- Architecture & Attack Surface
- MS13-092

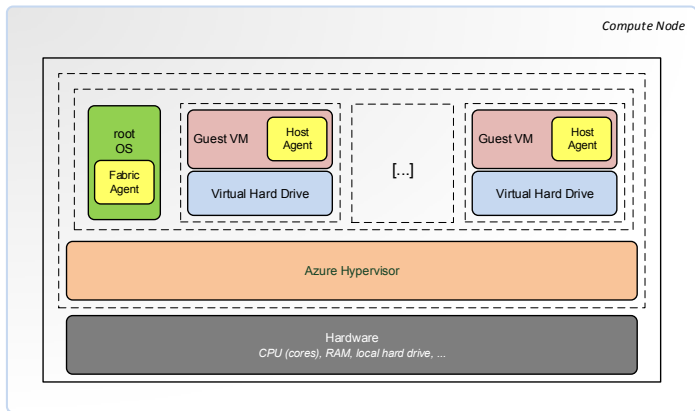


Prelude



- No complete security evaluation of Hyper-V
 - Still work in progress.
 - We will provide a detailed overview of the attack surface.
- Will talk about stuff that did not work (lots) and stuff that did work (some).
- Goal is to motivate and help other researchers.

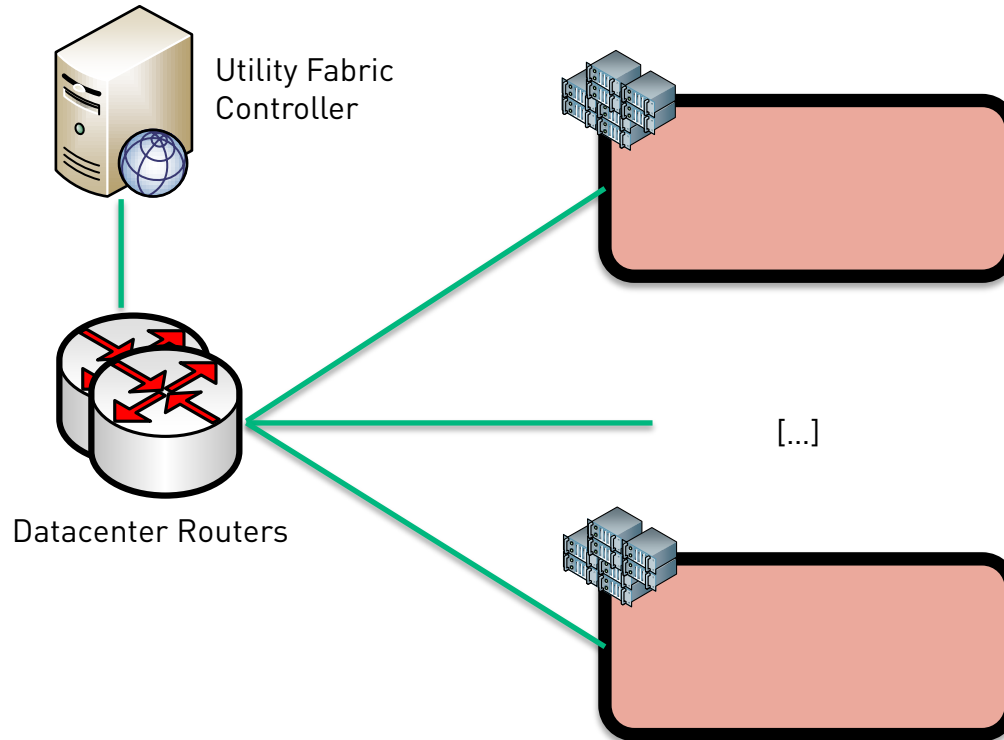
Background



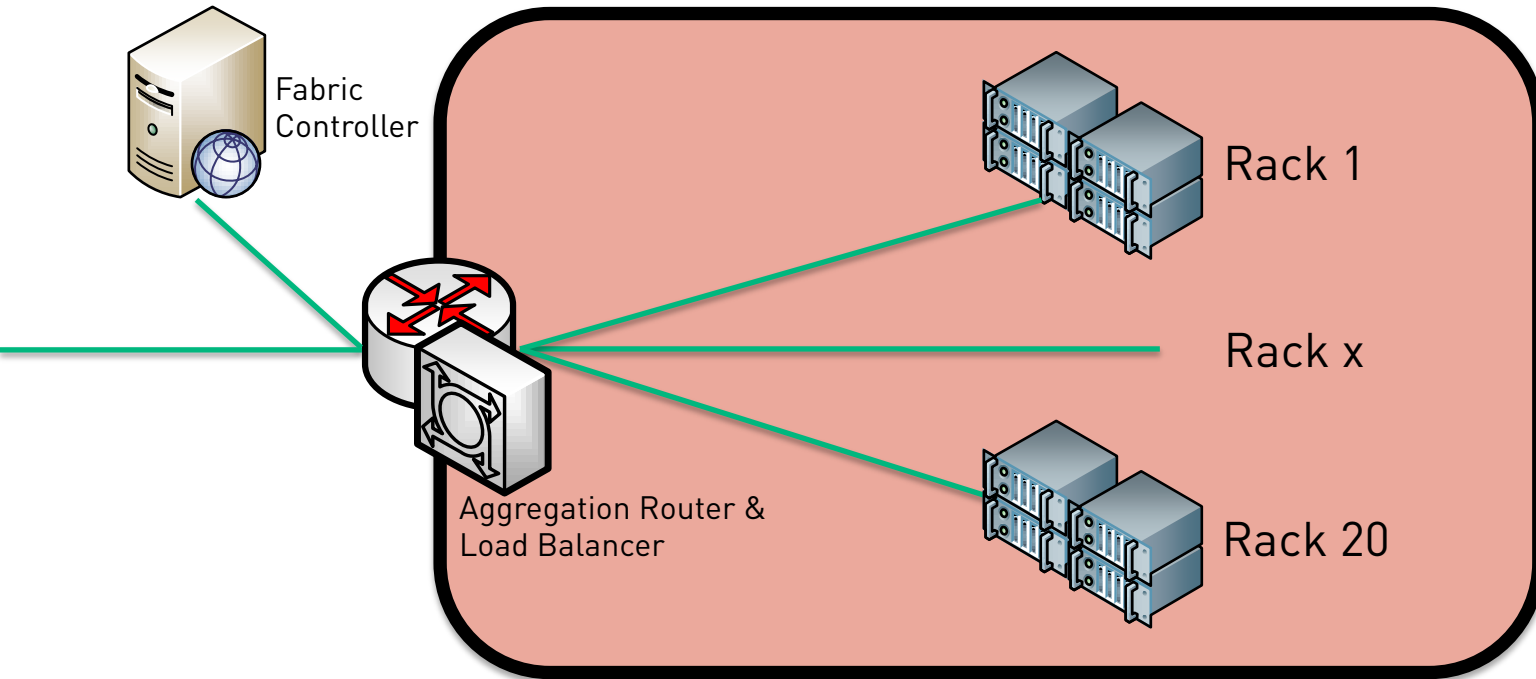
– Original research was part of German government research project.

