

# Bluepillling the Xen Hypervisor

Joanna Rutkowska & Alexander Tereshkin  
Invisible Things Lab

Black Hat USA 2008, August 7th, Las Vegas, NV

# **Xen 0wning Trilogy**

Part Three

Previously on Xen Owning Trilogy...

# Part I: “Subverting the Xen Hypervisor”

by Rafal Wojtczuk (Invisible Things Lab)

- ✓ Hypervisor attacks via DMA
  - ✓ TG3 network card “manual” attack
  - ✓ Generic attack using disk controller
- ✓ “Xen Loadable Modules” framework :)
- ✓ Hypervisor backdooring
  - ✓ “DR” backdoor
  - ✓ “Foreign” backdoor

# Part II: “Detecting and Preventing the Xen Hypervisor Subversions”

by Rafal Wojtczuk & Joanna Rutkowska

- ✓ Latest Xen security features
- ✓ How they fail: Q35 exploit
- ✓ How they fail: FLASK exploit
- ✓ The need for hypervisor integrity checks!
- ✓ Introducing HyperGuard!

Now, in this part...

1 **Nested virtualization** (“Matrix inside Matrix”)

2 **BluePillBoot**

3 **XenBP**: Bluepilling the Xen hypervisor **on the fly!**

4 Bluepilled Xen **detection**



# Nested Virtualization

Hypervisor (Primary)

VM<sub>1</sub>

VM<sub>2</sub> (Nested Hypervisor)

VM<sub>3</sub>

VM<sub>4</sub>

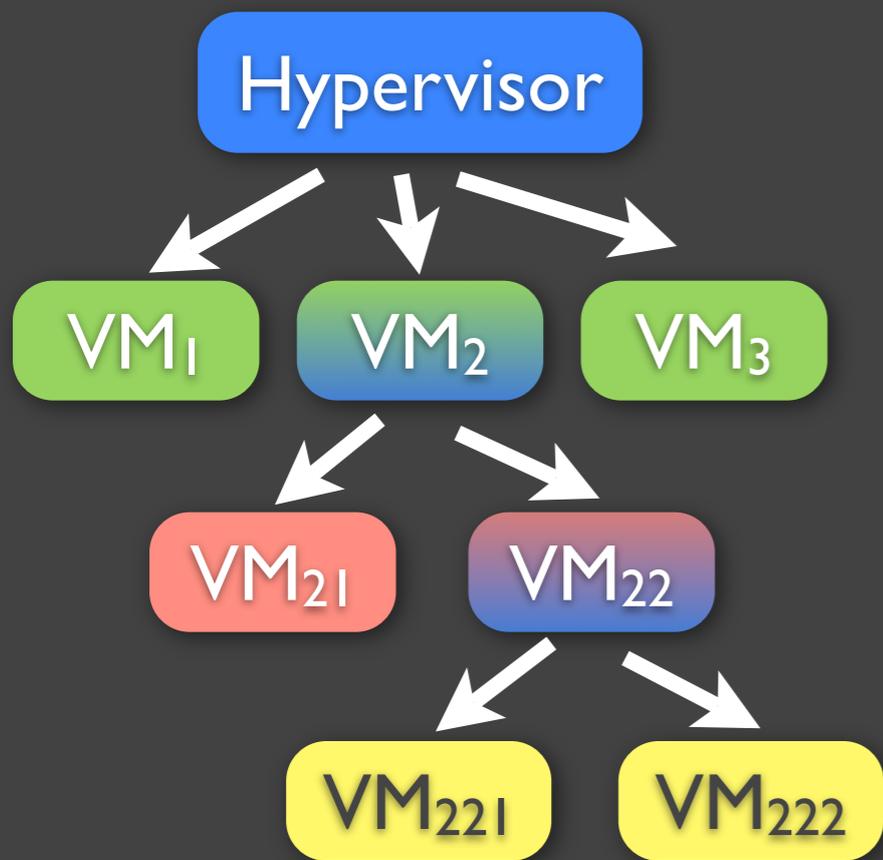
VM<sub>21</sub>

VM<sub>22</sub>

VM<sub>221</sub>

VM<sub>222</sub>

Idea of how to handle this situation...

