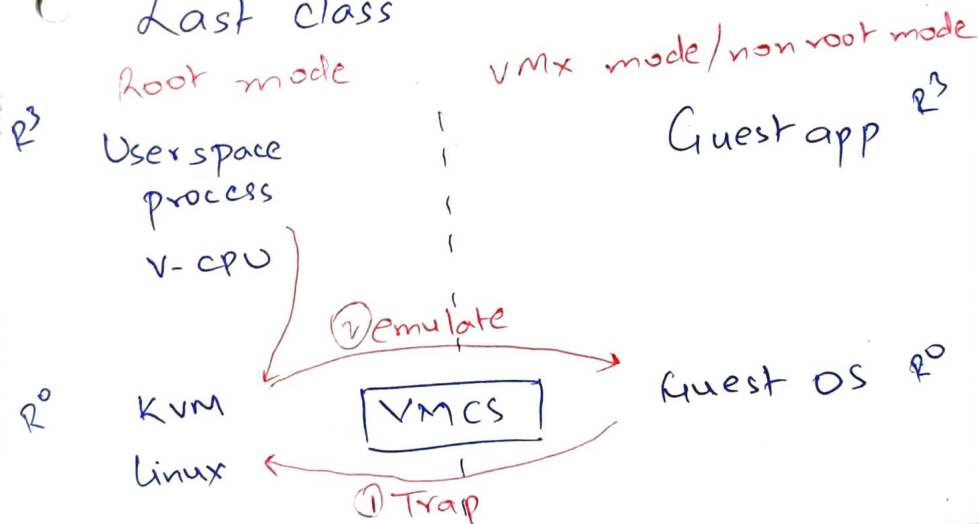


● Last class: Popek Goldberg Theorem (3)

Take control of resources ← Sensitive Instructions ⊂ Privileged Instruction
eg:- Change CR3
- Change IDT
- Disable Interrupts
↳ Trap

Problem with popf

● Last class (3)



● Great performance for CPU benchmarks (2)

- 0-9% slowdown on SPEC benchmark

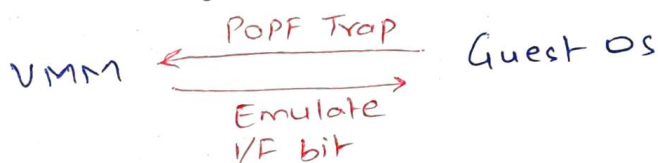
(Figure 2 ASPLOS-06. Comparison of sw hw tech for x86 virt)

- Housekeeping overhead Timer Interrupts

- Exits determine performance for trap-and-emulate hypervisors ^③

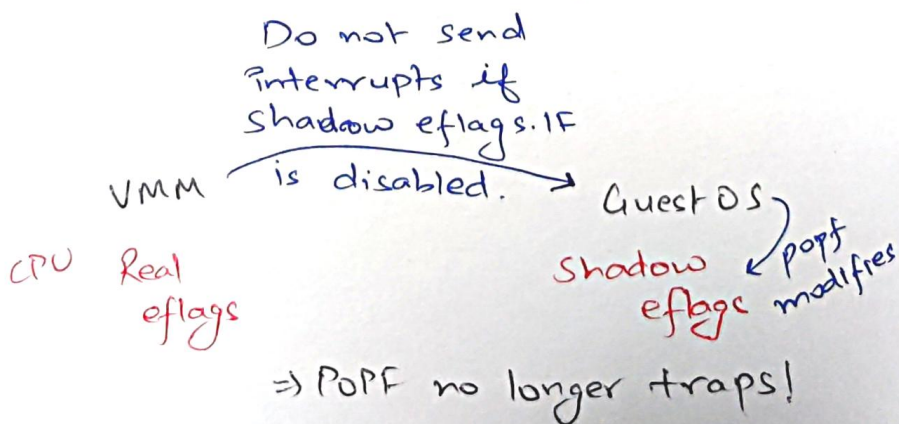
Micro Arch.	Launch date	VMM round trip cycles	Syscall round trip cycles
Nehalem	3Q09	1009	138
Sandy bridge	1Q11	784	134