- Overall security posture of the Azure Cloud
 - Network Security
 - VM Isolation
 - Management Interfaces
 -

Fabric

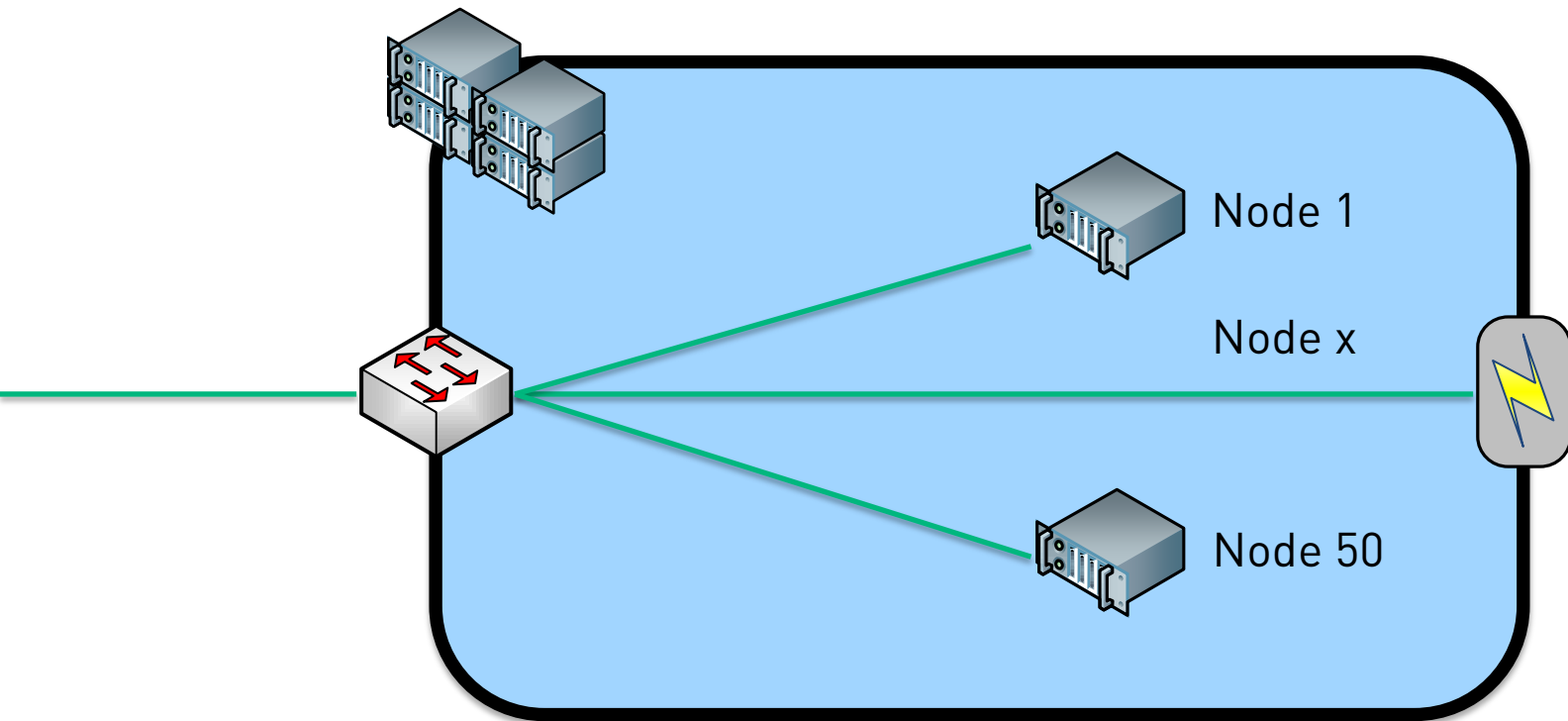


Cluster

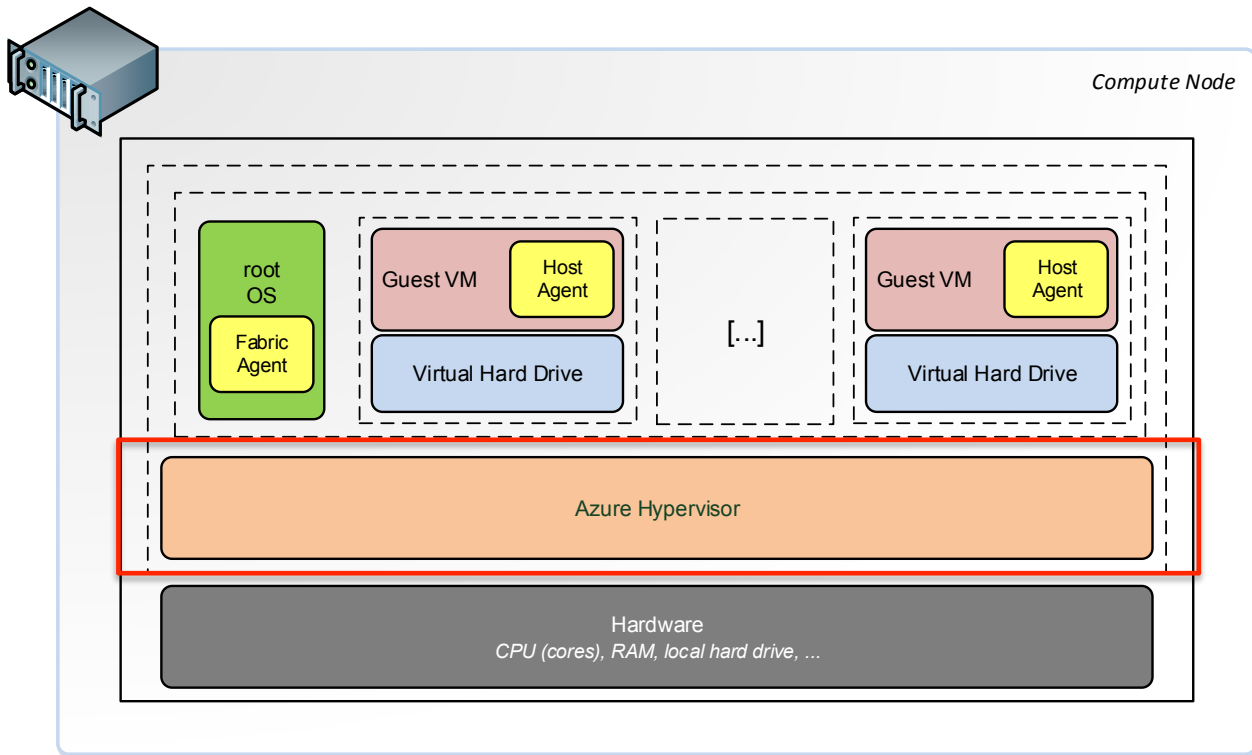




Rack



Compute Node



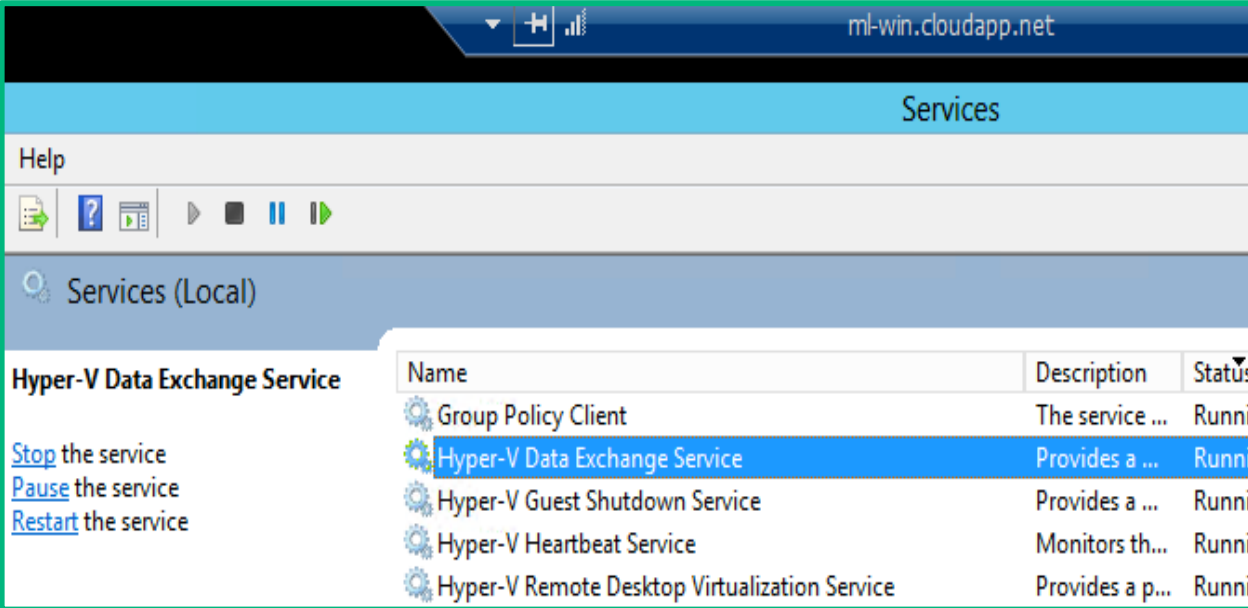


... the Cloud?

Why Hyper-V?



- Supposedly, very little research so far:
 - Four DoS vulnerabilities at all.
 - Almost no published 3rd party research.
- Used in a variety of corporate environments.
 - Including Azure!
 - Heavily marketed by Microsoft



The screenshot shows the Windows Services console. The title bar reads "ml-win.cloudapp.net". The main window title is "Services". Below the title bar is a "Help" section with a search icon and a toolbar with icons for help, refresh, stop, and start. The main area is titled "Services (Local)". A list of services is displayed, with "Hyper-V Data Exchange Service" selected and highlighted in blue. To the left of the list, there are links: "Stop the service", "Pause the service", and "Restart the service".

Name	Description	Status
Group Policy Client	The service ...	Runni
Hyper-V Data Exchange Service	Provides a ...	Runni
Hyper-V Guest Shutdown Service	Provides a ...	Runni
Hyper-V Heartbeat Service	Monitors th...	Runni
Hyper-V Remote Desktop Virtualization Service	Provides a p...	Runni

Azure Hypervisor

A.k.a. Hyper-V?



```
azureuser@ernw-s1:~$ ./cpuid
Checking for Hyper-V: True :)
Max leaf number: 0x40000006
Looks sane
```

```
Build number: 9200
Version: 6.2
Service Pack: 20
Service Branch: 20539
```

```
CreatePartitions: 0
AccessPartitionId: 0
AccessMemoryPool: 0
AdjustMessageBuffers: 0