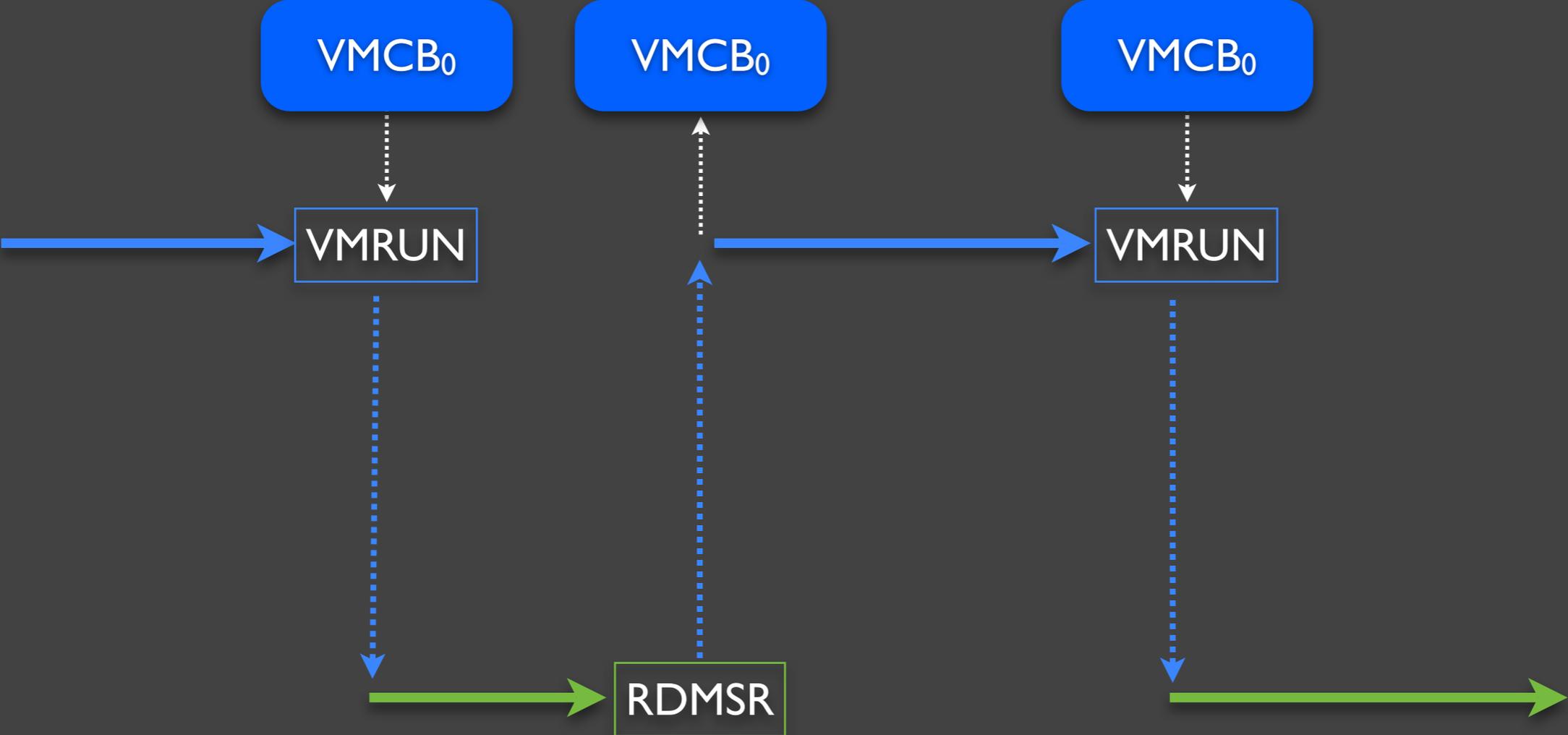


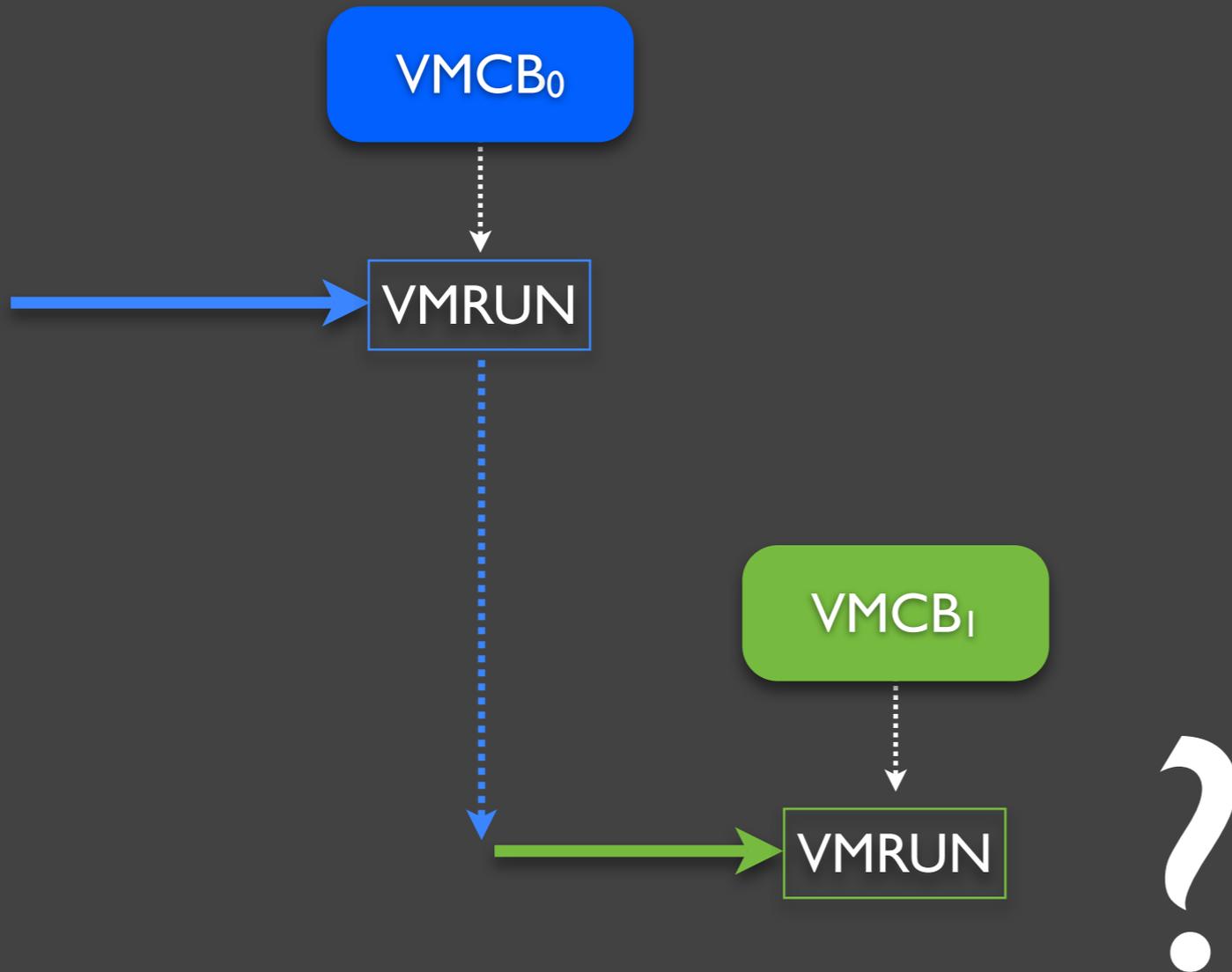
Hypervisor

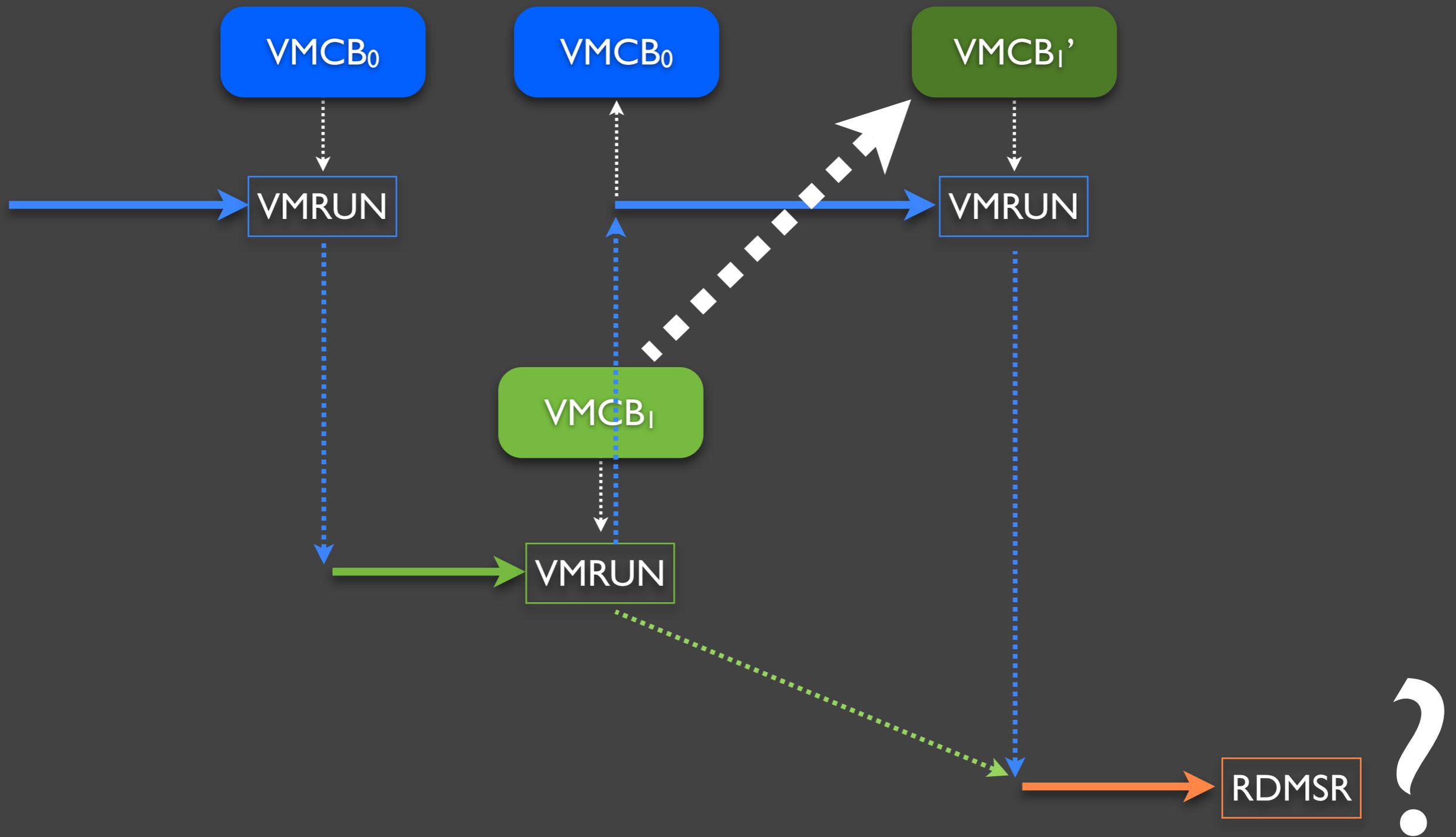


Now, lets look at the actual details :)

Let's start with AMD-V...

