



Nested EPT to Make Nested VMX Faster

Red Hat

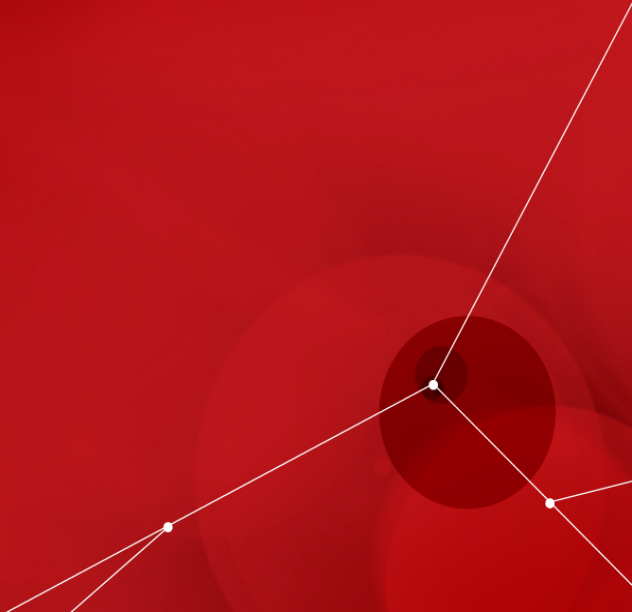
Author Gleb Natapov

October 21, 2013

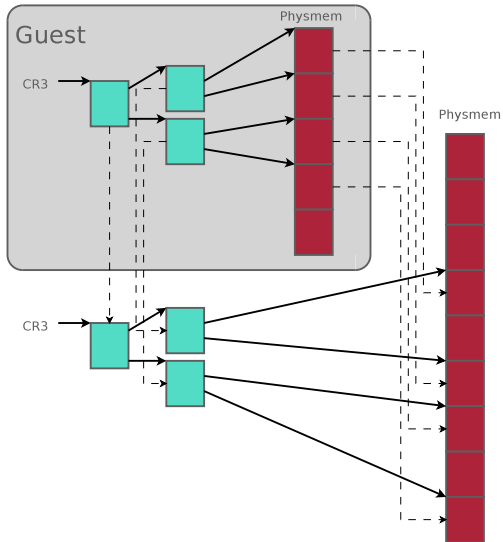


Section 1

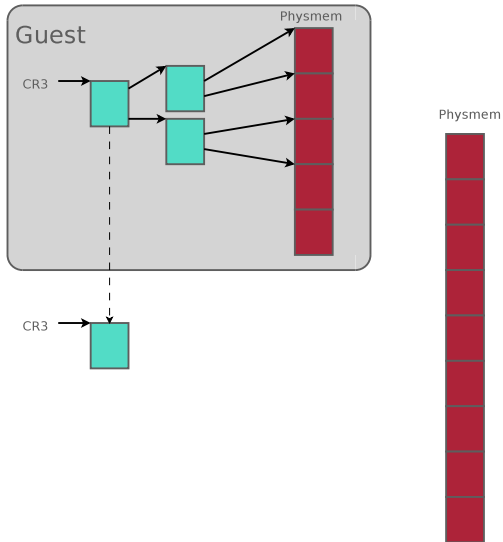
Background



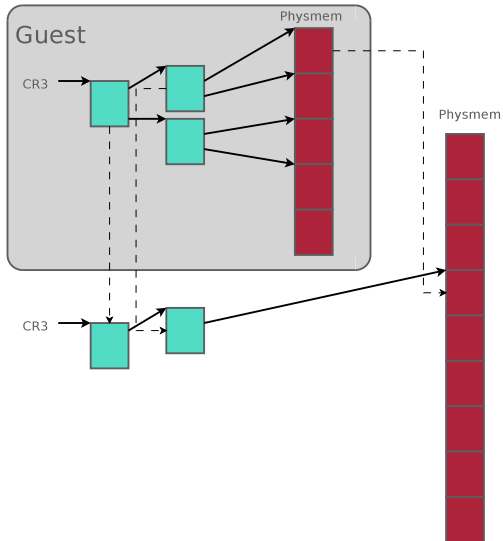
Shadow Paging



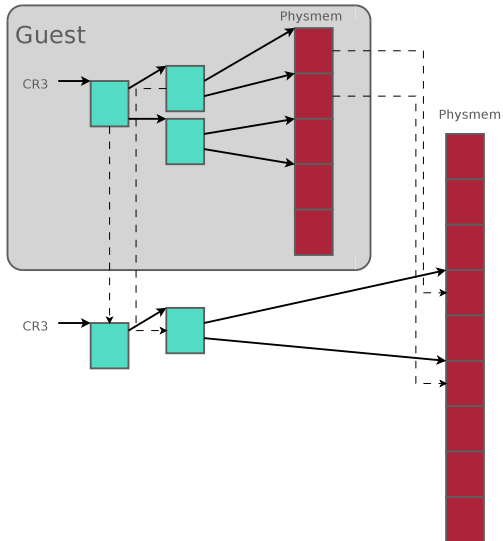
Shadow Paging



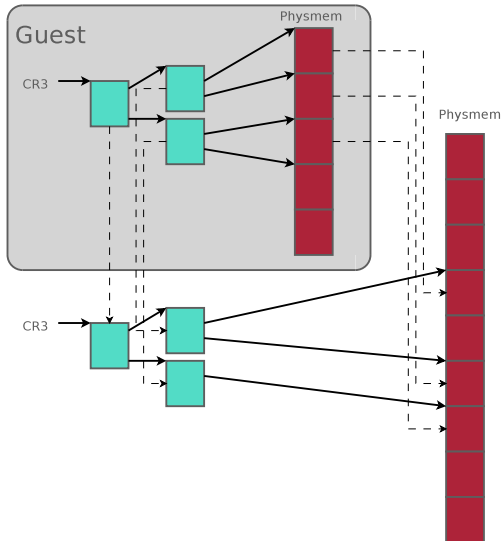
Shadow Paging



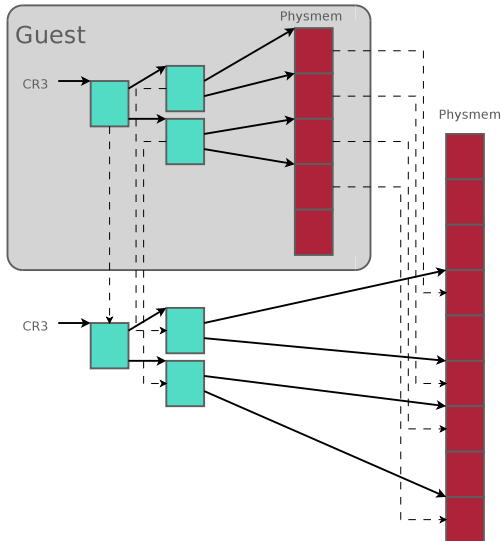
Shadow Paging



Shadow Paging