PostMessages: 1
SignalEvents: 1
CreatePort: 0
ConnectPort: 1
AccessStats: 0
RsvdZ: 0
RsvdZ: 0
Debugging: 1
CpuPowerManagement: 0
```

Azure Hypervisor

```
root@virtual-linux:/home/ernw# ./cpuid
Checking for Hyper-V: True :)
Max leaf number: 0x40000006
Looks sane
```

```
Build number: 9200
Version: 6.2
Service Pack: 16
Service Branch: 16384
```

```
CreatePartitions: 0
AccessPartitionId: 0
AccessMemoryPool: 0
AdjustMessageBuffers: 0
PostMessages: 1
SignalEvents: 1
CreatePort: 0
ConnectPort: 1
AccessStats: 0
RsvdZ: 0
RsvdZ: 0
Debugging: 1
CpuPowerManagement: 0
```

Hyper-V

Versions



Methodology

Or: We're searching for runtime breakouts, right?



HV Security Objectives



- Main and crucial objective: Isolation!
 - On all layers (network, runtime, storage, ...)!
- Potential attacks violating this objective:
 - Breakout attacks (of course ;-)
 - Side channels & timing attacks
 - Traditional attacks: Account compromises, misconfiguration, ...



Parts of an Hypervisor

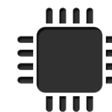
→ Device virtualization



→ Memory management & isolation



→ CPU virtualization



→ APIs



Vulnerability History

Hypervisor Breakout



- 2007: VMftp
- 2009: VMware Cloudburst
- 2011: Virtunoid – KVM Breakout
- 2012: VMSA-12-009
- 2012: VMDK Has Left The Building
- 2012: Xen SYSRET
- 2013: MS13-092
- 2014: VirtualBox HGCM
- 2014: VirtualBox Chromium Breakout



“Virtual Air Gap”



Hyper-V

~~VMware~~ *isn't an additional security layer:
It's just another layer to find bugs in"*

Kostya Kortchinsky/Immunity/Cloudburst, 2009



Hyper-V

Seriously, there is just no Hyper-V Logo available.