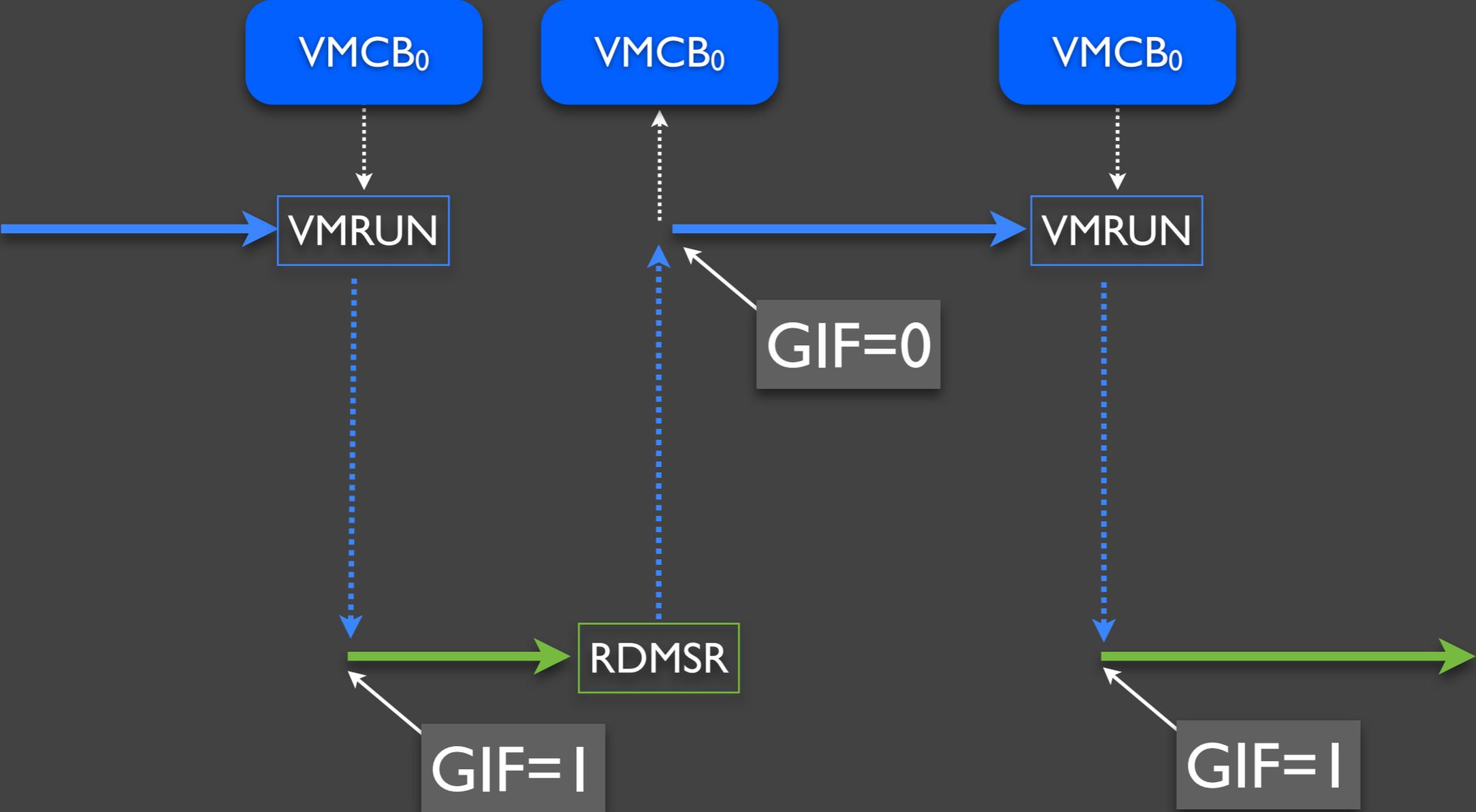


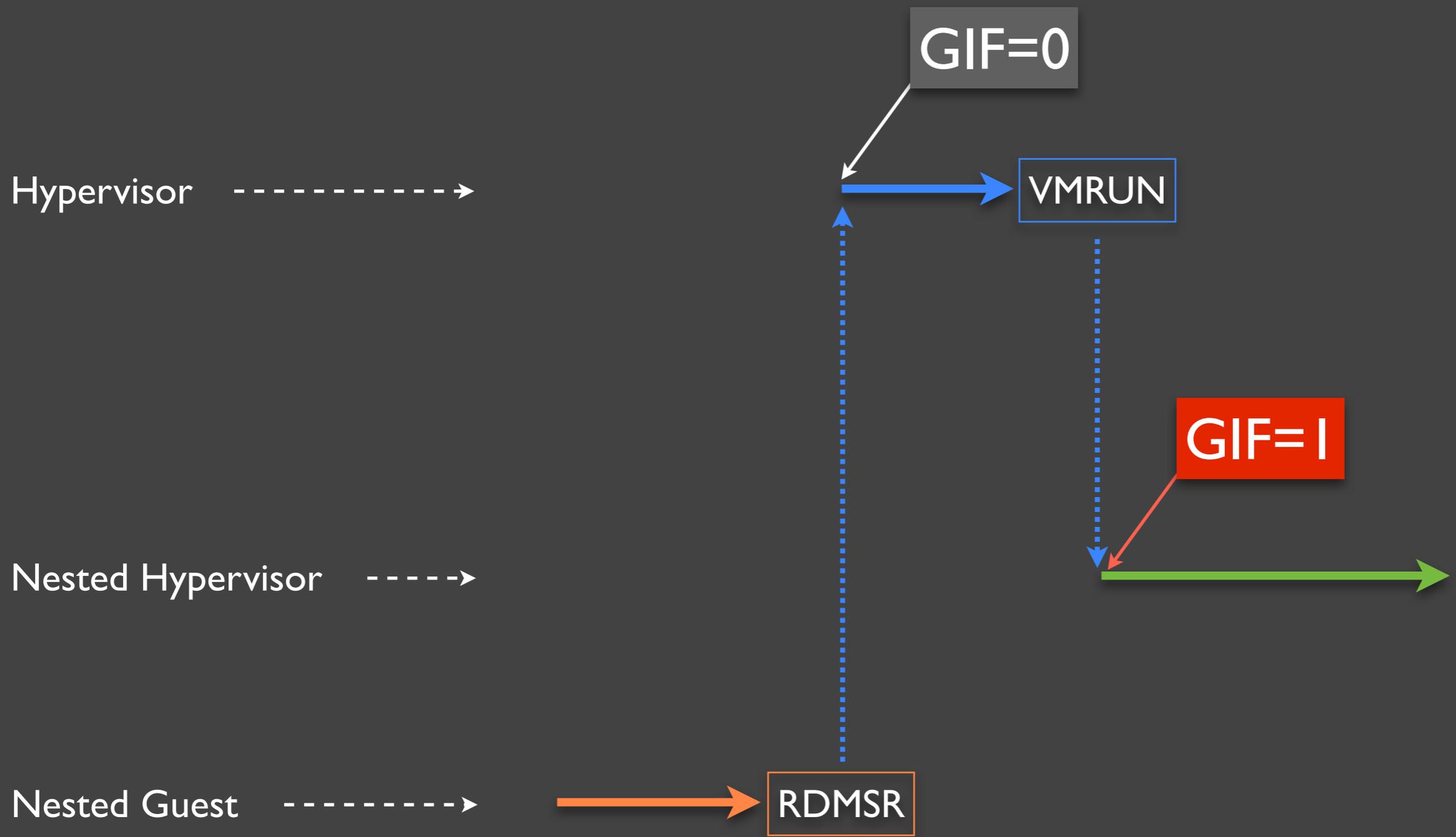






Looks convincing but we also need to take care about some technical details, that are not trivial...





- Hypervisors expect to have GIF=1 when VMEXIT occurs...
- They might not be prepared to handle interrupts just after VMEXIT from guests!
- ... but when we resume the nested hypervisor CPU sets GIF=1, because we do this via VMRUN, not VMEXIT...

# Getting around the “GIF Problem”

- We need to “emulate” that GIF is 0 for the nested hypervisor
- We stop this emulation when:
  - The nested hypervisor executes STGI
  - The nested hypervisor executes VMRUN
- How do we emulate it?

# GIF0 emulation

- $VMCB_i.V\_INTR\_MASKING = 1$
- Host's  $RFLAGS.IF = 0$
- Intercept NMI, SMI, INIT, #DB and held (i.e. record and reinject) or discard until we stop the emulation

# Additional details

- Need to also intercept VMLOAD/VMSAVE
- Need to virtualize VM\_HSAVE\_PA
- ASID conflicts

Hypervisor: ASID = 0

Conflicting ASIDs!

Nested Hypervisor: ASID = 1  
(but thinks that has ASID = 0)

Nested Guest: ASID = 1  
(assigned by the nested hypervisor)

But we can always reassign the ASID in the VMCB “prim”  
that we use to run the nested guest.

# Performance Impact

- One additional #VMEXIT on every #VMEXIT that would occur in a non-nested scenario
- One additional #VMEXIT when the nested hypervisor executes: STGI, CLGI, VMLOAD, VMSAVE
- Lots of space for optimization though

Intel VT-x

# Nested virtualization on VT-x

- No GIF bit - no need to emulate “GIF0” for the nested hypervisor :)
- No Tagged TLB - No ASID conflicts :)
- However:
  - VMX instructions can take memory operands - need to use complex operand parser
  - No tagged TLB - potentially bigger performance impact

# Nested VT-x: Status

- We have that working!
- The VT-x nesting code cannot be published though :(

Who else does Nested (hardware-based) Virtualization?

# IBM z/VM hypervisor on IBM System z™ mainframe

*“Running z/VM in a virtual machine (that is, z/VM as a guest of z/VM, also known as “second-level” z/VM) is functionally supported but is intended only for testing purposes for the second-level z/VM system and its guests (called “third-level” guests).”*

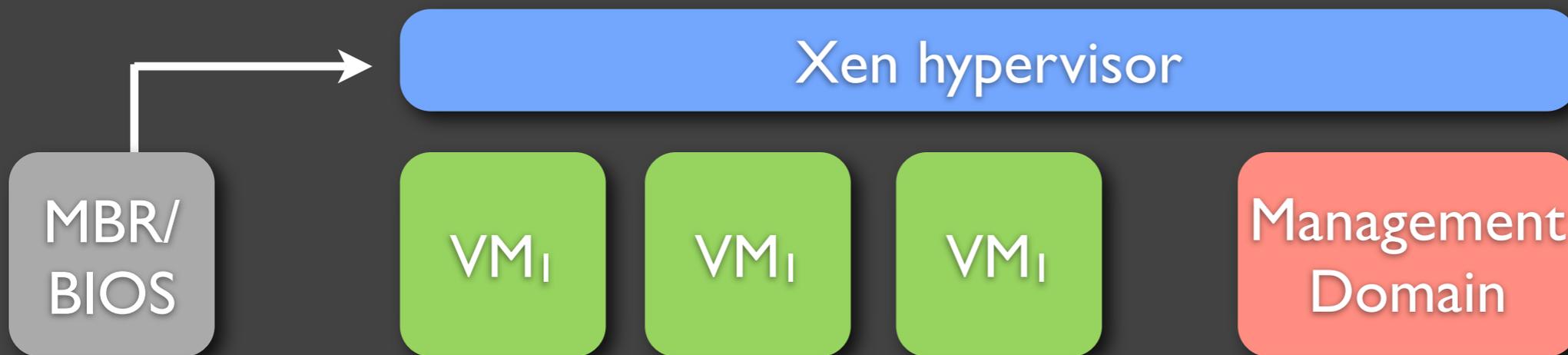
-- <http://www.vm.ibm.com/pubs/hcsf8b22.pdf>