- Example: guest running popf frequently ^①

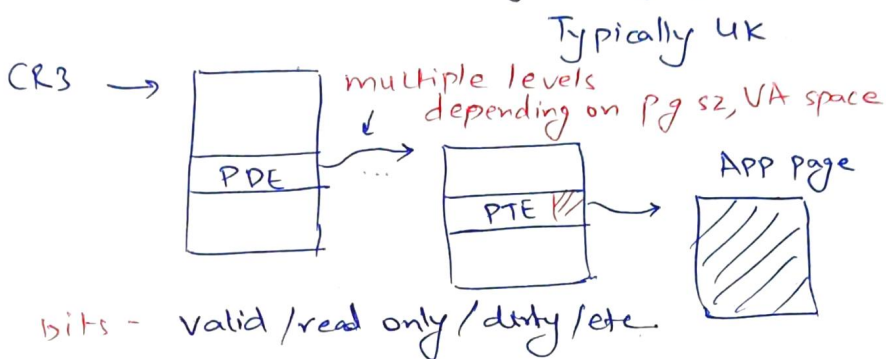


Solution: Avoid exits as much as possible!

- Dealing with popf: Shadow EFLAGS register in VMCB ^②



Virtualize memory! Paging (3)



Metrics of success

Transparency: Process thinks it owns memory

→ Guest OS, Guest process (2)

Isolation: One process can not read another's memory

→ Guest OS, Guest process

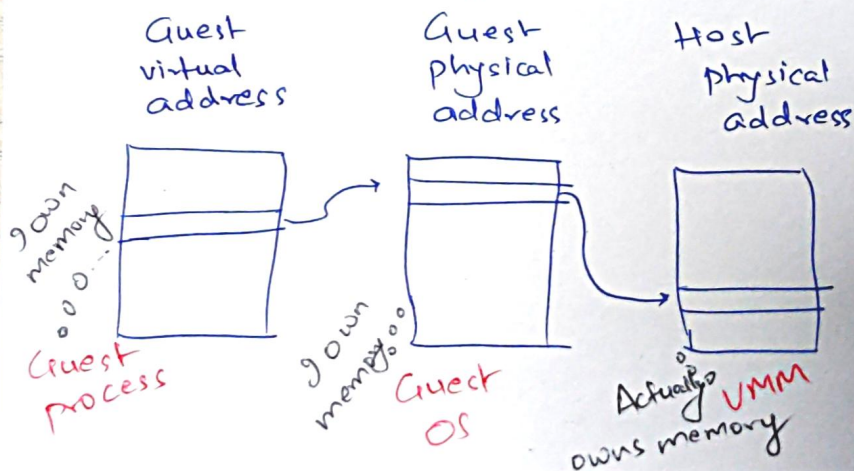
Performance: Once mapped, address translation need not involve OS

↑ Guest OS, VMM

Fast load/stores

General idea! (3)

Two levels of address translation



● Idea 1: Interpretive Execution (Lab) ②

30xx → mem[xx] = AC

array owned by Interpreter
Virtual address

✓ Isolation

✓ Transparency

✗ Performance

eg: IBM System 370
SIE instruction → Start Interpretive Exec

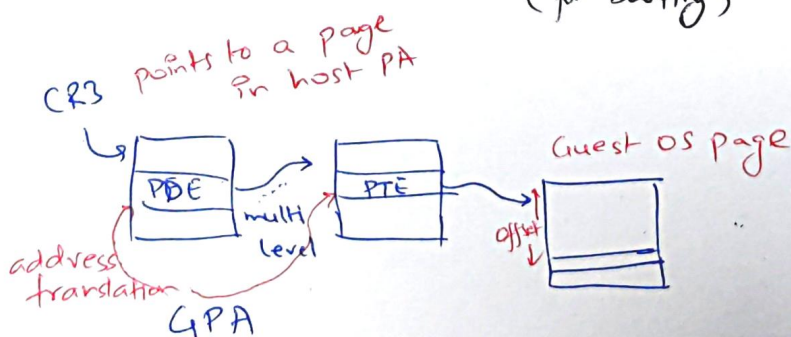
● Idea 2: Reuse paging hardware ②

Shadow Page table:

map directly from GVA → HPA

⇒ Most memory ops LD/ST run at native speed

● Step 1: Prepare page table for GPA (for booting) ②



Now guest OS can boot. (3)

- Initially, address translation is disabled on real hardware

• CRO, bit 31 can enable-disable paging.

Guest OS tries to disable paging $\xrightarrow{\text{Trap}}$ VMM changes CR3 to the prev. page table
 $\xleftarrow{\text{emulate}}$

Guest "thinks" it has disabled paging. (3)

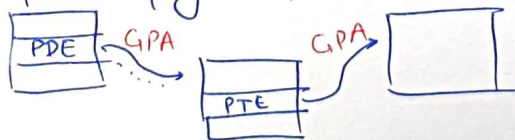
Reading CRO. 31 shall { trap or return false } shadow register

But address translation is active.