Hyper-V

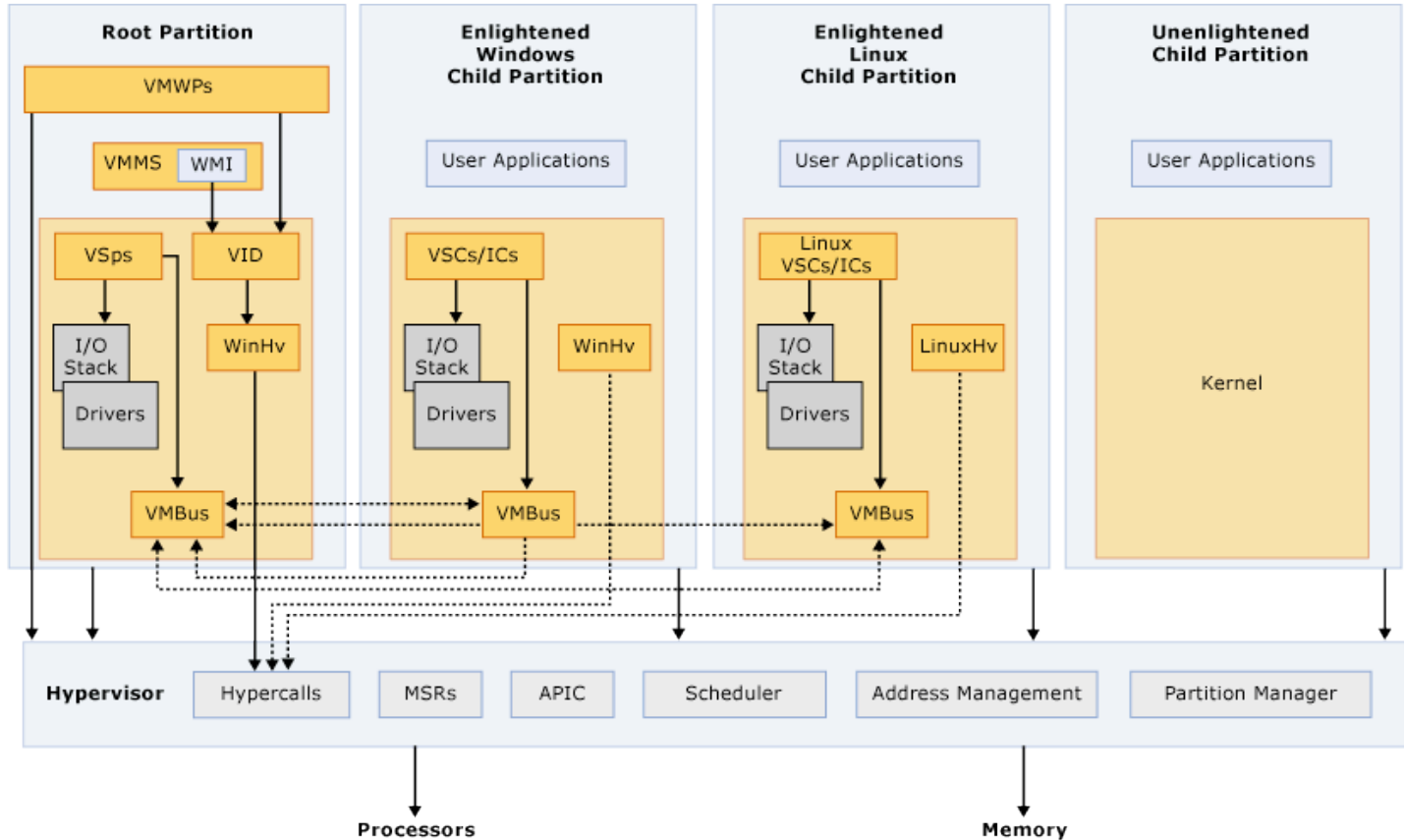


Windows Server®
Hyper-V™

- Type 1 hypervisor
 - „Bare metal“
- VMs are called „partitions“
 - Root Partition performs management duties.
 - Creation, Configuration, Destruction
 - Logging
- Support for „unenlightened“ + „enlightened“ systems
 - Enlightened = Explicit support for Hyper-V.
Requires guest modification.
- Uses Intel VT instruction set for hardware based virtualization



Hyper-V High Level Architecture



Source: [http://msdn.microsoft.com/de-de/library/cc768520\(en-us\).aspx](http://msdn.microsoft.com/de-de/library/cc768520(en-us).aspx)



Hyper-V vs. VMware ESXi

	Hyper-V	ESXi
Architecture	Micro Kernel	Monolithic
Most interesting attack surface	Parent Partition	Hypervisor
Platform to exploit	Windows	Proprietary Unix-like environment with a lot of custom libraries and characteristics

The micro kernel improves the security posture of the hypervisor, but moves attention to the parent partition – which simplifies post-exploitation!

Hyper-V Information Sources

19) **United States** 12) **Patent Application Publication** **Post et al.**

54) **SHARED MEMORY BETWEEN CHILD AND PARENT PARTITIONS**

75) Inventors: **Bradley Stephen Post**, San Francisco, CA (US); **Ed Cox**, Morgan Hill, CA (US)

73) Assignee: **Microsoft Corporation**, Redmond, WA (US)

21) Appl. No.: **12/894,896**

22) Filed: **Sep. 30, 2010**

- Hypervisor Top Level Functional Specification
 - Detailed documentation of Hypercall API
- Linux Integration Services
 - Open Source Hypercall / VMBus / VSC implementation
- Patent Applications
 - See picture
- Singularity Header Files
 - <https://singularity.svn.codeplex.com/svn/base/Windows/Inc/>
- Common Criteria Certification Documents
 - http://www.commoncriteriaportal.org/files/epfiles/0570b_pdf.pdf
- Binaries
 - ...

Reverse Engineering - Hypervisor



- Hypervisor itself is implemented in Hvix64.exe (Intel) and Hvax64.exe (AMD)
- Unfortunately no public debugging symbols available ☹️
- However some symbols can be ported from winload.exe and hvloader.exe
 - Networking Code, USB Stack, Debugger implementation
 - Full credit to Gerhart from securitylab.ru for this idea!



Reverse Engineering – VMX



Intel® 64 and IA-32 Architectures
Software Developer's Manual

Volume 3 (3A, 3B & 3C):
System Programming Guide

```
FFFFF8000064EA8D mov    eax, VMCS_Guest_SS_access_rights
FFFFF8000064EA92 vmread rax, rax
FFFFF8000064EA95 mov    [rsp+78h+arg_8], rax
FFFFF8000064EA9D shr    al, 5
FFFFF8000064EAA0 and    al, 3
FFFFF8000064EAA2 cmp    al, 3
FFFFF8000064EAA4 jz     loc_FFFF8000064EB72
```

- No documented system libraries
- Almost no strings
- Intel VMX instructions and VMCS are biggest help
 - Semi-automated approach: Locate vmx instructions and access to VMCS properties
 - Translate VMCS property values into corresponding name
 - Implemented using IDAPython script

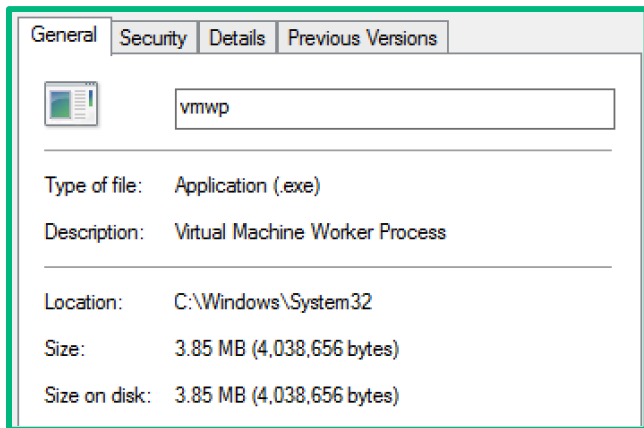


Reverse Engineering - Debugging

- Debugging via WinDBG is supported
 - Serial Port, Firewire, Ethernet...
- Fortunately no dedicated hardware is needed
 - VMware supports nested virtualization
- Hyper-V specific functionality requires hvexts.dll
 - Not publicly available
 - Have a copy? Drop us a mail 😊



Reverse Engineering - VSP



- Virtualization Service Provider
 - Implemented in Kernel (Storage, Networking) or Userspace (Framebuffer, HID)

- Public debugging Symbols are “sometimes” available
 - vmswitch.sys 😊
 - storvsp.sys 😞
 - vmwp.exe 😞 (Many debugging strings)



Attack Surface



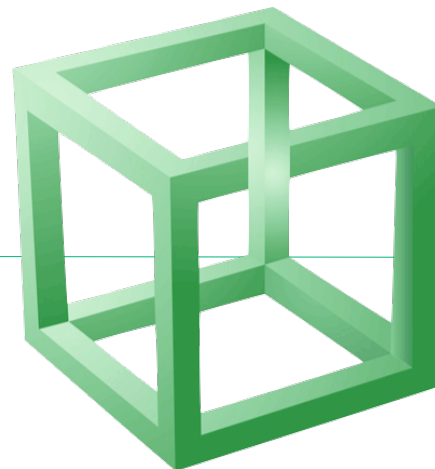
Attack Surface

- Based on the taxonomy presented earlier:
 - CPU: VM Exits
 - Device Virtualization: VMBus, Emulated Devices, RemoteFX
 - APIs: VMBus, Hypercalls