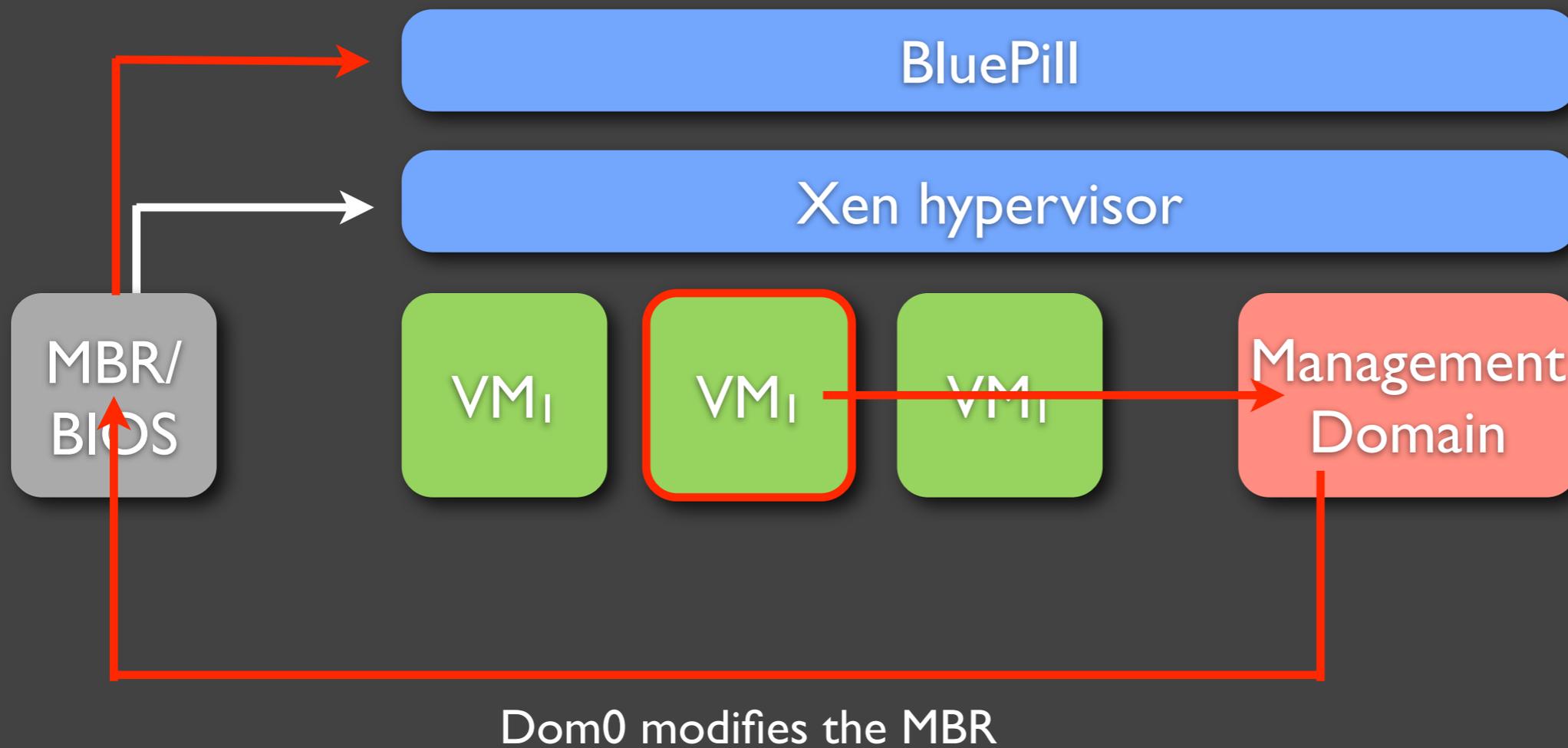


IBM System z10, source: ibm.com



Blue Pill Boot



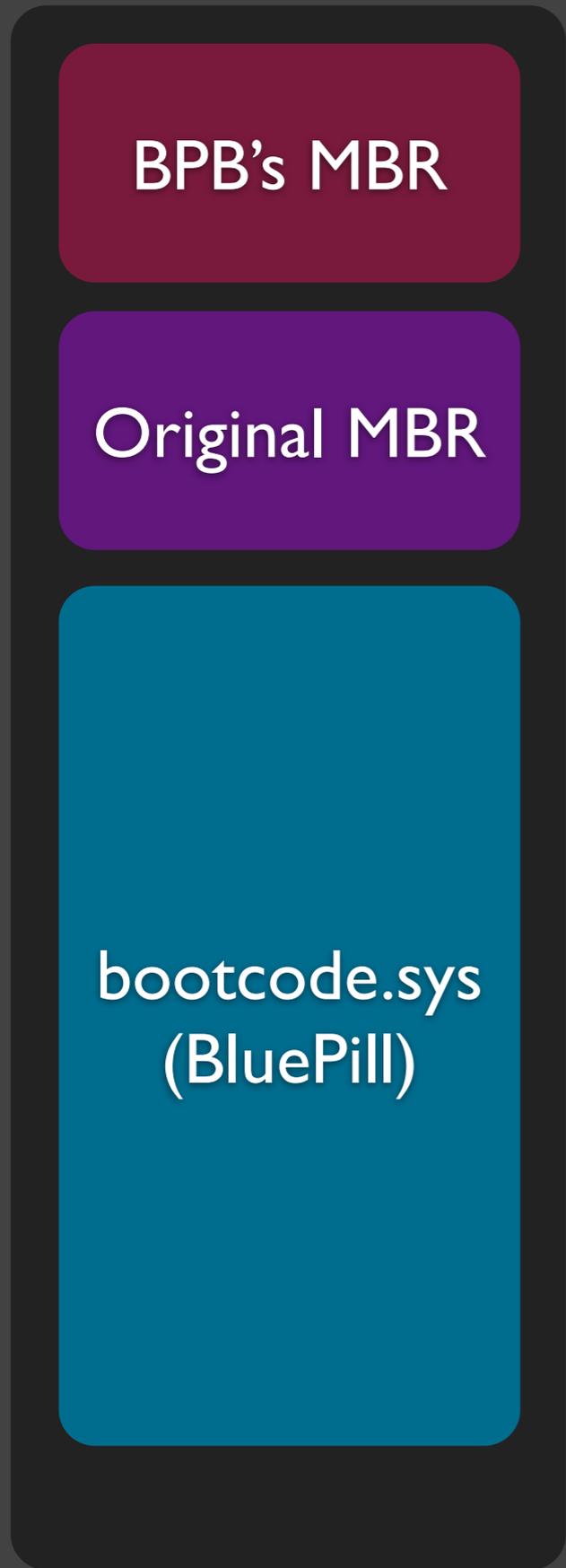


**Blue Pill Boot =**

MBR infector +

Blue Pill loader +

Blue Pill that supports nested virtualization



BPB's MBR

Original MBR

bootcode.sys  
(BluePill)

