

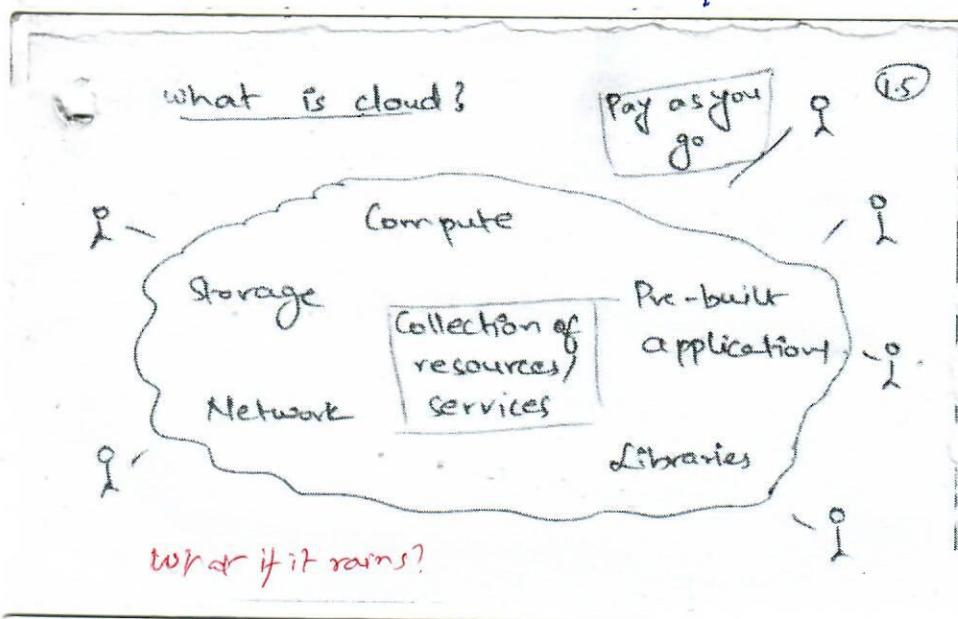
why should I take this course?

①

why should I care about cloud?

virtualization?

rust?



why cloud?

①

Don't want to manage resources

② Hardware faults are common

Assume your laptop crashes after 5 years

$\approx 1800$  days

If you have 200 of them, one is crashing 9 days

why cloud? Elasticity 2m

memes app

- My blog today got very popular

100 → 100,000 users



Laptop crashed. my chance of becoming popular is missed.

- Solution: use 1000 laptops. High upfront cost

why cloud? Provides fault tolerance 3m

- Doesn't lose data even when individual disks might fail
- Moves your computation to other machine if CPU overheats and dies.

→ AWS SLA

- Downtime is expensive > \$1M/hour ①

- Crashed server
- Disconnected wifi/electricity
- Tornado/Earthquake
- Loss of trust
- Bank lost last 100 transactions

## ① Why cloud? Proximity

1m

My memes app is getting popular  
in Canada.

- Ping time  $\approx 300$  ms
  - Every 100 ms delay drops 7% users
  - Rent a VM next door to user  $\approx 20$  ms
- Region-wise ping times

## ② Virtualization

What?  
- Act of virtualizing physical resources

Why?  
relevance to cloud

Challenges?

①

③

Investment



20 yrs sell

- Management, security overhead

Buy virtual gold: Gold bonds

GOLD 10g
Govt INDIA

① Why? (Buyer) back to VMs ②

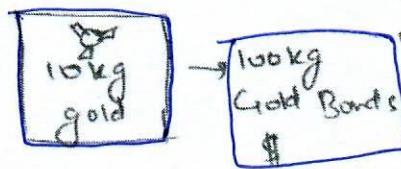
- No management
- Little upfront cost
- Elasticity
- Trusted



- DDoS attack
- Power outage
- Earthquake
- SW update

① Why? (Seller) ②

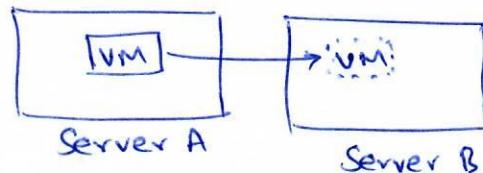
- Overprovision



- Anyways buying physical resources for self
- Amazon, Google
- Rent unused.
- Incremental migrant overhead

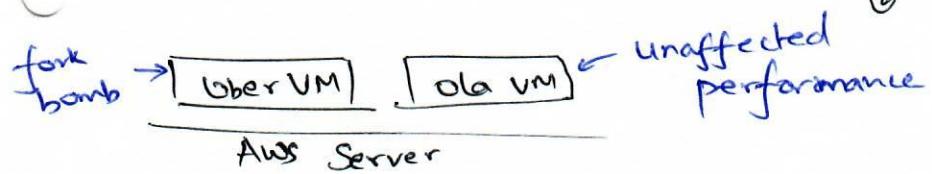
① Virtualization supports cloud (1) ②

- Live migration / checkpoint-restore



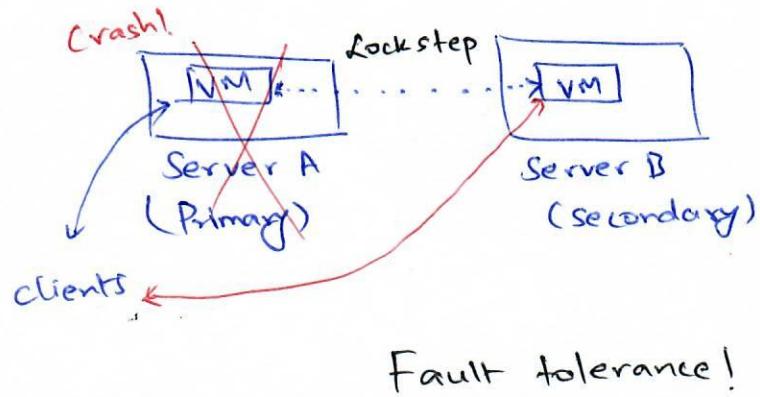
- Load balancing
- Management: restart/upgrade Server A (supercomputers)
- Proximity