Shadow Paging



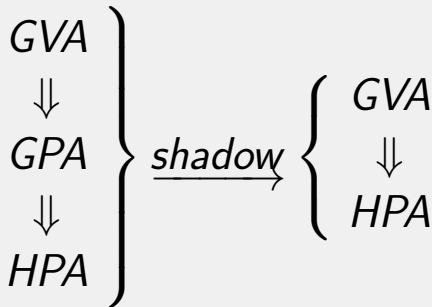
Shadow Paging (Cont.)

Slow!

- CR3 change traps to hypervisor
- Page table modification by a guest traps to hypervisor
- New address space creation (fork) requires new shadow page table to be created

Shadow Paging (Cont.)

What actually happens



EPT Saves the Day

Two level paging in HW so shadow is not needed!

GVA



Guest Page Table

GPA



Extended Page Table

HPA



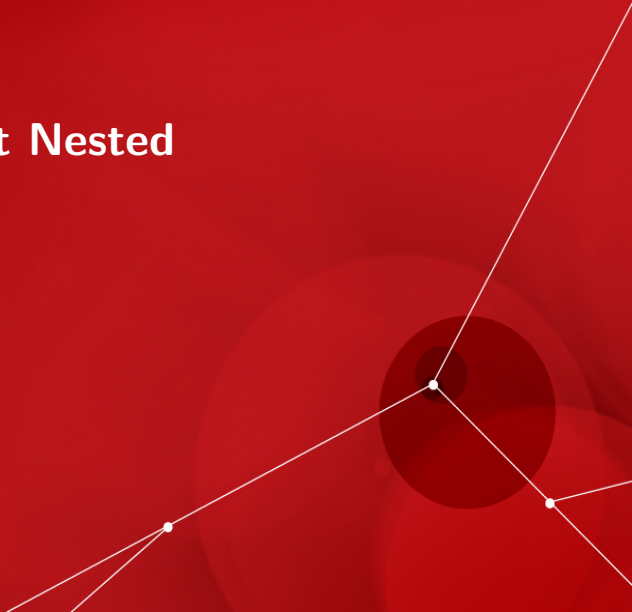
EPT Saves the Day (Cont.)