Addresses treated as GPA

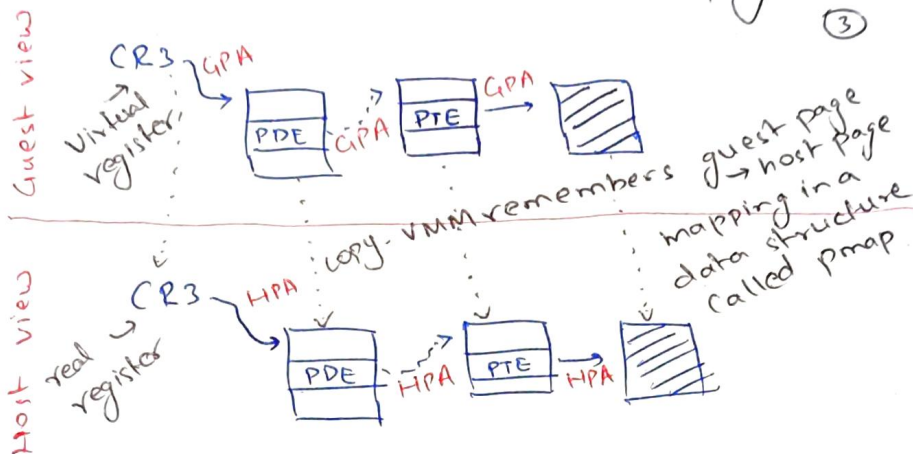
Guest tries to run a process (2)

① Prepares page table

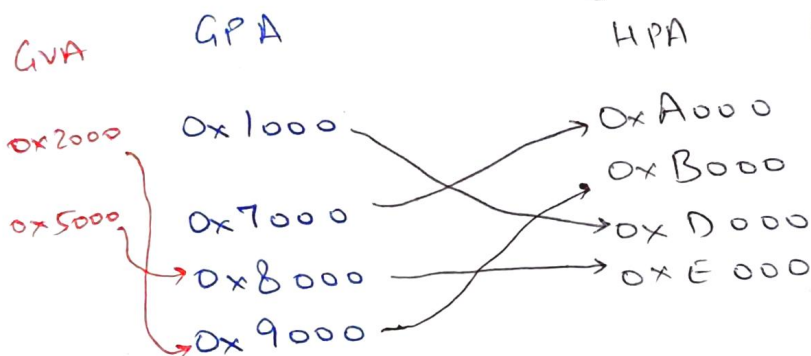


② Changes CR3 to top-level GPA

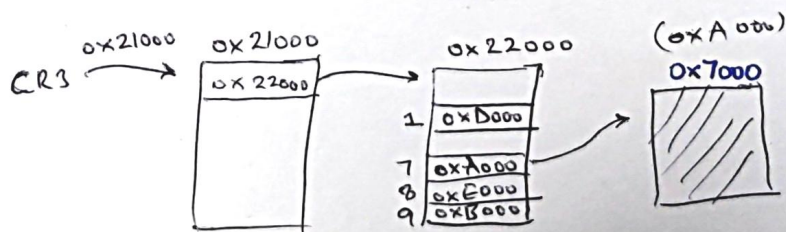
VMM constructs shadow page table ③



Example: memory mapping setup for vm ②

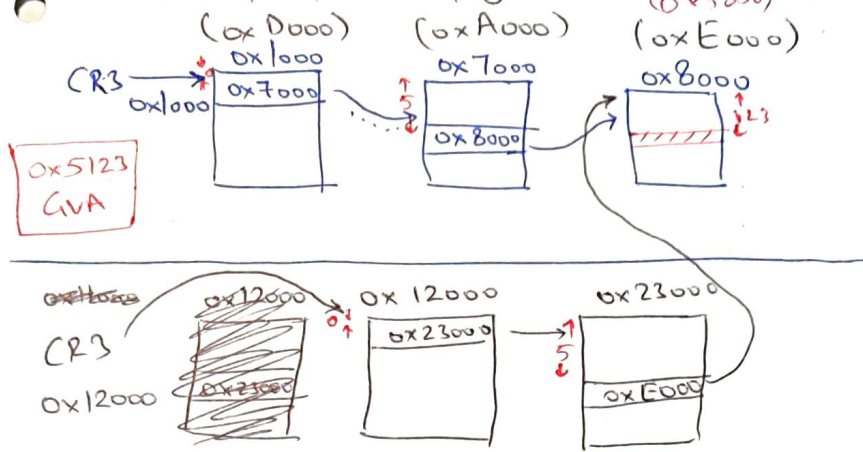


Example: memory setup for vm



When CR0.31 bit is disabled by guest, Hypervisor converts GPA → HPA

Guest prepares page table (3)