- \* Strong Isolation (even performance) ②

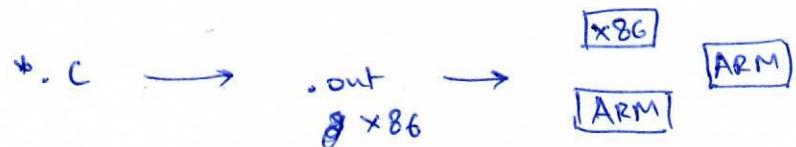


- \* Auditing/monitoring / Intrusion analysis
  - What all did the VM do?

### State machine replication

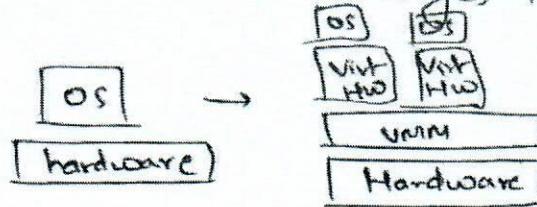


- \* Binary for one architecture can run on another ②



- Diversified placement
- Support legacy software

## Virtualization challenges: Feasibility ②



- zero/minimal modification to OS, apps  
(OS thinks it owns HW)
- Performance

## What's new? ②

Traditional IaaS - Rent VMs

- 1 How many VM?

Dropping requests  
Too few ←



Idle ⇒ losing money  
→ Too many

Sophisticated capacity planning, elasticity

## New cloud. ①

FaaS: specify functions, triggers

: Auto scaling

: **Do not pay** when function is not running > AWS Pricing <

New challenges: fast boot, dynamic provisioning, sharing data, etc

## Pricing comparison

(3)

- \* FaaS. Mobile application backend

\$ 2.33

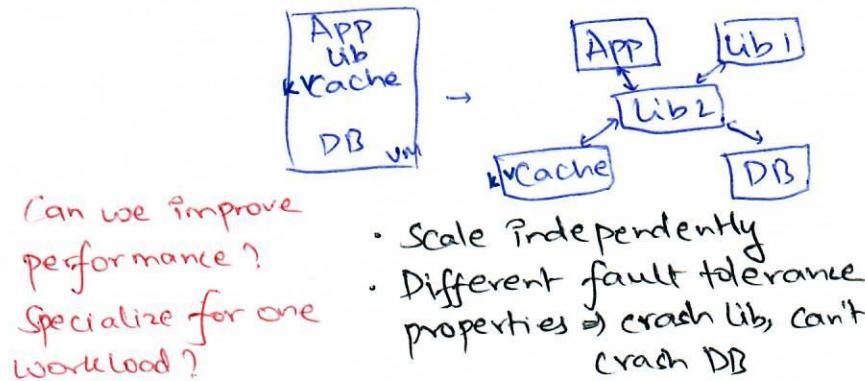
- \* IaaS. Al.medium  $\rightarrow 0.0255/\text{hr} \times 24 \times 30$   
 $= \$ 18.36$

- \* Ignoring data transfer/storage costs

## What's new (2)

(2)

Monolithic  $\rightarrow$  Microservices



## Why take this course?

(1.5)

### 1. You'll work with cloud

- Startup: using cloud infra  
(personal projects)

- AWS / Google / Azure: Managing

- ML experiments

- Research: interesting problems

2. Hands on

virtualization: beautiful

3. Study real systems  $\rightarrow$  firecracker history

## Highly collaborative project ①

- Collaboration is hard
- ⇒ CI/CD, PRs, unit tests
- Design is hard  $\Rightarrow$  interfaces/abstraction  
but fun
- Real end-to-end product

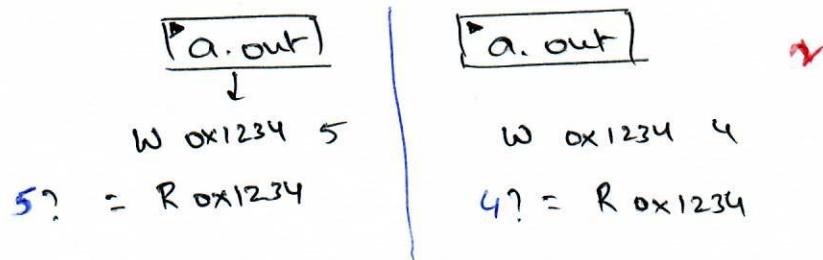
## why rust? ②

- Fastest growing language "most-loved" according to stack overflow
- Fast. Systems language C/C++
- Automatic memory management Java/Python/C#
- Modern paradigms: Iterators, closures without GC
- Safe: no null pointers, no double mutation,  
no pointer arithmetic

# OS background: Process virtualization

(2)

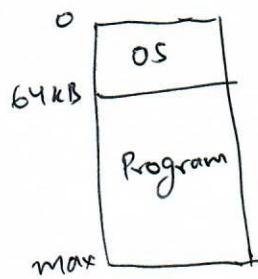
main.c → a.out



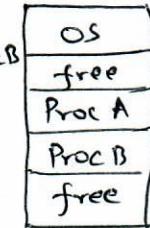
Memory of Processes are completely isolated

## Early systems

(3)