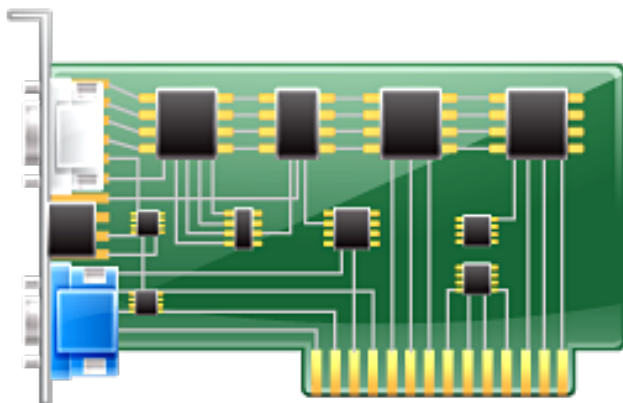




Emulated Devices

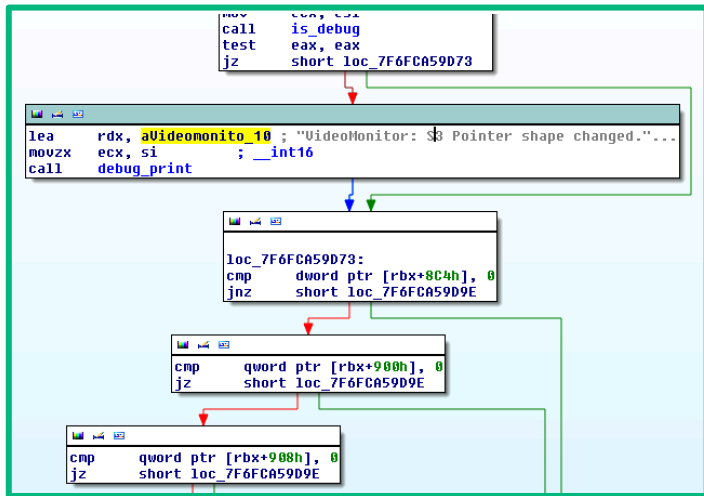


Emulated Devices



- Emulation of standard devices to allow unenlightened guest partitions
- Includes
 - Network adapter
 - S3 Trio graphic card
 - Keyboard / Mouse
 - IDE Controller
- Implemented in vmwp.exe

Emulated Devices

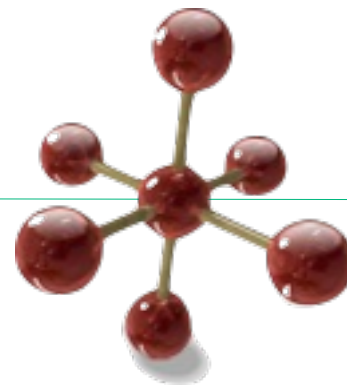


- We performed basic fuzzing (think IOFUZZ)
 - Permanent reboots
 - Worker process eats 100% CPU and requires hard kill
 - ... no interesting results besides that.

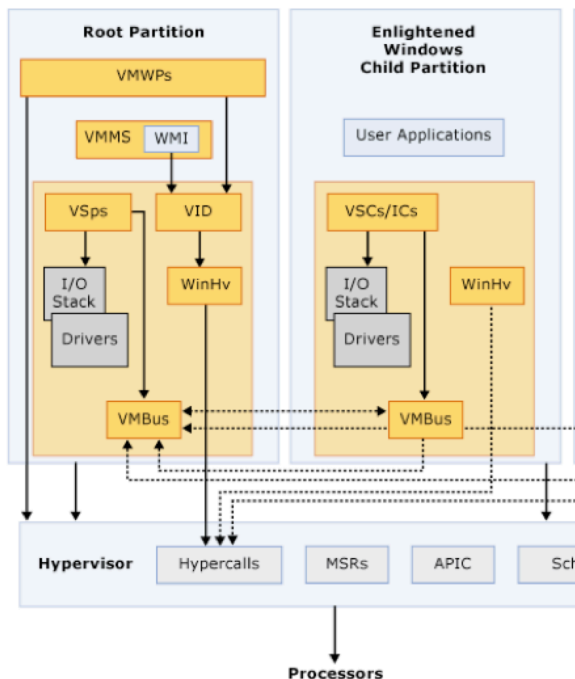
- Basic static analysis
 - No symbols, but debug strings.
 - Lots of sanity checks.