Guest manages its address space by itself



Section 2

What About Nested



Nested Guest is Running

Three levels of address translation!

nGVA



nGPA



GPA



HPA



Nested Guest is Running (Cont.)

But HW has only two levels!



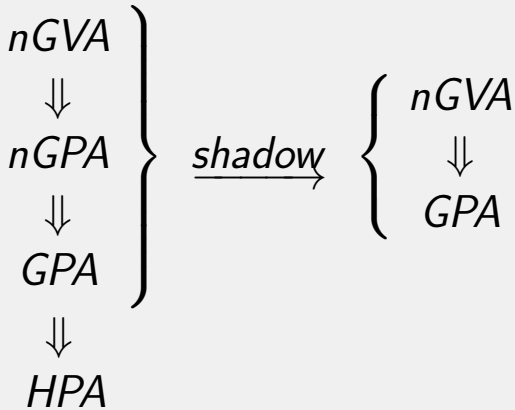
Nested Guest is Running (Cont.)

Something has to be shadowed



Shadow on EPT

What actually happens





Shadow on EPT (Cont.)

Slow for all the same reasons as regular shadowing

Plus each L2's #PF and CR3 access traps to L0
and forwarded to L1



Shadow on EPT (Cont.)

Slow for all the same reasons as regular shadowing
Plus each L2's #PF and CR3 access traps to L0
and forwarded to L1

Nested EPT

Key observation

Guests are created/destroyed much less frequently than processes

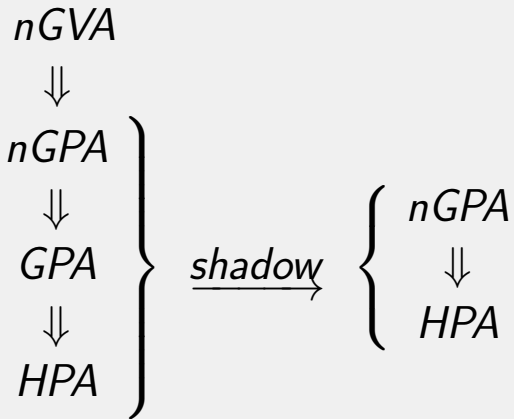


Nested EPT (Cont.)

Why not shadow nGPA to HPA translation instead?

Nested EPT (Cont.)

What actually happens





Nested EPT (Cont.)

Nested guest manages its address space by itself



Section 3

Implementation



Good

KVM already has shadow paging code

Good (Cont.)

KVM shadow code understands all guest's paging modes

- 32-bit Paging
- PAE Paging
- IA-32e Paging



32-bit Paging

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
Address of page directory ¹												Ignored						P C D	PW T	Ignored			CR3									
Bits 31:22 of address of 4MB page frame				Reserved (must be 0)		Bits 39:32 of address ²				P A T	Ignored	G	1	D	A	P C D	PW T	U / S	R / W	1	PDE: 4MB page											
Address of page table										Ignored	0	I g n	A	P C D	PW T	U / S	R / W	1	PDE: page table													
Ignored																	0	PDE: not present														
Address of 4KB page frame										Ignored	G	P A T	D	A	P C D	PW T	U / S	R / W	1	PTE: 4KB page												
Ignored																	0	PTE: not present														

PAE Paging

6	6	6	6	5	5	5	5	5	5	5	5	5	5	M ¹	M-1			3	3	3	3	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	0	9	8	7	6	5	4	3	2	1	0		
Ignored ²															Address of page-directory-pointer table															Ignored					CR3																					
Reserved ³															Address of page directory															Ign.	Rsvd.	P C W D T	P R s v d	1	PDPTE: present																					
Ignored																									0	PDPTE: not present																														
X D 4	Reserved															Address of 2MB page frame										Reserved					P A T	Ign.	G	1	D A	P C W D T	P U / S /	R /	1	PDE: 2MB page																
X D	Reserved															Address of page table															Ign.	0	I g n	P C W D T	P U /	R /	1	PDE: page table																		
Ignored																									0	PDE: not present																														
X D	Reserved															Address of 4KB page frame															Ign.	P G A T	D A	P C W D T	P U /	R /	1	PTE: 4KB page																		
Ignored																									0	PTE: not present																														