- Time sharing / multiprogramming
- expensive computers
- efficiency

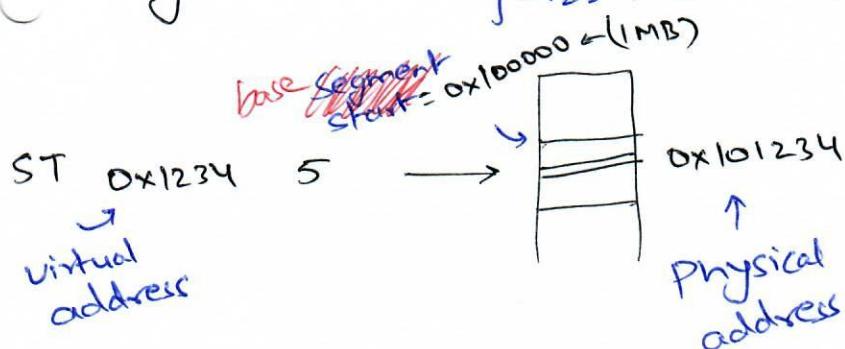


Physical memory

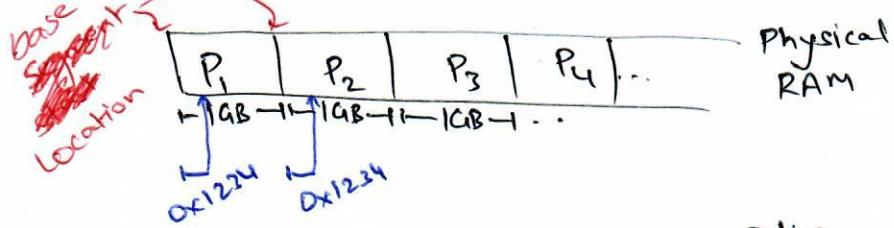
~~Base and bound~~

~~Segmentation: Performance~~

(2)

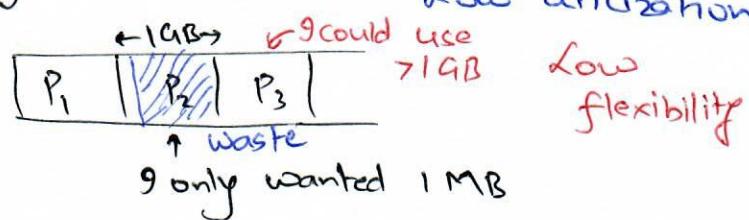


## Memory virtualization; ~~segmentation~~ ③



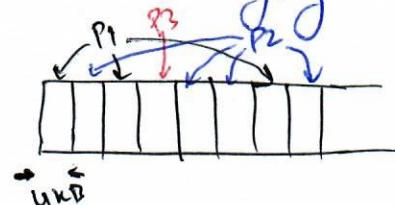
- + Virtualization: Process thinks it owns ~~the~~ entire memory
- + Performance: Just adder
- + Isolation/Security - P<sub>2</sub> cannot see P<sub>1</sub>'s memory
- Utilization

## ~~Segmentation~~ Problems - Low utilization ③



?: I can only start 4 applications?

## Solution: Paging ③



- + Flexibility: Dynamic mapping
- Allocate a page only when required
- Send a page to disk if not used (demand paging)

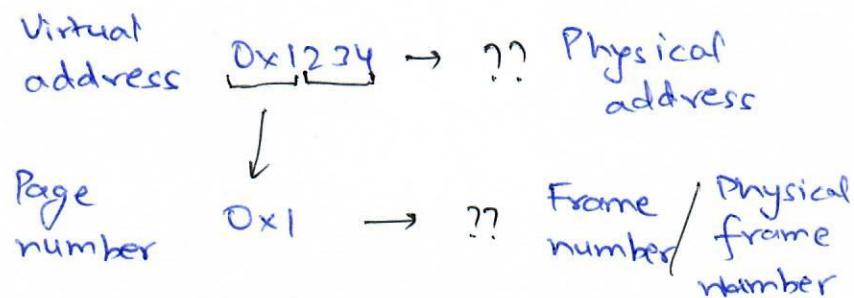
+ Dynamic mapping  $\Rightarrow$  Over provisioning ①

③: can start many more processes

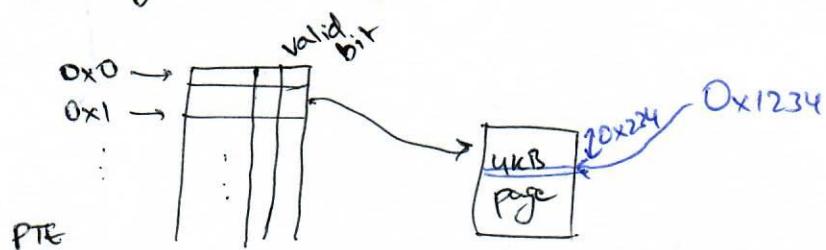
+ over-provisioning  $\Rightarrow$  better resource utilization

less free/unused memory

- flexibility  $\Rightarrow$  complex address translation ②



Page table: Page#  $\rightarrow$  Frame# ③



Read-only bit: Don't accidentally write code section

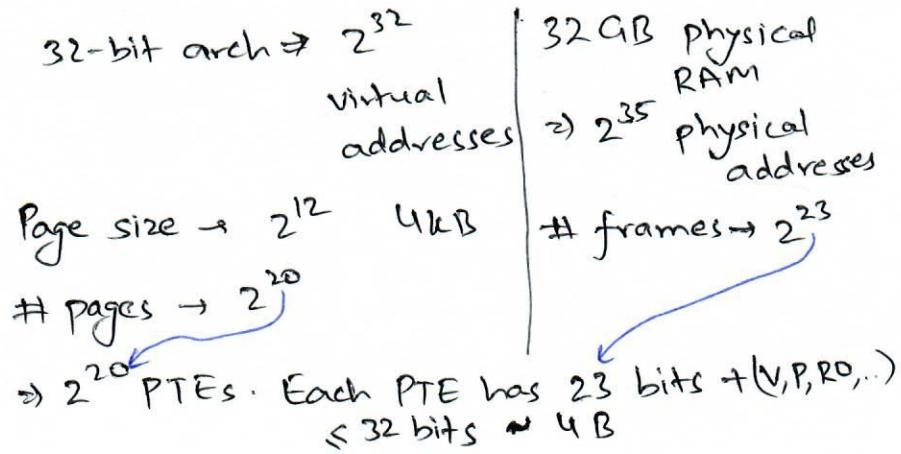
Dirty bit: Was the page written to

Present bit: Allow sending to disk

(Over-provisioning)

