

Shadow page table overheads ③

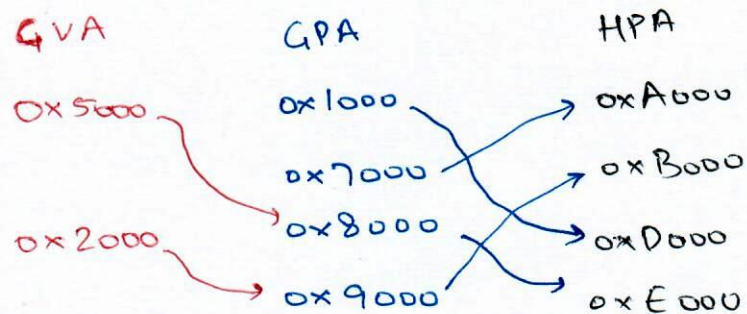
- 1 fault for every PTE R/w → R/o
- Changing CR3 to child
 - : Copy page table to shadow page table
- At a write, protection fault at VMM → forwarded to guest → guest copies page, updates PTE → again a fault → hypervisor updates shadow page table.

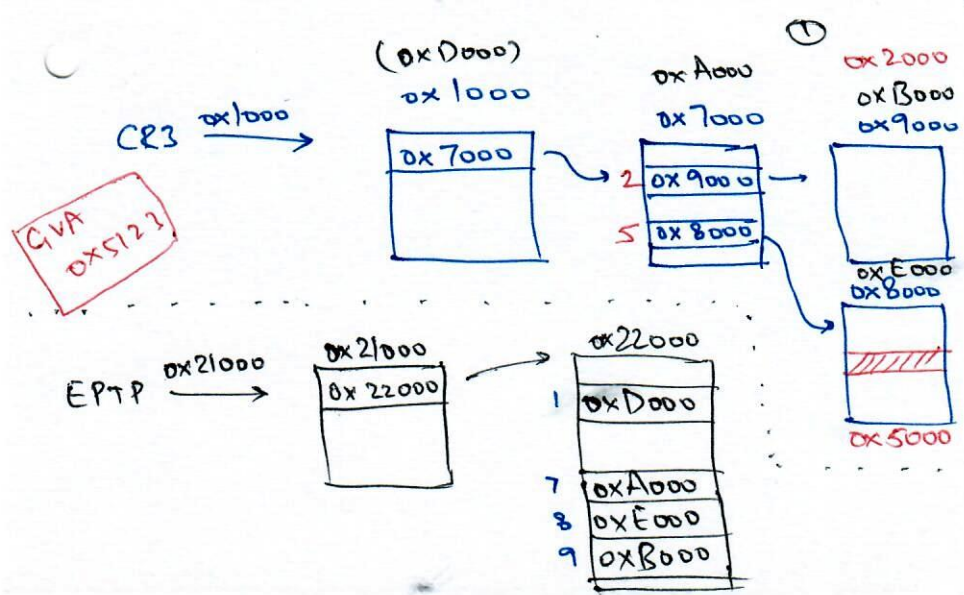
Overheads:

Pentium 4 672 ②

- Fork a proc, wait for child to end, fork again → 4.4x slowdown
- PTE modification
 - native → single cycle store
 - Shadow PT → 12733 cycles

Extended Page Tables ①





Page walk GVA 0x5123 ③

$0x1000? \rightarrow 0x21000 \rightarrow 0x22000 \rightarrow 0xD000$
 $0x7000? \rightarrow 0x21000 \rightarrow 0x22000 \rightarrow 0xA000$
 $0x8000? \rightarrow 0x21000 \rightarrow 0x22000 \rightarrow 0xE000$

8 memory references

m level from GPA \rightarrow HPA

n level for GVA \rightarrow GPA

$0x5 \rightarrow 0xE$
cached in TLB

$\Rightarrow m+n+m \cdot n$ references

- Hardware support for 2-D address translations ②

- Changing CR3 need not trap \Rightarrow fast ctx sw

- Modifying PTEs need not trap \Rightarrow fast fork

④
Table 1:- Instruction and data translations

3.9x - 4.57x slowdown

Figure 3:- Percentage of unique page entries for each 2D page walk reference
most PTEs are common on higher levels