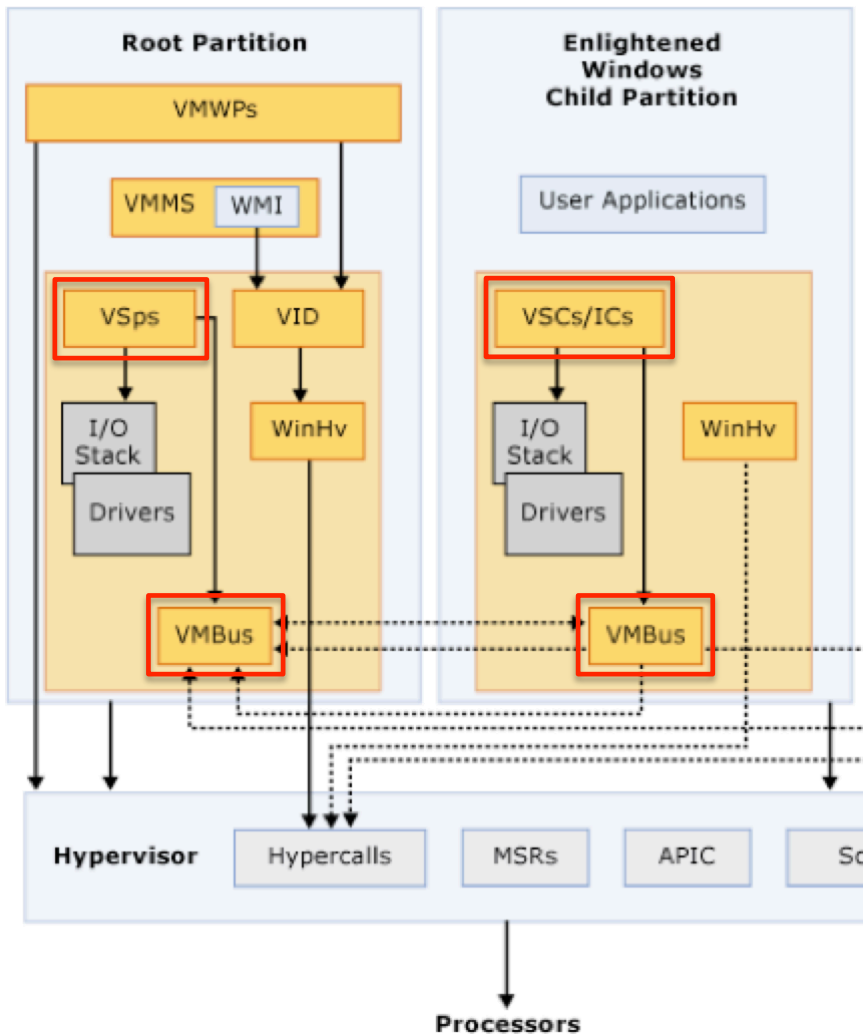
VMBus



VMBus

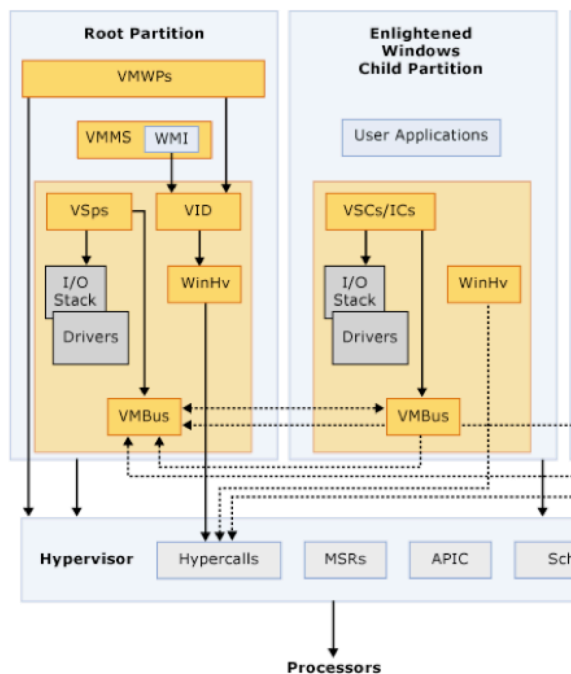


- Message bus for communication between partitions
- Default configuration:
 - Only communication to and from root partition allowed
- Memory pages that are mapped for multiple partitions
 - Heavily used for performance critical tasks in enlightened partitions
- Large attack surface



VMBus Components

VMBus Architecture



Initialization:

- PostMessage HyperCall
 - Is used to create VMBus between Partitions
- Negotiates/sets up shared memory sections

Communication:

- One or more channels are used to isolate communication between different functionality
 - Think of IP transport and TCP connections
- Each channel has two ring buffers, one for inbound and one for outbound communication
- Large transfers are performed via GPADLs



VMBus-based Services

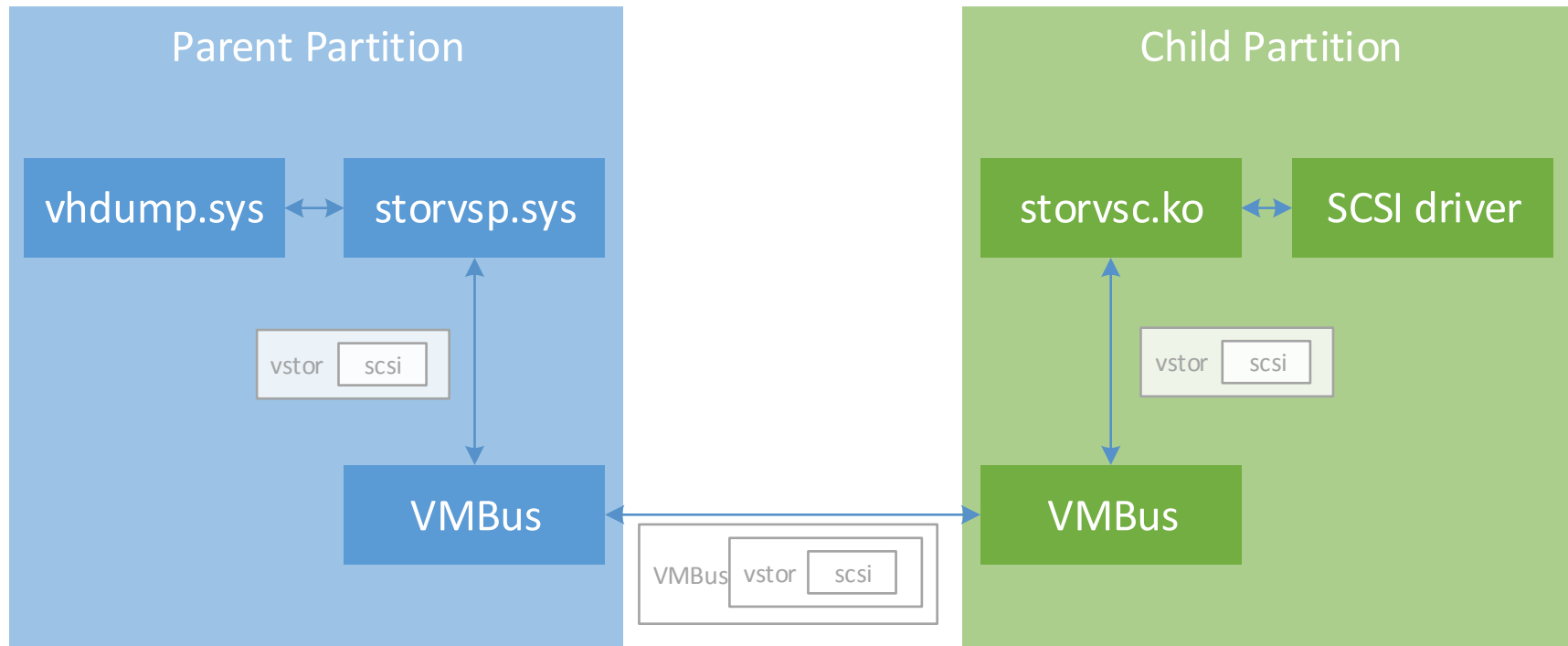
Term	Description	Driver
Networking	VM traffic is sent via VMBus to the Hyper-V Virtual Switch running in the parent partition.	Child: hv_netvsc.ko Parent: vmswitch.sys
Storage	Child partitions send (mainly) SCSI commands encapsulated in so-called <i>vstor</i> packets.	Child: hv_storvsc.ko Parent: storvsp.sys, vhdump.sys
Utilities	Different minor functionality such as heartbeats, time synchronization, or shutdown is available via SyncIC messages	Child: hv_util.ko Parent: vmwp.exe
Framebuffer	Virtual child console in parent Hyper-V Manager tool	Child: hyperv_fb Parent: vmwp.exe

Analysis

- Network and Storage: Interesting attack surface
 - Parent endpoints in kernel space
 - In contrast to VMWP.exe running under non-admin user
- Potential vulnerabilities:
 - Protocol/logical flaws
 - Memory corruption in parent partition drivers/processes