Disk



Sector 1

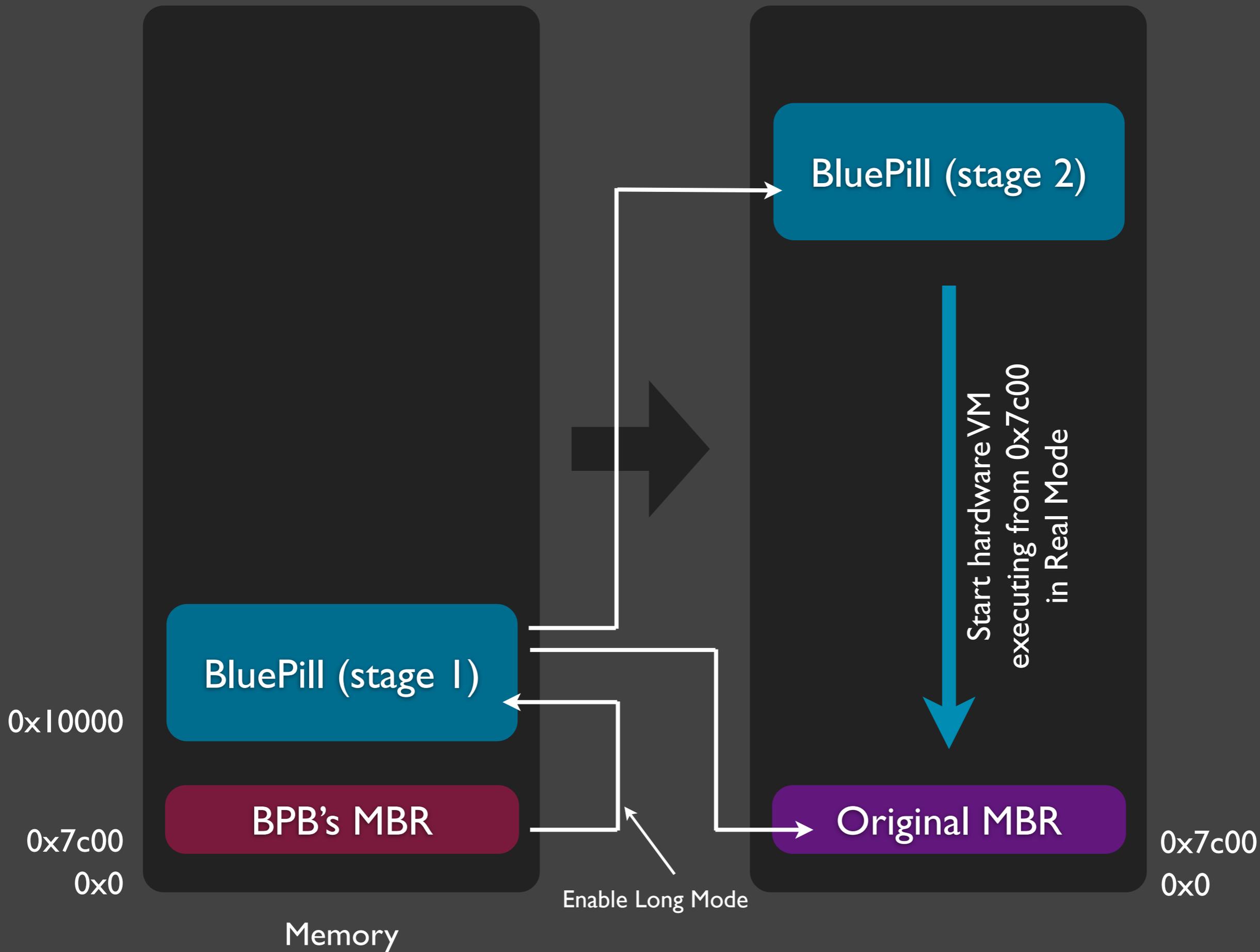


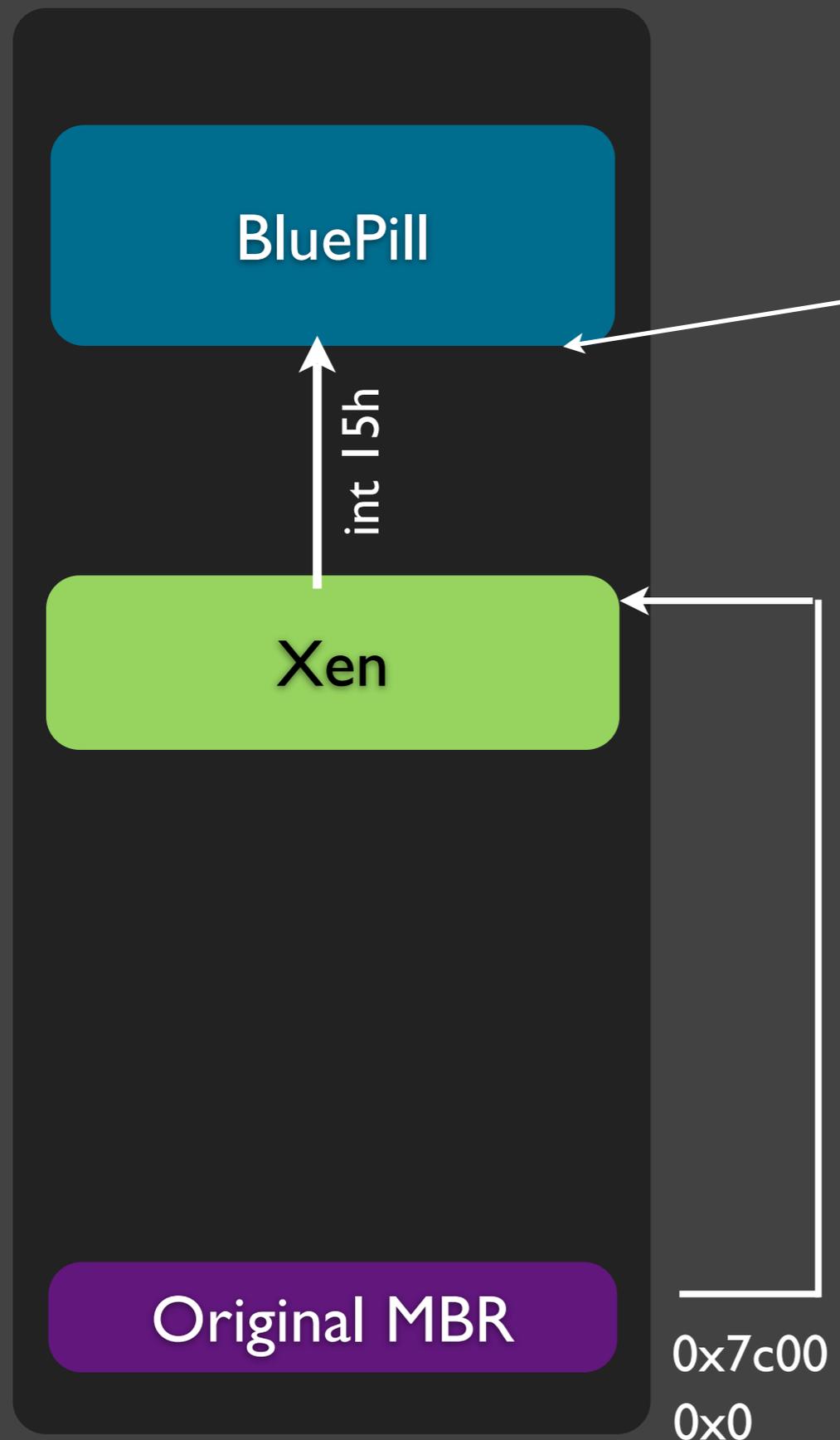
Sector 2



Sectors 3...n





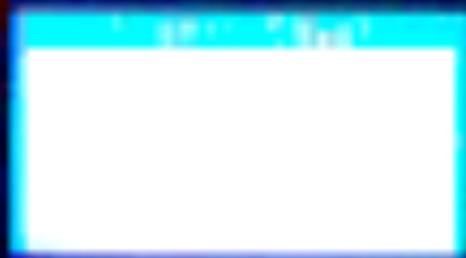


Int 15h/e820h queries are intercepted by BluePill

MBR starts Xen which now runs in a hardware virtual machine controlled by the BluePill

# Demo: BluePillBootting the Xen

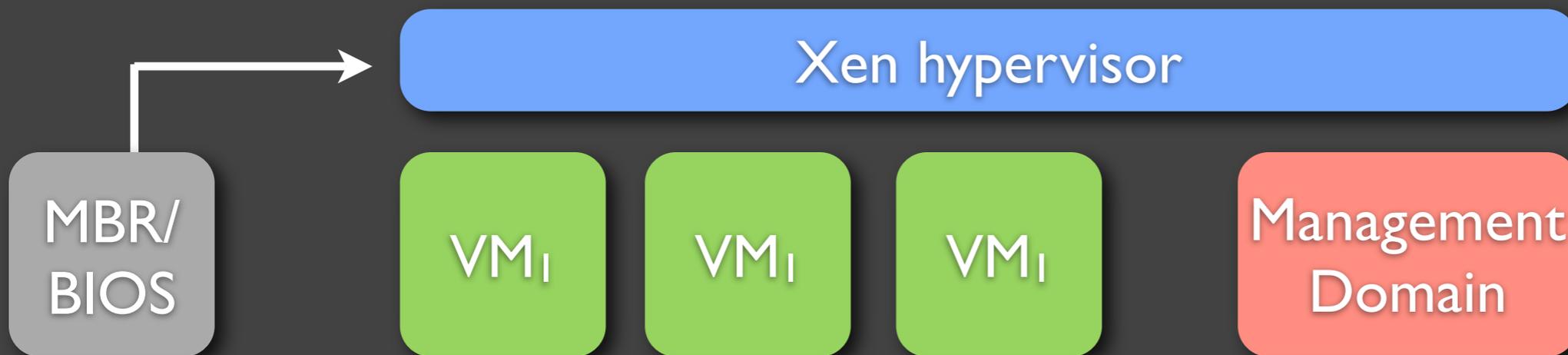
(please excuse the recording quality)