Page table is huge :-

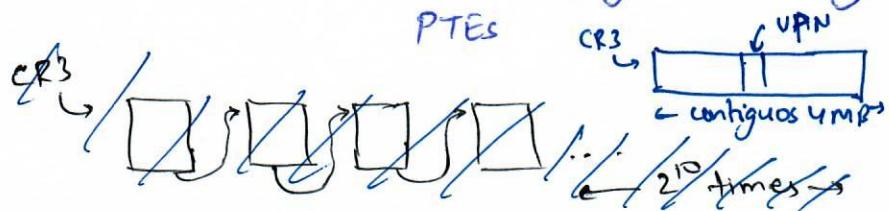
⑤



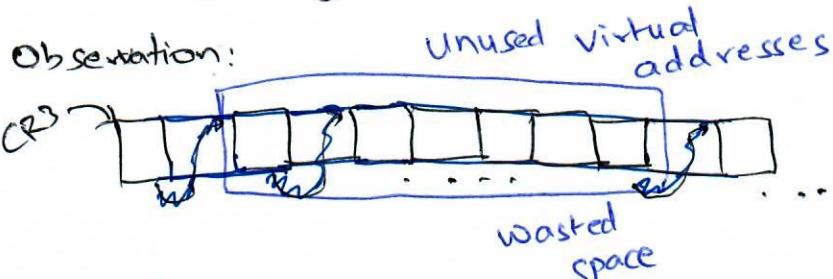
⇒ One 4KB page can have ③  
 $(2^{12})$   $2^{12}/4 = 2^{10}$  PTEs

Total PTEs required  $\rightarrow 2^{20}$

Naive solution: Use  $2^{10}$  pages containing

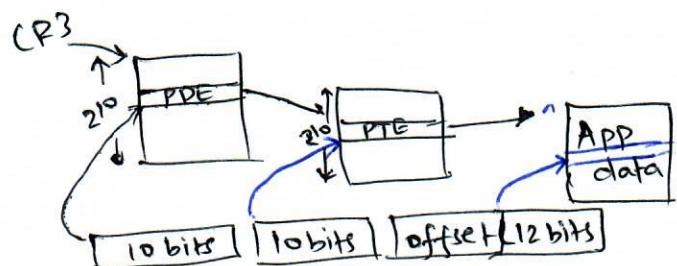


⇒  $2^{10} \times 4\text{MB} = 4\text{MB}$  of pages for holding page tables of 1 process! ③



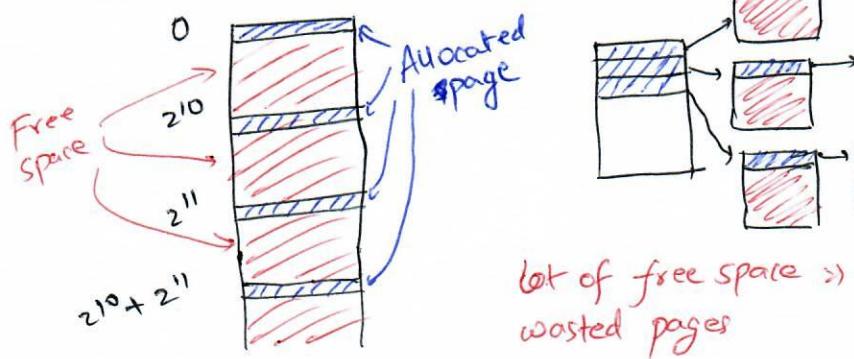
Some observation as segmentation

Some solution: Lazily allocate page to page table  
Hierarchical page table

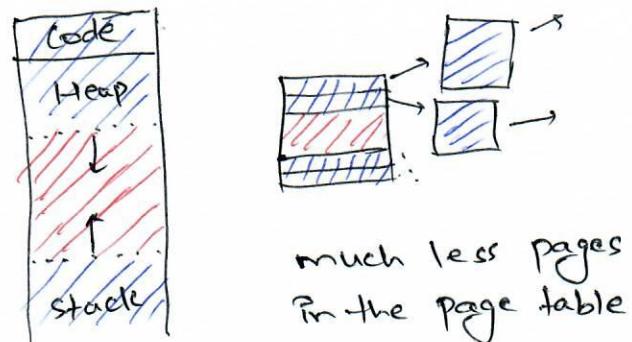


Effect on memory layout

Bad layout:



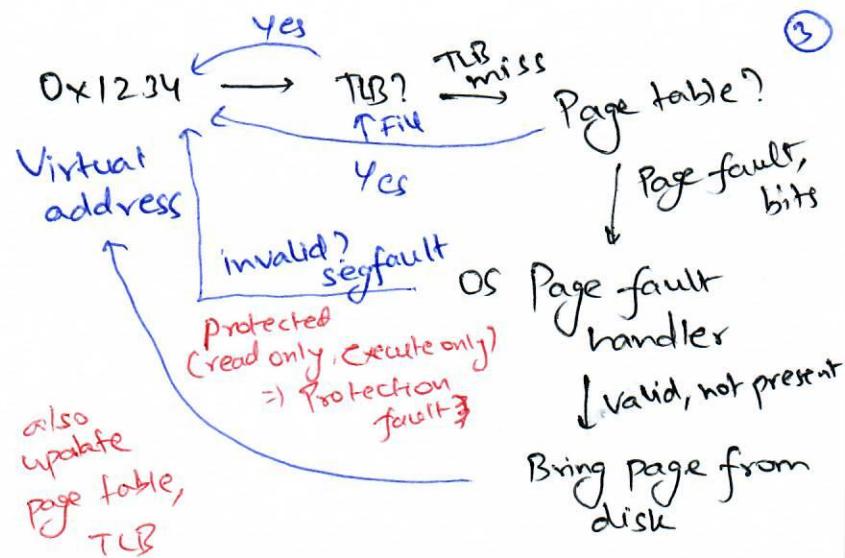
Good layout: segments



Problem: Every load/store needs 3 ②  
memory accesses 2 page table + 1 actual  
for 64 bit architectures

$$3 \rightarrow 6$$

Soln: TLB remember page # → frame #  
mapping



Memory virtualization ②  
Paging vs static segmentation

- + Each process thinks its own entire (virtual) address space → **Transparently move pages to disk**
- ~ Security
- + Overprovisioning ⇒ better resource utilization
- lot more complexity in hardware/os

## >Last lecture: Memory virtualization

②

Transparency: I think I own memory

Base & bound:

addr translation → adder  
+ check bounds

✓ Perf

✗ Flexibility

✓ Hardware complexity

✗ Utilization

③



## Paging

②

- Bits: valid, present, dirty, r/o

- Page Table

• Large ⇒ Hierarchical

• TLB - spatial locality for TLB hits

• Demand paging

Perf ✓

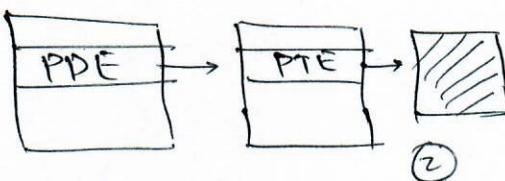
Flexibility ✓

HwComplexity ✗

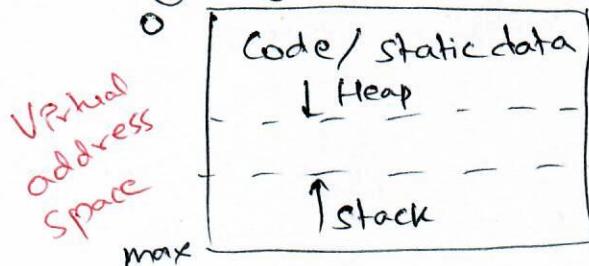
Resource util ✓

## HPT

CR3 →



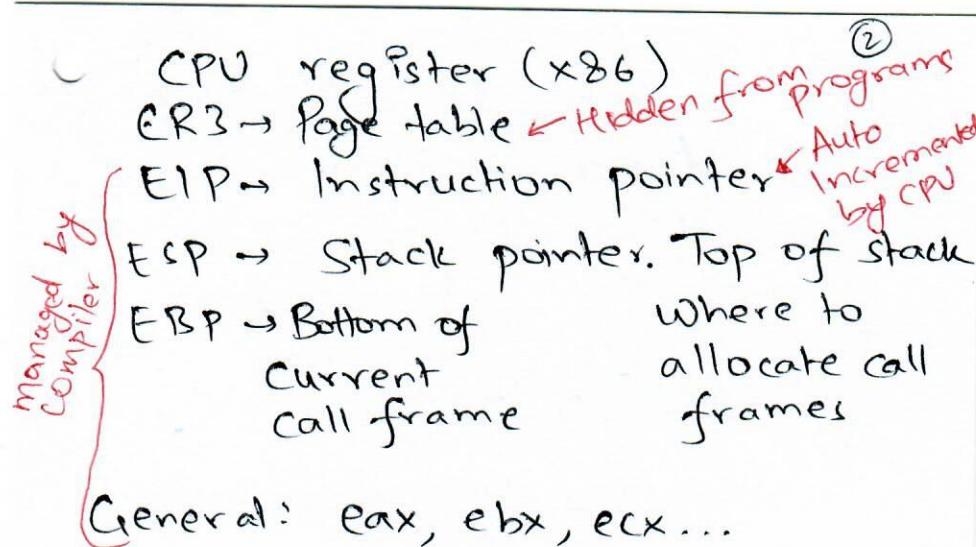
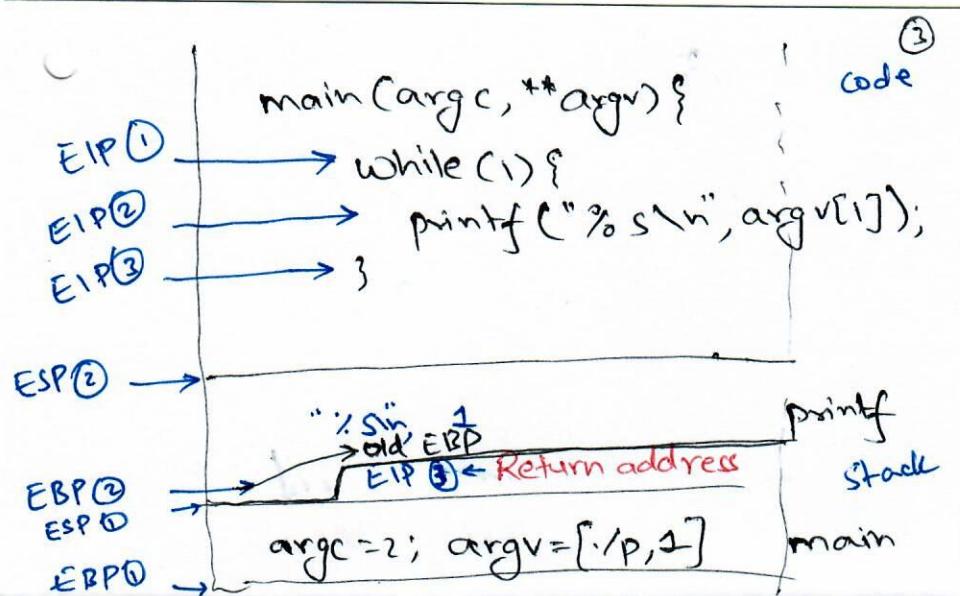
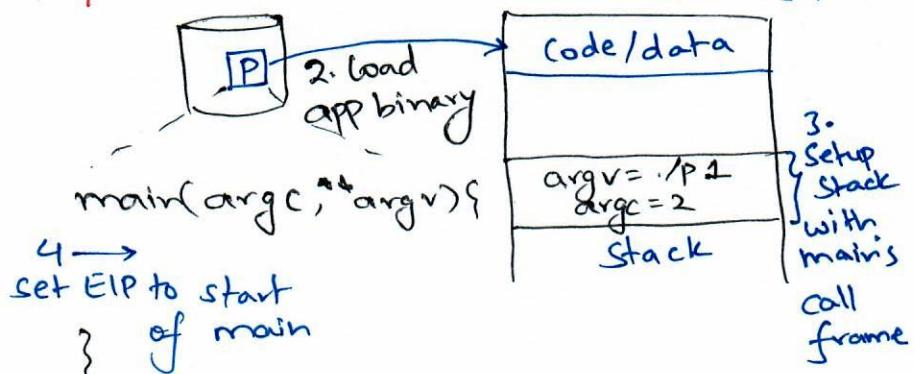
## Memory layout