IA-32e Paging

6	5	4	3	2	1	0	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
Reserved ²											Address of PML4 table											Ignored											P	P	C	W	D	I	Ign.	CR3																														
X	D	3	Ignored						Rsvd.	Address of page-directory-pointer table											Ign.											R	s	v	d	I	g	n	A	P	P	C	W	D	I	U	S	/	R	/	W	1	PML4E: present																	
Ignored																													Q	PML4E: not present																																								
X	D	Ignored						Rsvd.	Address of 1GB page frame						Reserved						P	A	T	Ign.	G	1	D	A	P	P	C	W	D	I	U	S	/	R	/	W	1	PDPTPE: 1GB page																												
X	D	Ignored						Rsvd.	Address of page directory											Ign.											Q	g	n	A	P	P	C	W	D	I	U	S	/	R	/	W	1	PDPTPE: page directory																						
Ignored																													Q	PDPTPE: not present																																								
X	D	Ignored						Rsvd.	Address of 2MB page frame						Reserved						P	A	T	Ign.	G	1	D	A	P	P	C	W	D	I	U	S	/	R	/	W	1	PDE: 2MB page																												
X	D	Ignored						Rsvd.	Address of page table											Ign.											Q	g	n	A	P	P	C	W	D	I	U	S	/	R	/	W	1	PDE: page table																						
Ignored																													Q	PDE: not present																																								
X	D	Ignored						Rsvd.	Address of 4KB page frame											Ign.											G	A	T	D	A	P	P	C	W	D	I	U	S	/	R	/	W	1	PTE: 4KB page																					
Ignored																													Q	PTE: not present																																								

What is Common?

- bit 0 - Present
- bit 1 - R/W
- bit 2 - User
- bit 5 - Accessed
- bit 6 - Dirty
- bit 7 - Large Page
- bit 63 - Execute Disabled (PAE & IA-32e)

What is Different?

- PTE size (32bit vs 64bit)
- Number of page table levels

How Differences are Handled

- Shadow paging code is a template
- All differences are template parameters
- Template code is compiled for each paging mode
- `vcpu->mmu` is initialized according to current guest mode

Bad

EPT page table format is very different

Find the Differences

Bit	Regular Paging	EPT
0	present	readable
1	writable	writable
2	user	executable
5	accessed	memory type
6	dirty	ignore pat
7	large page	large page
8	ignored	accessed
9	ignored	dirty
63	XD	Suppress #VE

Step One: Make PTE handling parameterizable

- Reserved bits
- Present
- Dirty
- Accessed
- Permission

Step Two: Teaching Shadow About EPT

```
arch/x86/kvm/mmu.c          |    5 +++++
arch/x86/kvm/paging_tmpl.h  |   37 ++++++-----
2 files changed, 41 insertions(+), 1 deletion(-)
```

Step Three: Switch to Shadow EPT

On nested guest entry switch `vcpu->mmu` to EPT

But...

KVM uses `vcpu->mmu` for two purposes:

- 1 Virtualize guests memory
- 2 Translate GVA to GPA during instruction emulation

But... (Cont.)

What if L0 wants to emulate L2's instruction?

It needs to translate an address from nGVA to GPA

EPT vcpu->mmu translates from nGPA to GPA

But... (Cont.)

What if L0 wants to emulate L2's instruction?

It needs to translate an address from nGVA to GPA

EPT vcpu->mmu translates from nGVA to GPA

But... (Cont.)

What if L0 wants to emulate L2's instruction?

It needs to translate an address from nGVA to GPA

EPT vcpu->mmu translates from nGVA to GPA

Solution

Nested MMU

- Pointed to by `vcpu->nested_mmu`
- Translates nested guest's address twice:
 - 1 nGVA \rightarrow nGPA
 - 2 nGPA \rightarrow GPA

Numbers

Kernel compile

Shadow-on-EPT: 33m22s

Nested EPT: 9m46s



The end.

Thanks for listening.