Figure 6:- 86-93% of native performance

Accelerating 2D page walks for Paper
(ASPLOS08) virtualized systems

Two main ideas: 1) Nested TLB ②

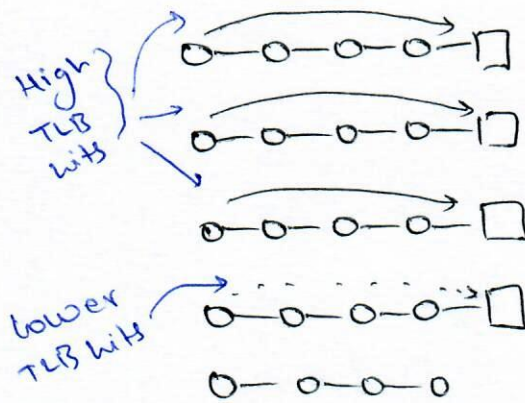


Figure 8

+7.3% performance improvement end-to-end

2) Large pages ②

- \rightarrow Increases TLB reach \rightarrow Also helps in shadow page tables
 - \rightarrow Reduces m and n
 - \Rightarrow Less number of references in 2D page walk
- Reduces PFs

2MB nested pages give 21-43% fewer TLB misses than 4K nested pages

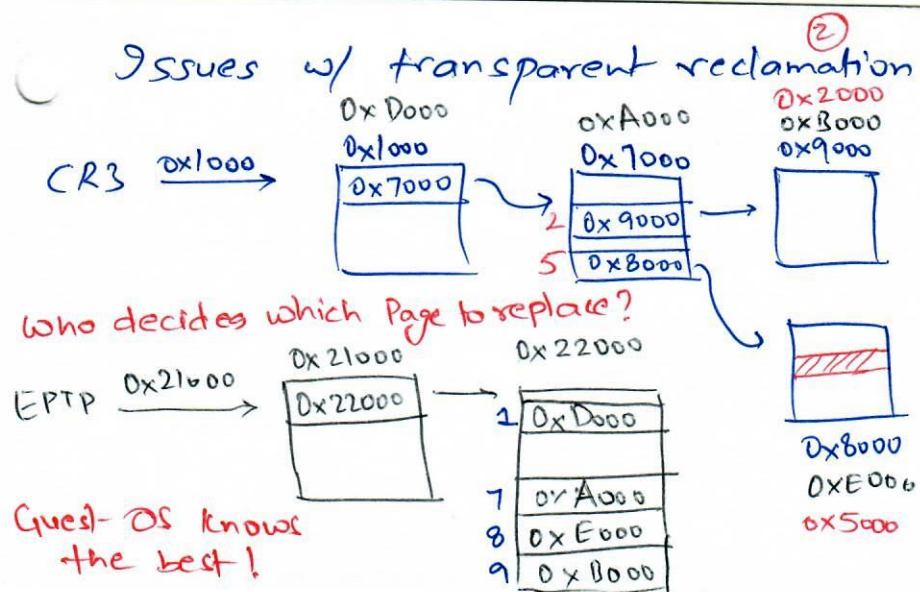
Memory management

Reclamation

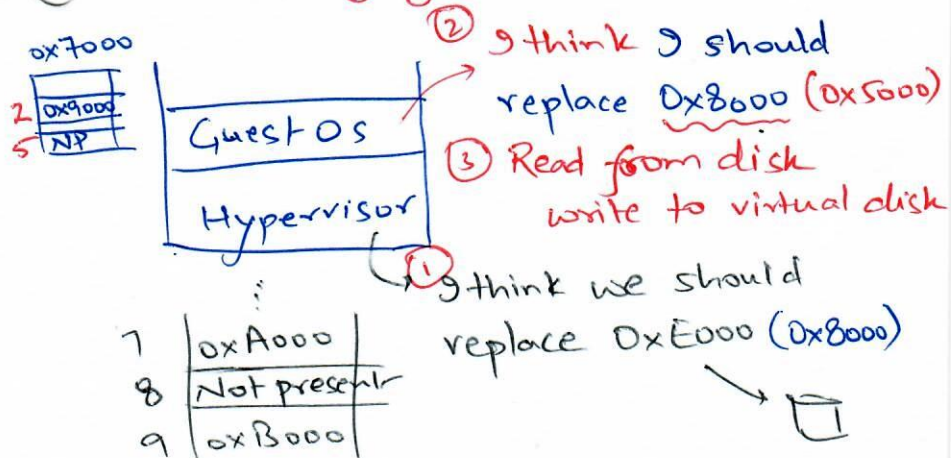
- Ballooning

Sharing memory

- Content based page sharing

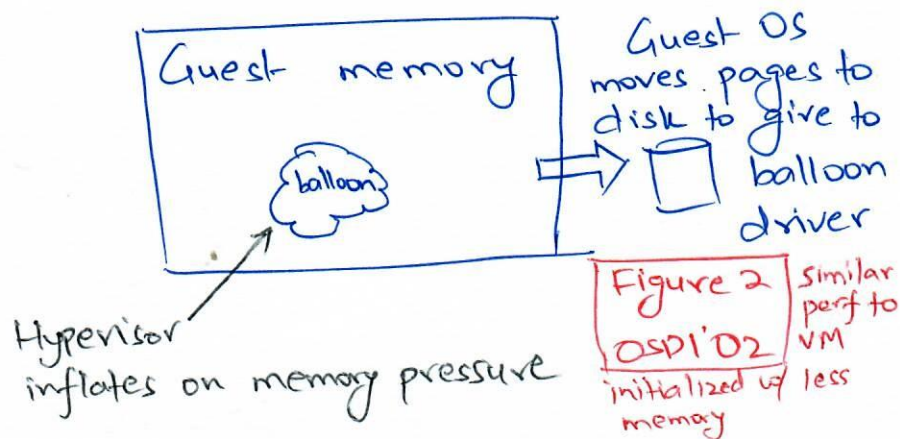