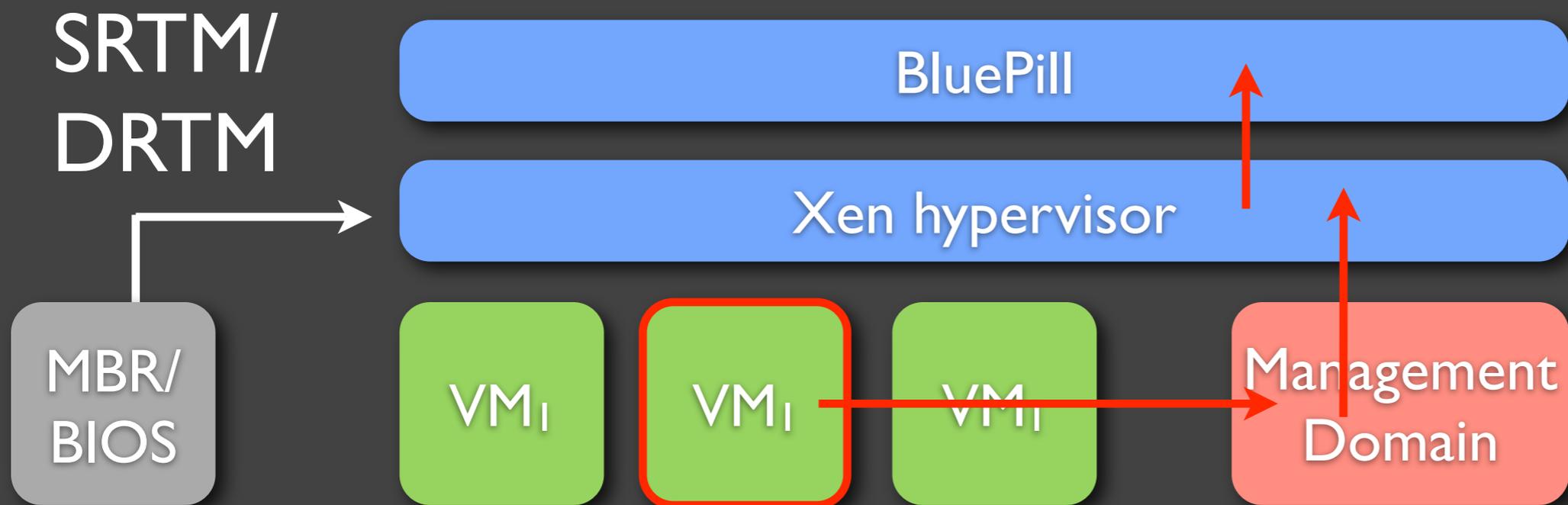


Ensure hypervisor integrity via SRTM or DRTM

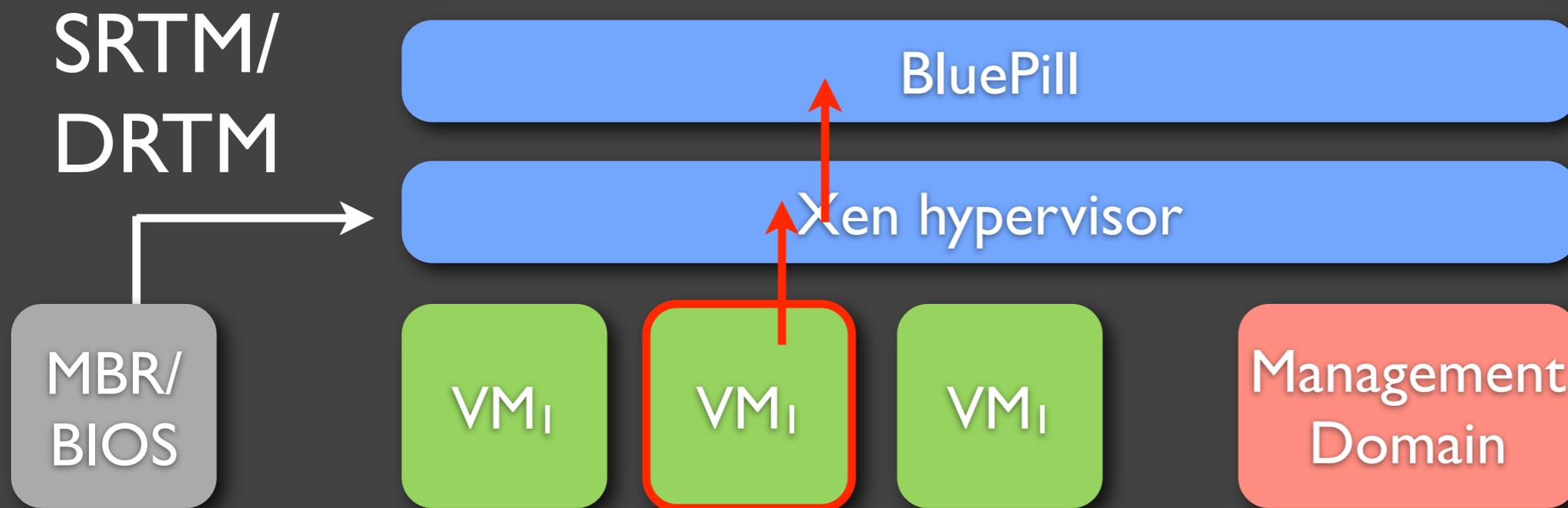


Xen Blue Pill





SRTM/DRTM do not protect the **already loaded hypervisor!**

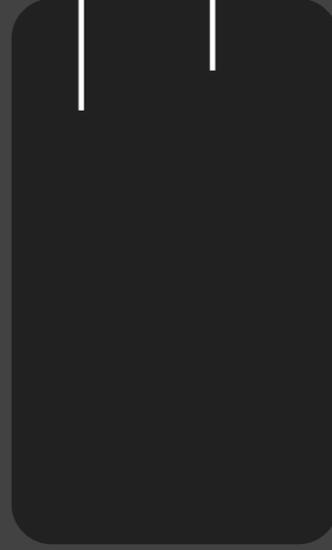
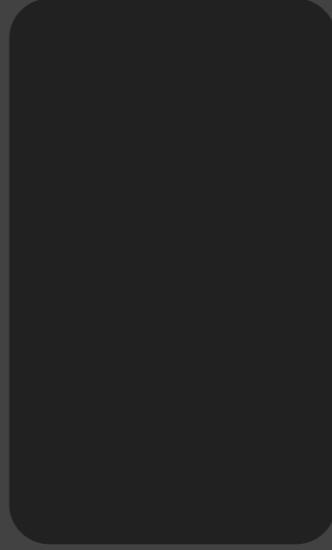


SRTM/DRTM do not protect the **already loaded hypervisor!**

The details

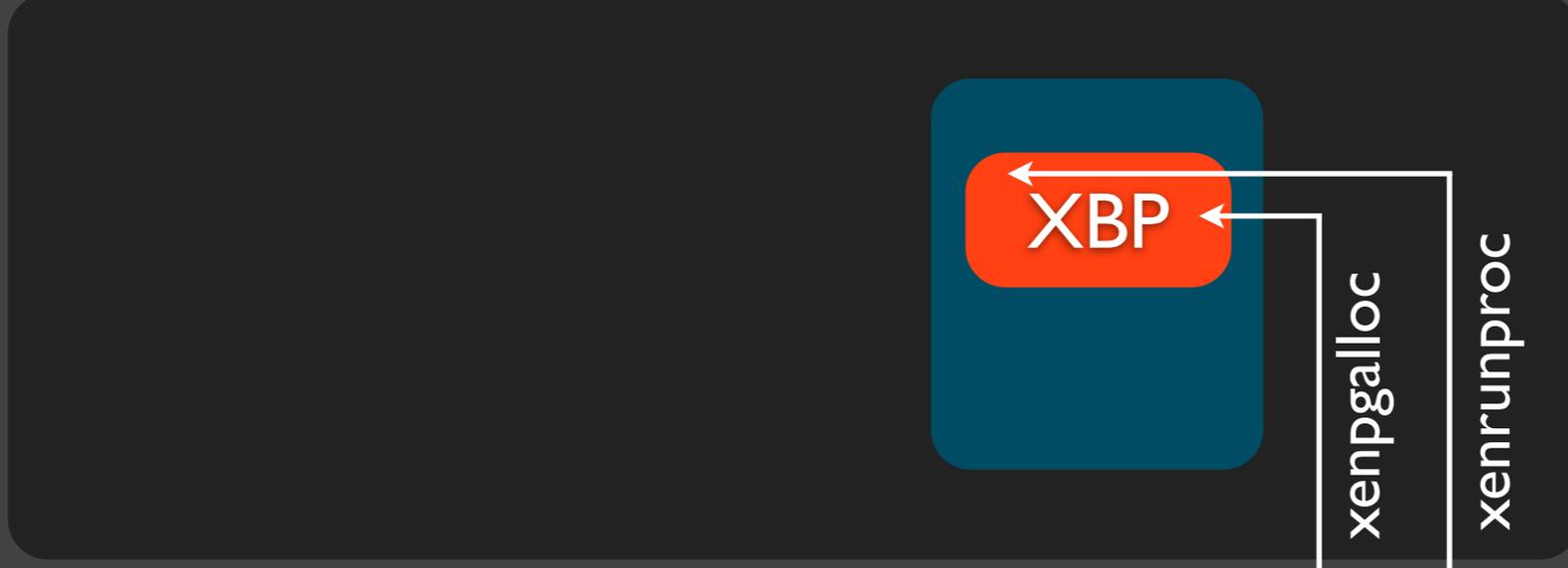
Loading using Rafal's XLM framework...