○ Transparency: If guest reads CR3, it
is emulated
or
kept as shadow CR3

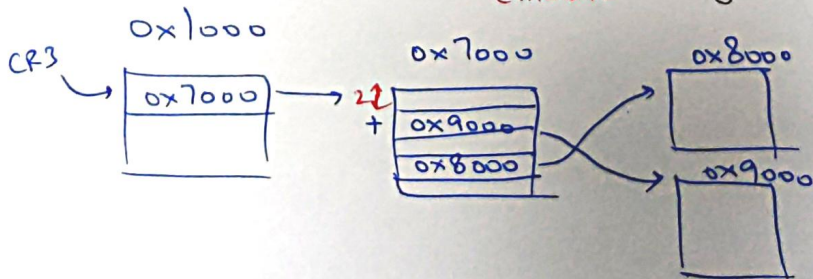
Performance: Directly use MMU once page table is setup

Safety/ Isolation: Guest can't directly touch shadow page table. 0x12000, 0x23000 are **NOT** in guest's address space

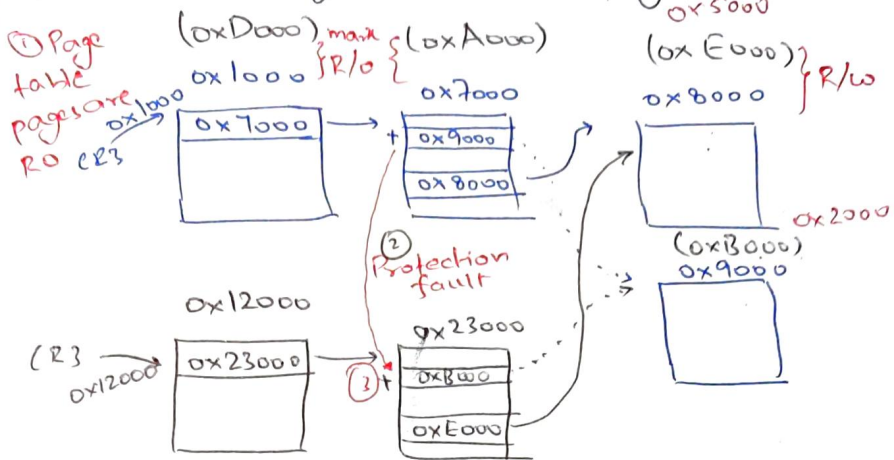
○ ~~Guest~~ Guest app mmaped at GVA. (2)

Guest OS did not allocate (lazy)

LD 0x2123 → PF in VMM → PF in guest OS
emulate



● Maintaining shadow page table ③



● Memory tracing - ②

- Trace writes - mark R/w
- Trace read/writes - mark invalid

Another use - memory mapped I/O

● → Memory / performance optimizations ②

Need 1 shadow page table per guest page table.

① VMM deletes a shadow page table

OK - Rebuild when CR3 is updated back to 0x1000.

- Hidden Page fault \rightarrow Transparent to guest
- ② Move 0xE000 page to disk. Set "Not Present" in shadow page table. ③

App runs LD 0x5123

\rightarrow Page fault

Bring from disk to 0xF000.

Update ALL shadow page tables pointing to 0xE000.

\rightarrow Maintain backward mappings in "pmap"

- ③ Move 0xA000 to disk. When ② guest OS tries to update page table.
 \rightarrow Page fault

Bring to 0xC000

Update GPA \rightarrow ~~GPA~~ HPA mapping
0x7000 \rightarrow 0xC000

Summary: ②

- ① Trap on mv CR3
- ② Copy PT to shadow PT
- ③ Give shadow PT directly to h/w for address translation
- ④ Memory tracing for maintaining consistency
- ⑤ Pagefault, hidden page faults \rightarrow backward mappings

○ Transparency ✓ Isolation ✓ ②

- Performance: Native LD/ST in common case ✓
 - Overhead on process creation
 - " " new PTE
 - " " page faults
 - Additional overhead - hidden page fault
- Resource utilization: maintain copy of page table for every guest virtual address space