Double paging problem



Memory ballooning

(2)



Memory ballooning

(2)

- Gives control to Guest OS
- Effective in taking memory

If Guest OS refuses/disables ballooning driver

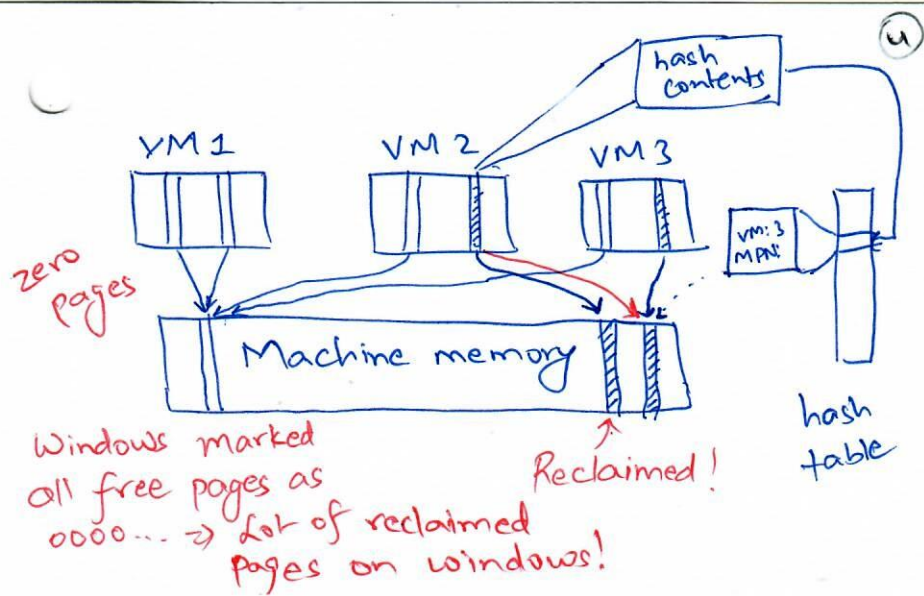
⇒ hypervisor forcibly takes pages
Guest OS may suffer from double paging

Content-based page sharing (4)

- Scan a random page periodically
- Hash ↳ opt: R/O page marked by OS
- Match ⇒ fully match page
- Match ⇒ Mark low. Reclaim

⇒ Transparent to Guests!

Figure 5 {OSDI'02} 7-32% reclaimed memory



Interesting trivia

①

Linux KSM (kernel Samepage merging)
/ " shared memory

- Similar hash-based implementation [Nov '09]
- VMware Patent!
- Red-black tree based implementation [Apr '09]