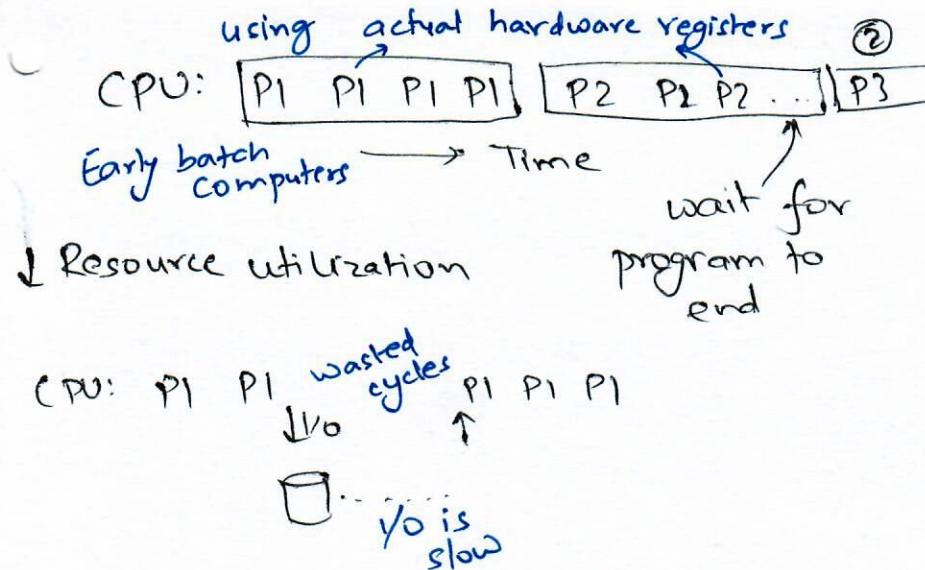


# How does a program run? ③

\$ ./p1

1. New VA is allocated





### Operator (manual): ①

- Kill **run away process**
- Prioritize process
  - manage list
- Only used during office hours
- Only way to intervene is to **kill** a process

P: main(argc, \*\*argv) { ① }

modifies actual hardware registers      while(1)  
     printf ("%s\n", argv[1]);

\$ ./P A &; ./P B &; ./P C & ...

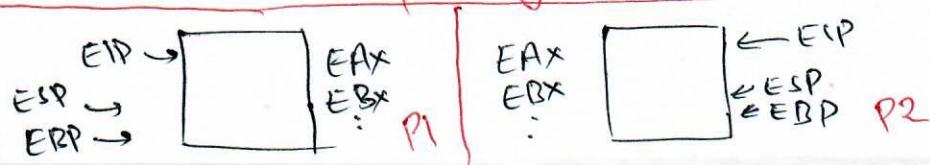
A  
B  
A  
C  
B

multiple processes are able to run even though there is single CPU

## ✓ CPU virtualization:

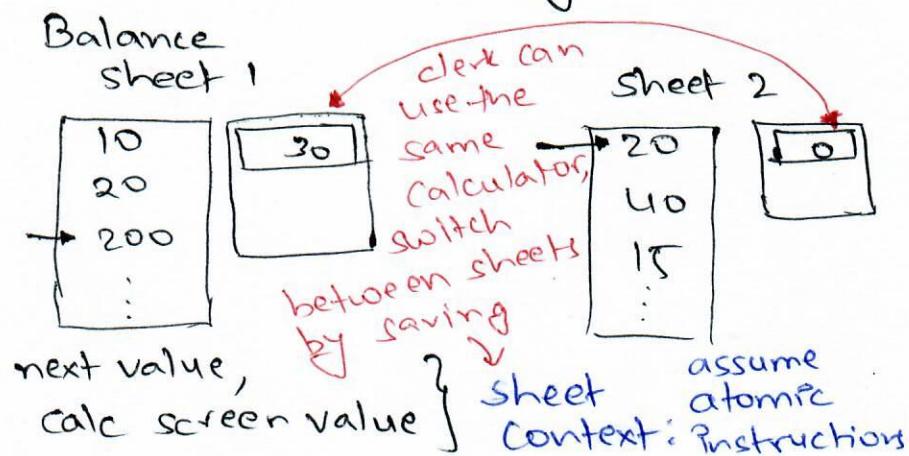
(2)