Storage



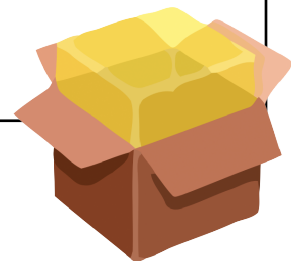
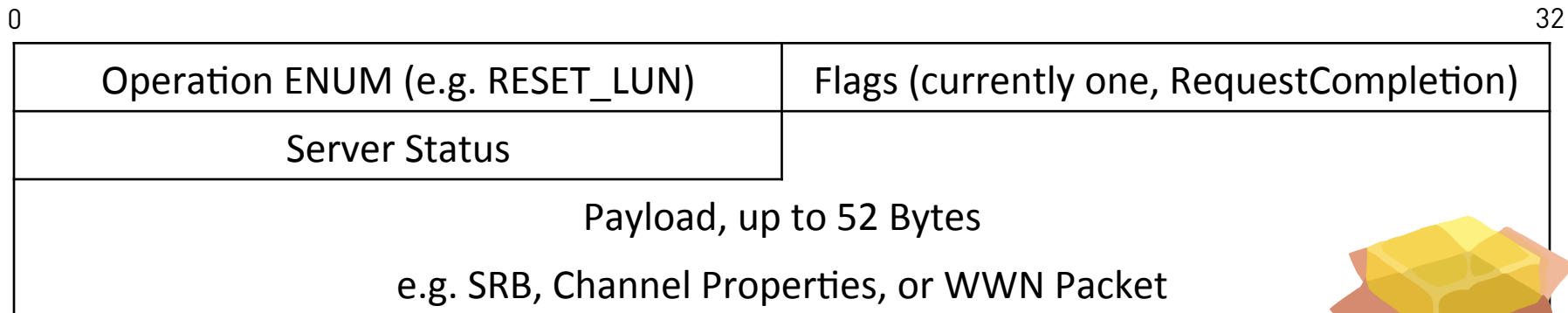


vstor_dmp

```
# ./vstor_dmp
0x000000: 03 00 00 00 01 00 00 00 .....
0x000008: 00 00 00 00 34 00 00 00 ....4...
0x000010: 02 00 00 00 0a 14 00 00 .....
0x000018: 00 20 00 00 2a 00 07 c5 . . .*...
0x000020: ee 70 00 00 10 00 00 00 .p.....
0x000028: 00 00 00 00 00 00 .....
Operation:      03 00 00 00 (VSTOR_OPERATION_EXECUTE_SRB)
Flags:          01 00 00 00
Server Status:  00 00 00 00
```



vstor_packet





0 VMBus Packet

32

Type	DataOffset
Length	Flags
Transaction ID	
Reserved	
Rangecount	
Page Buffer Range (array of PFNs, multiples of 16 bytes)	
Payload e.g. vstor_packet	

Analysis



- Performed so far:
 - Dumb fuzzing
 - → Guest VM crashes due to VSC failures
 - Manual protocol analysis
 - No results, but difficult to do due to high amount of VMBus packets
 - → Debugging is PITA
 - Basic static analysis
 - No interesting results
 - Checked for banned functions (strcpy, sprintf, you name it)
 - Not too surprising, SDL works here

Basic Fuzzing



- Implemented using kprobes
 - <https://www.kernel.org/doc/Documentation/kprobes.txt>
- Allows to hook (almost) arbitrary kernel functions and tamper with existing data.
- Approach:
 - Hook `packet_send` functions and tamper with `packet` arguments.



Hypercalls



Hypercall Interface

5.6.1 HvCreatePartition

The HvCreatePartition hypercall allows an authorized guest to create a new partition.

Wrapper Interface

```
HV_STATUS  
HvCreatePartition(  
    __in  UINT64          Flags,  
    __in  HV_PROXIMITY_DOMAIN_INFO ProximityDomainInfo,  
    __out PHV_PARTITION_ID NewPartitionId  
);
```

Native Interface

HvCreatePartition	
	Call Code = 0x0040
➔ Input Parameters	
0	Flags (8 bytes)
8	ProximityDomainInfo (8 bytes)
◀ Output Parameters	
0	NewPartitionId (8 bytes)

- Like system calls but for communication between VM kernel and hypervisor
- Documented interface
 - And known security boundary
 - RCX = Call Number
 - RDX = Input GPA
 - R8 = Output GPA
- Used by enlightened partitions to improve performance
- Root partition manages other partition using hypercalls

Hypercall Interface

```
FFFFF800034C5B7E mov     rdx, cr3
FFFFF800034C5B81 mov     ecx, VMCS_Host_CR3
FFFFF800034C5B86 call    vmwrite_wrapper
FFFFF800034C5B8B test    eax, eax
FFFFF800034C5B8D jnz     loc_FFFF800034C5CDE
```

```
FFFFF800034C5B93 mov     rdx, cr4
FFFFF800034C5B96 mov     ecx, VMCS_Host_CR4
FFFFF800034C5B9B call    vmwrite_wrapper
FFFFF800034C5BA0 test    eax, eax
FFFFF800034C5BA2 jnz     loc_FFFF800034C5CDE
```

```
FFFFF800034C5BA8 lea    rdx, vm_exit_handler
FFFFF800034C5BAF mov     ecx, host_rip
FFFFF800034C5BB4 call    vmwrite_wrapper
FFFFF800034C5BB9 test    eax, eax
FFFFF800034C5BBB jnz     loc_FFFF800034C5CDE
```

- Auditing hypercalls requires finding the handler functions
- First Step:
 - Finding main VM Exit handler which is stored inside VMCS field

Hypercall Interface

- VM Exit handler looks similar in all VMM implementations:
 - Store Guest State
 - Perform action depending on VMCS Exit Reason
 - Restore Guest State

- Hypercalls are executed using “vmcall” instruction.

```

FFFFF8000342596D
FFFFF8000342596D loc_FFFF8000342596D:
FFFFF8000342596D mov     eax, VMCS_Exit_reason
FFFFF80003425972 vmread  rax, rax
FFFFF80003425975 movzx   edx, ax
FFFFF80003425978 mov     rcx, 0C00191D000F600h
FFFFF80003425982 mov     [rbp+40h+var_A8], rax
FFFFF80003425986 bt      rcx, rdx
FFFFF8000342598A jnb     short loc_FFFF8000342599C
  
```

```

FFFFF8000342598C mov     eax, VMCS_VM_exit_instruction_length
FFFFF80003425991 vmread  rax, rax
FFFFF80003425994 mov     [rsp+128h+var_D0], rax
FFFFF80003425999 mov     [rdi+4], al
  
```

```

FFFFF8000342599C
FFFFF8000342599C loc_FFFF8000342599C:
FFFFF8000342599C mov     [rsp+128h+var_8], rbx
FFFFF800034259A4 mov     [rsp+128h+var_10], rsi
FFFFF800034259AC mov     [rsp+128h+var_18], r13
FFFFF800034259B4 mov     [rsp+128h+var_20], r15
FFFFF800034259BC test   edx, edx ; test on exit reason
FFFFF800034259BE jnz     loc_FFFF800034259AE3
  
```



Hypercall Interface

```
dq offset HvCallSwitchVirtualAddressSpaceFunc
dw HvCallSwitchVirtualAddressSpace ; call
dw 0 ; rep
dw 8 ; inpi
dw 0 ; unki
dw 0 ; outj
db 2 dup(0), 43h, 3 dup(0) ; unki
dq offset HvCallFlushVirtualAddressSpaceFunc;
dw HvCallFlushVirtualAddressSpace ; call
dw 0 ; rep
dw 18h ; inpi
dw 0 ; unki
dw 0 ; outj
db 2 dup(0), 43h, 3 dup(0) ; unki
dq offset HvCallFlushVirtualAddressListFunc;
dw HvCallFlushVirtualAddressList ; call
dw 1 ; rep
dw 18h ; inpi
dw 8 ; unki
dw 0 ; outj
db 2 dup(0), 43h, 3 dup(0) ; unki
dq offset HvCallGetLogicalProcessorRunTimeFunc;
dw HvCallGetLogicalProcessorRunTime ; call
dw 0 ; rep
dw 8 ; inpi
dw 0 ; unki
dw 20h ; outj
.. ; ..
```