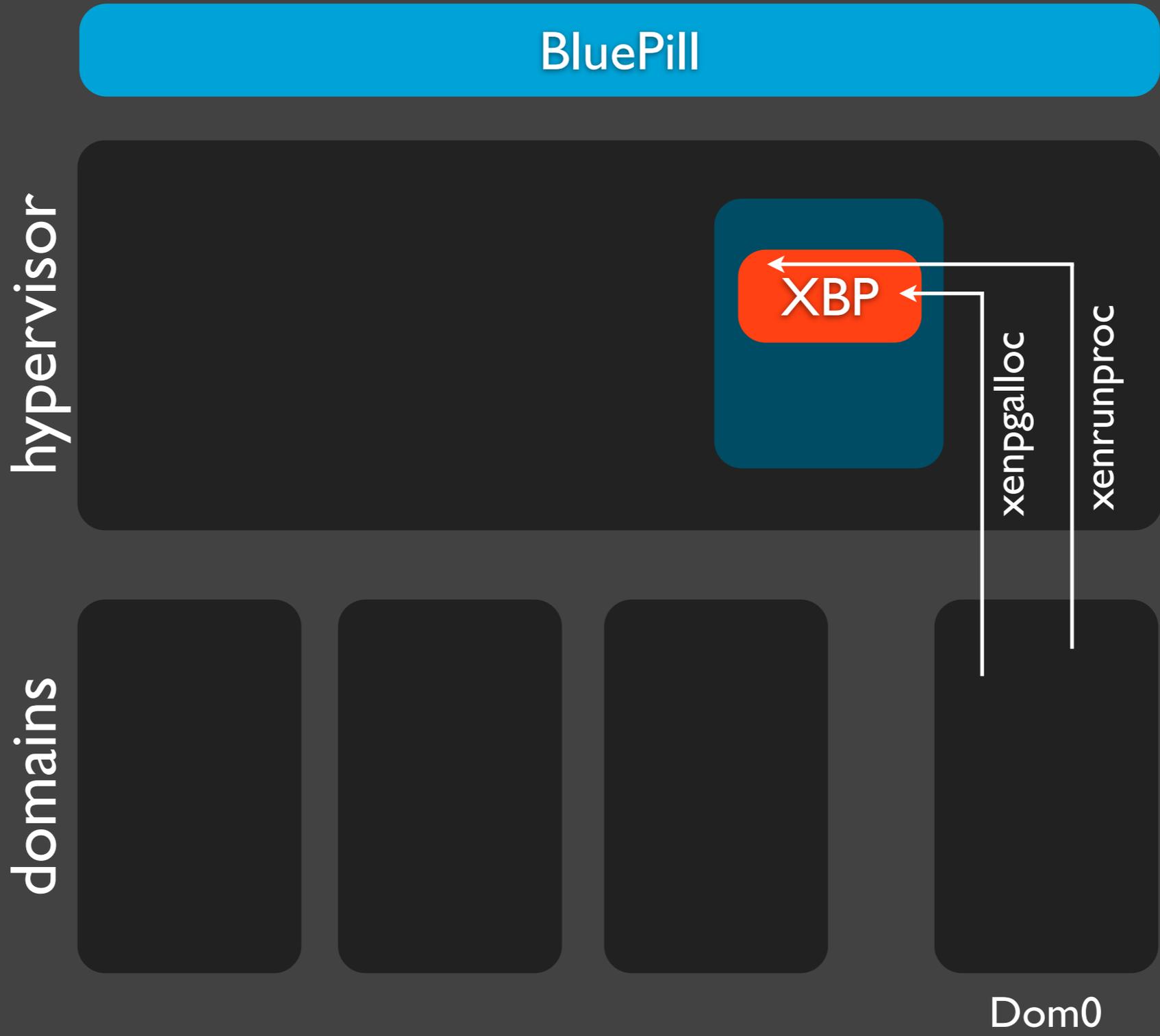
domains



Dom0

hypervisor





We allocate a block of memory for XBP inside Xen hypervisor -- this memory is used for both the XBP's code and data and heap

Demo: Bluepilling the Xen on the fly...

```
[root@turion64 ~]#
```

On Xen 3.3 we need  
to use Q35 exploit  
instead of direct hdd  
(see the talk #2)



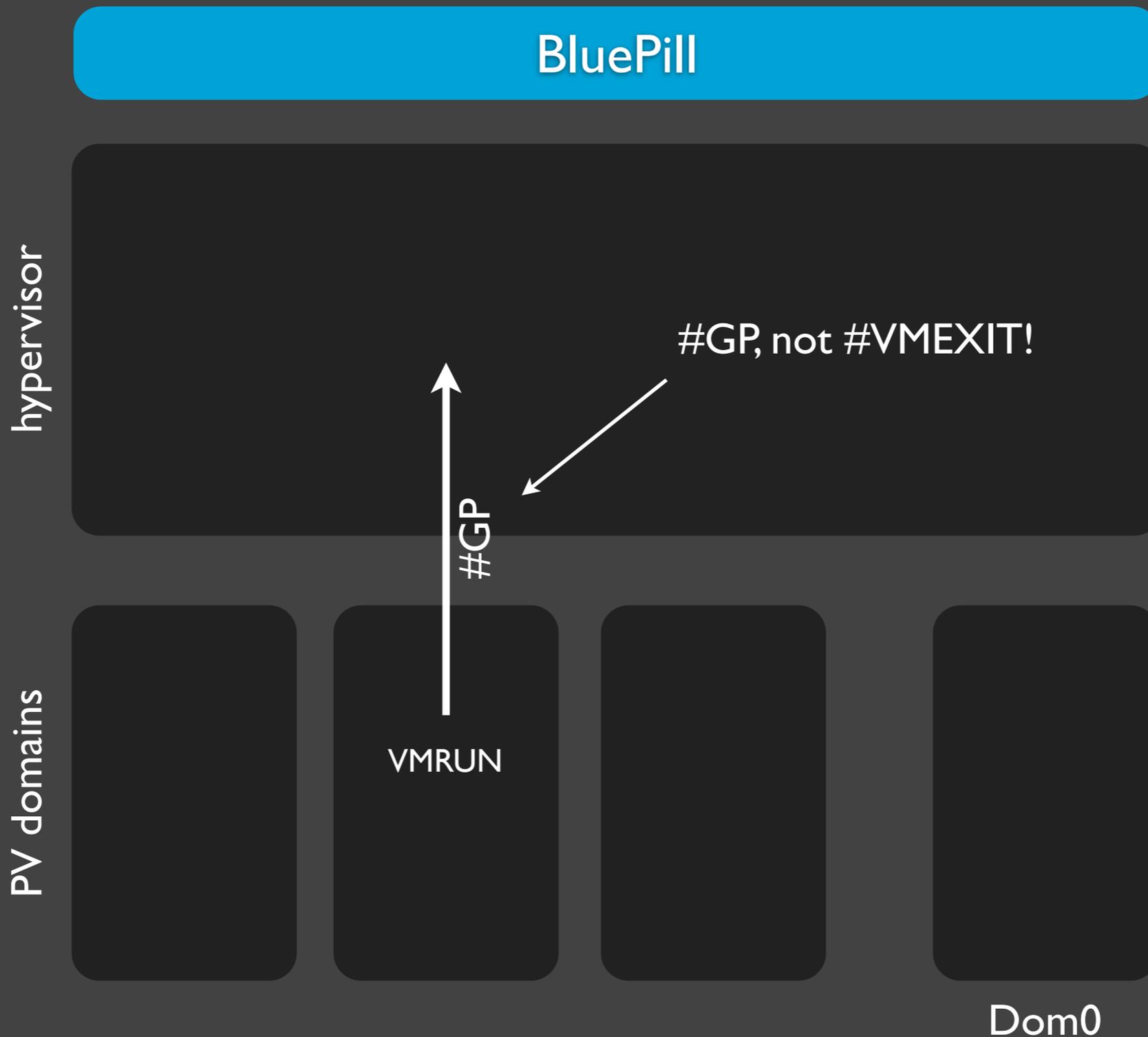
Detecting a VMM is now not enough...

... as we know there is already one VMM in the system  
already (i.e. the Xen)...