Every process THINKS that it owns CPU registers and directly changes them when it runs! Transparency



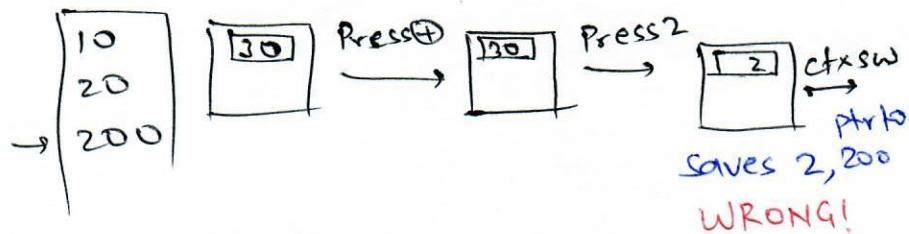
## ✓ Bank clerk analogy

(1)



## ✓ Context sw can't happen in middle of instruction

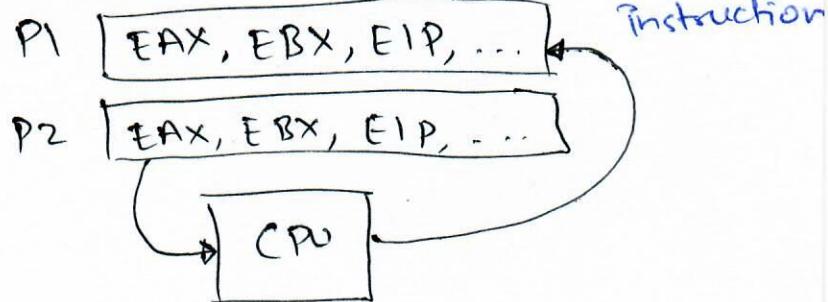
(2)



Sheet context assumes 1 instruction → adding num  
Each instruction → ATOMIC

## Context switch

Process Control block



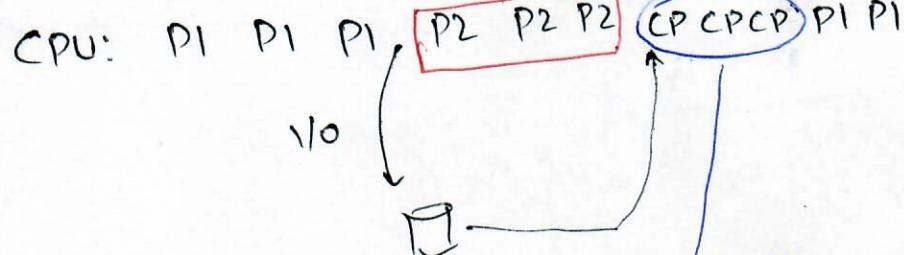
## Comparison w/ operator

Flexibility ↑

- Run multiple processes simultaneously
- Insert high-priority job without rebooting
- Dynamically Kill / spawn processes
- No need for operator → OS!

## Performance benefit

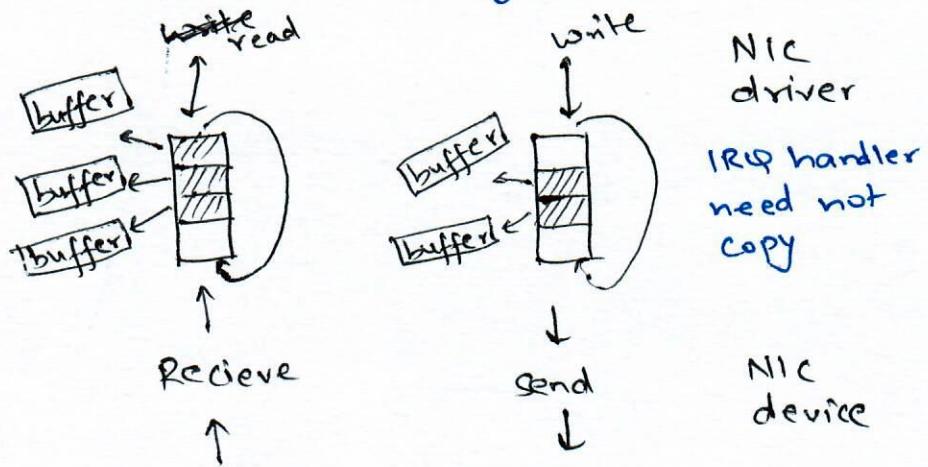
overprovision → saved wasted cycles



Further optimization

Improved resource utilization

## Aside: DMA ring buffers ②



Problem: Context switch is light ②  
→ few registers

Impact on performance is bad.