- RCX register stores hypercall number
- Is used as index into handler table
- Handler Table allows quick identification of all handler functions

Hypercall Interface

```

FFFFF800007043B8
FFFFF800007043B8
FFFFF800007043B8
FFFFF800007043B8 HvCallGetPartitionIdFunc proc near
FFFFF800007043B8 sub     rsp, 28h
FFFFF800007043BC mov     rcx, gs:82E8h
FFFFF800007043C5 mov     r9, rdx
FFFFF800007043C8 mov     rdx, 200000000h
FFFFF800007043D2 call    hypercall_check_privileges
FFFFF800007043D7 test    eax, eax
FFFFF800007043D9 jnz     short loc_FFFF800007043EE
  
```

```

FFFFF800007043DB mov     rcx, gs:82E8h
FFFFF800007043E4 mov     r8, [rcx+400h]
FFFFF800007043EB mov     [r9], r8
  
```

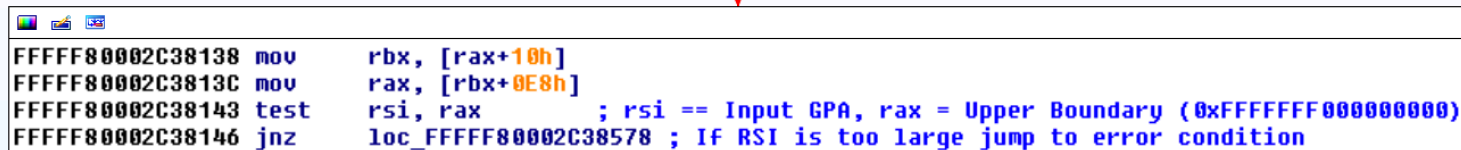
```

FFFFF800007043EE
FFFFF800007043EE loc_FFFF800007043EE:
FFFFF800007043EE add     rsp, 28h
FFFFF800007043F2 retn
FFFFF800007043F2 HvCallGetPartitionIdFunc endp
FFFFF800007043F2
  
```

- Hypercall handler are quite simple
- RCX = Input / RDX = Output
 - Copied from original input pages
- Early permission checks guard most complicated handler functions
- Performed fuzz testing but no interesting results.

However...

- Sanity checks are executed before calling the actual hypercall handler:
 - Hypercall comes from Ring 0?
 - Input and Output GPA are aligned?
 - Input and Output GPA look sane?



A screenshot of assembly code with annotations. A red arrow points to the first instruction. A red arrow points to the jump instruction. A green line connects the jump instruction to the text 'If RSI is too large jump to error condition'.

```

FFFFF80002C38138 mov    rbx, [rax+10h]
FFFFF80002C3813C mov    rax, [rbx+0E8h]
FFFFF80002C38143 test   rsi, rax           ; rsi == Input GPA, rax = Upper Boundary (0xFFFFFFFF00000000)
FFFFF80002C38146 jnz    loc_FFFF80002C38578 ; If RSI is too large jump to error condition
  
```

The Bug

```
mov     r8, [rcx+ept_object.pte_table]
mov     rax, rdx
shr     rax, 12
mov     rbp, rdx
mov     rsi, rcx
mov     r8, [r8+rax*8] ; rax = rdx = r
                        ; r15 = Hyperca
mov     rax, 8000000000000000h
test    rax, r8
jz     return_loc
```

- Guest GPAs $\geq 0x1000000000$ and $< 0x10000000000000$ set special error condition
- Error function tries to map GPA to system PTE
 - r8 points to table
 - rax = input_gpa $\gg 12$;
- Classic out-of-bound access because GPA can be almost unlimited large.
- The trigger?
 - hypercall(0xXY, input_address = 0x4141414141414141...) ☺

Fixed with MS13-092



- Denial of Service of complete hypervisor
- Potentially privilege escalation to other virtual machines
- Azure affected !

Acknowledgments

Microsoft **thanks** the following for working with us to help protect customers:

- thinktecture (www.thinktecture.com) & ERNW (www.ernw.de; Felix Wilhelm) on behalf of Bundesamt für Sicherheit in der Informationstechnik (BSI, German Federal Office for Information Security) for reporting the Address Corruption Vulnerability (CVE-2013-3898)



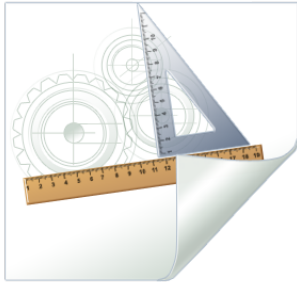
Some thoughts on exploitation...

- DoS trivially possible by accessing invalid memory.
- Privilege Escalation possible?

```
1 long pte = gpa_to_spa[guest_physical_addr >> 12]
2 if (!(pte & 0x8000000000000000) || (pte & 0x1E00000000000000))
3     return 0;
4 if (pfn(pte) != SPECIAL_ADDR)
5     return 0;
6 if (last_three_bits(pte) & 0b111)
7     return 1;
8 acquire_locks();
9 func_x(pte | 5, guest_physical_addr)
10 return 1;
```



Some thoughts on exploitation...



- Patch adds additional range check in front and force-returns 0
 - Code flow called after vuln function returns 0 is uninteresting
- Can't read attacker controlled value
 - Upper limit + 64bit address space
- However, offset to “next” virtual machine PTE table is constant
 - Possible to access PTEs of other VMs



Some thoughts on exploitation...

```
if (pfn(pte) != SPECIAL_ADDR)
    return 0;
```

- Limited to a “special” PFN number
 - Could not clearly identify its use case + not used by VMs in our lab :/
- Code flow after successful checks is interesting
 - But can only be minimally influenced by attacker
 - Two values used: PTE and INPUT_GPA (both are more or less fixed)
- Classic Memory corruption not possible
 - But triggering a PTE map from/to a different partition would be game over as well.



Some thoughts on exploitation...

- Patch is noisy
 - *Pretty sure*TM that we did not miss any additional security relevant modifications

- Might have missed something completely obvious..





Your PC ran into a problem and needs to restart. We're just collecting some error info, and then we'll restart for you. (40% complete)

If you'd like to know more, you can search online later for the error: `HYPervisor_ERROR`

Demo

Future Work



- Improve fuzzing
 - Think of analysis approach that allows high coverage and automation without constantly crashing the guest.

- Look at new stuff
 - RemoteFX
 - Gen2 VMs
 - Xbox One

- Go deeper
 - VM exit handling

Conclusions



- Hypervisor security & research is not hard.
 - Just software and big attack surface.
- Hypervisor security & research is hard.
 - .. but you learn a lot!
- Hyper-V is solid software
 - .. but still has critical security bugs.
- Hypervisors are getting faster/more features but not more secure.
 - .. still not enough public research!
- *Make the theoretical practical.*



There's never enough time...

THANK YOU...



@_fel1x
@uchi_mata
@Enno_Insinuator



fwilhelm@ernw.de
mluft@ernw.de
erey@ernw.de



...for yours!

Code & Slides:
<https://www.insinuator.net>
(..soon)



Sources & Mentions

- Gal Diskin, Virtually Impossible
- Gerhart,
<http://www.securitylab.ru/contest/444112.php>
- Hypervisor Functional Specification, Microsoft





Disclaimer

All products, company names, brand names, trademarks and logos are the property of their respective owners!