We can only try direct timing analysis to see if #VMEXITs  
will take longer time to execute...  
(then on “non-bluepilled” Xen)

Impact on PV domains

ring 3



We don't need to intercept anything besides VMRUN (and optionally VMLOAD, VMSAVE, STGI, CLGI) -- all those instructions cause #GP when executed in PV guests (including Dom0)

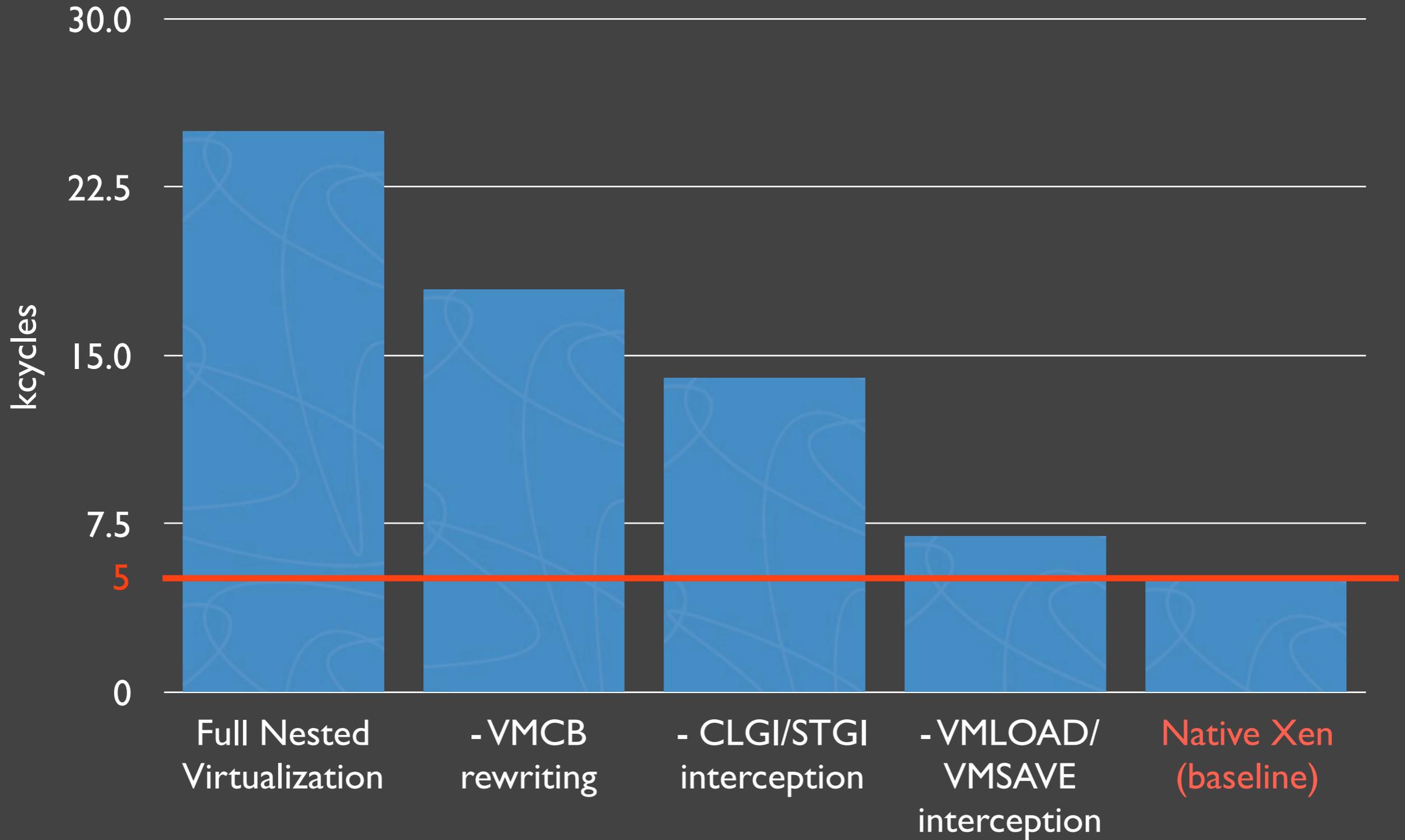
On AMD!

0

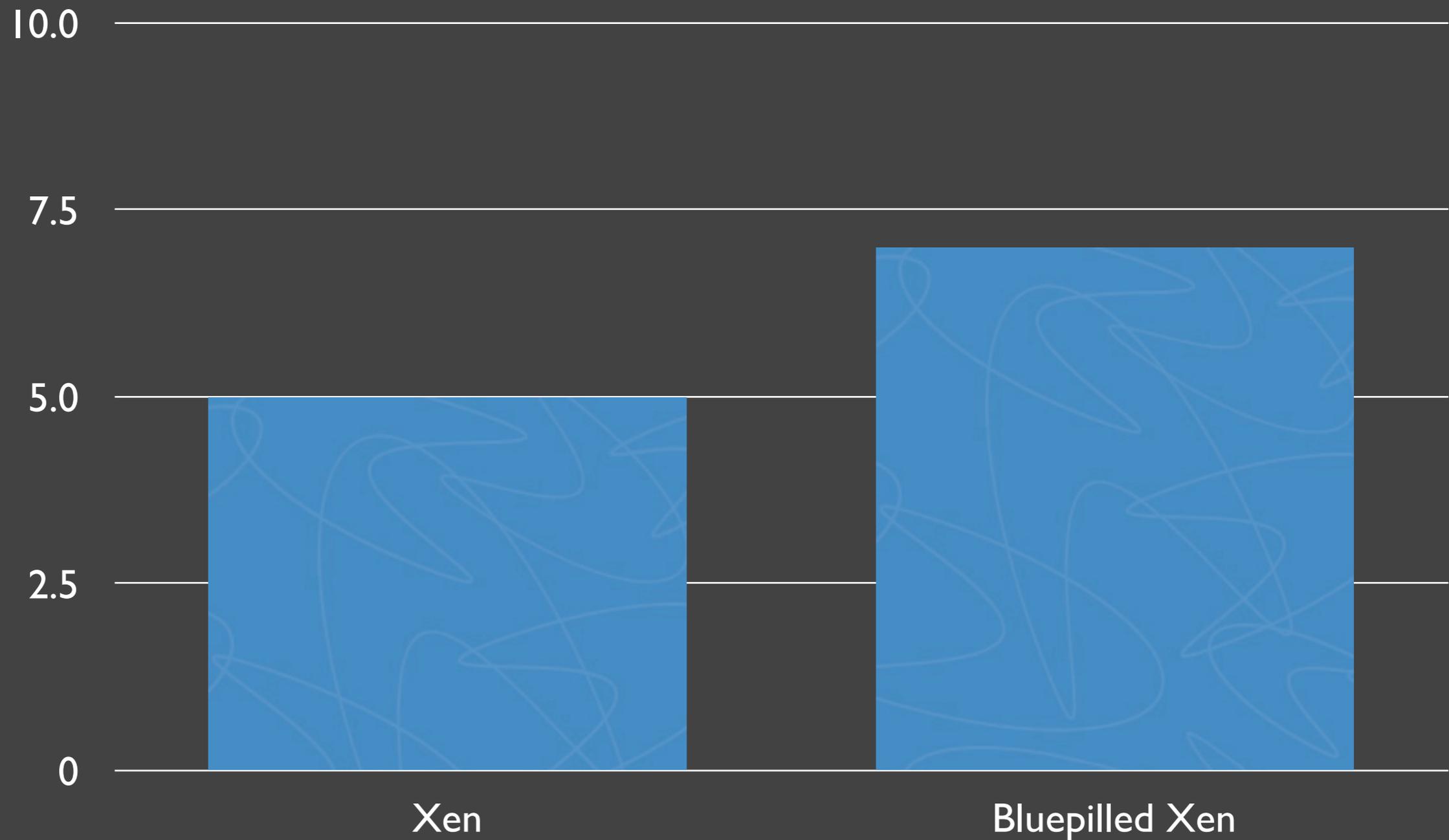
On Intel we have obligatory intercepts  
(CPUID, INVD, MOV CR3).

Impact on HVM domains

# HVM domains: impact on #vmexit time (RDMSR intercept on AMD)



# 5k cycles (Native Xen) vs. 7k cycles (Bluepilled Xen)



2000 cycles from the Holy Grail ;)



But that you can observe only in a HVM domain;  
on PV domains it is: 0 cycles (on AMD)!

HyperGuard vs. BluePill?

Summary  
(of the whole trilogy)

Talk #1 (Rafal)

Modifying Xen via DMA attacks

# “Xen Loadable Modules” Framework

Hypervisor Rootkits/Backdoors for Xen  
(don't confuse with virtualization-based rootkits!)

# Talk #2 (Joanna & Rafal)

DMA protections (IOMMU/VT-d) on recent Xens

Getting around VT-d Xen protection

(BONUS: on the fly SMM modification, despite D\_LCK set)

Other Xen protection mechanisms...

... and how they sometimes might be bypassed...

Exploiting a heap overflow in Xen hypervisor

HyperGuard - integrity scanner for a hypervisor

# Talk #3 (Alex & Joanna)

# Hardware Nested Virtualization

Blue Pill Boot

Xen Blue Pill: Bluepillling the Xen **on the fly**

Discussed the XBP detection

Slides available at:  
<http://invisiblethingslab.com/bh08>

Demos and code will be available from the same address after Intel releases the patch.

Thank you!