P1 ↔ P2: different address spaces

D-cache, I-cache, TLB misses

Flushed in some architectures

Improve perf. ①

- Only change registers

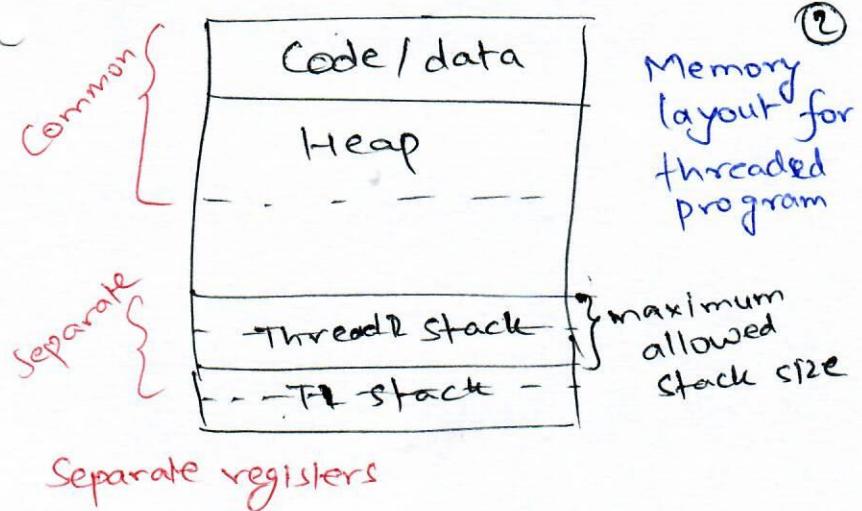
- Don't change address space

- Threads

- T1 T1 T1 T2 T2 T2 T1 T1 T1

↑  
fast

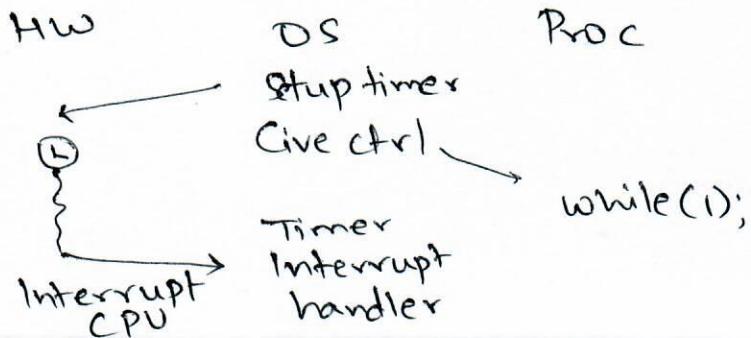
- very effective for I/O heavy processes



- Prob Give control, keep ability to take back control

`while(1);` Single CPU ①

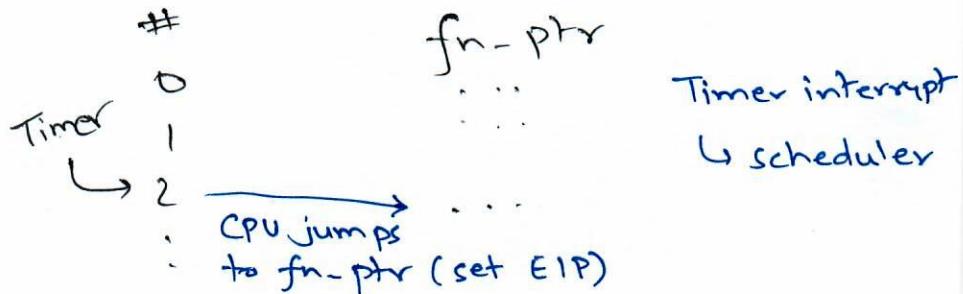
Sol<sup>n</sup>



- which handler? ②

- Interrupt descriptor table

Set up at boot

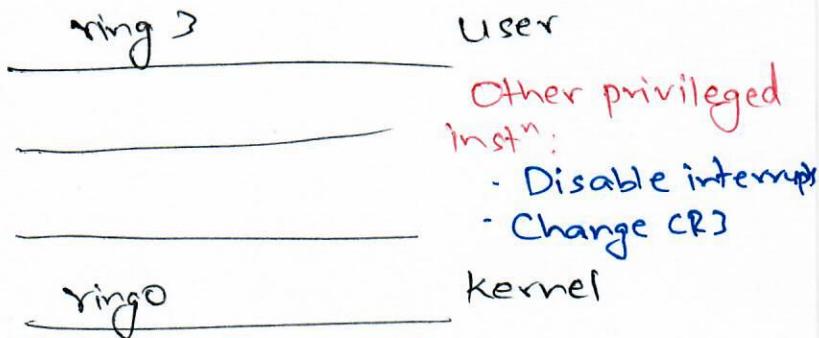


①: ~~Modify~~ IDT ②

- # fm-ptr
- 0 ---
- 1 ---
- 2 → does not call sched

Privileged Instruction  $\Rightarrow$  Traps to OS  
 $\Rightarrow$  Kill proc.

① x86 rings



② Interesting programs Need to use privileged instructions

- Sending data on network
- Reading disk contents (100x heavier than function calls)
- System calls

execute INT N instruction

③ IDT

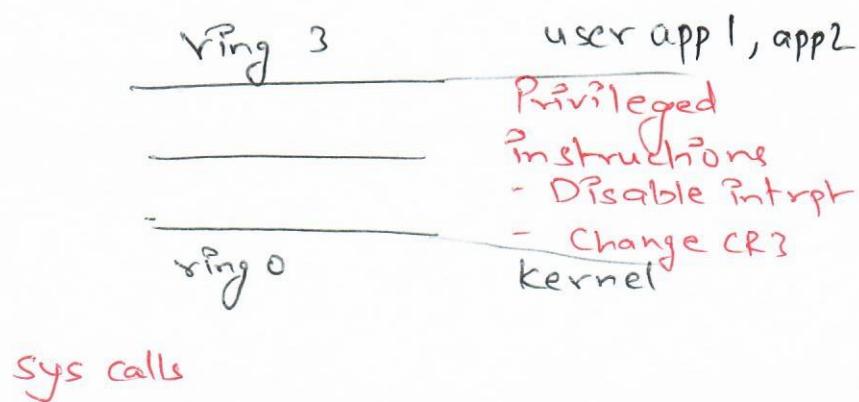
## CPU registers ( $\times 86$ )

②

- CR3 → Page table *Hidden from programs*  
EIP → Instruction pointer *Autoinc by CPU*  
ESP → Stack pointer. Top of stack  
EBP → Bottom of call frame  
General: EAX, EBX, ECX, ...  
*Managed by Compiler*

## $\times 86$ rings

②

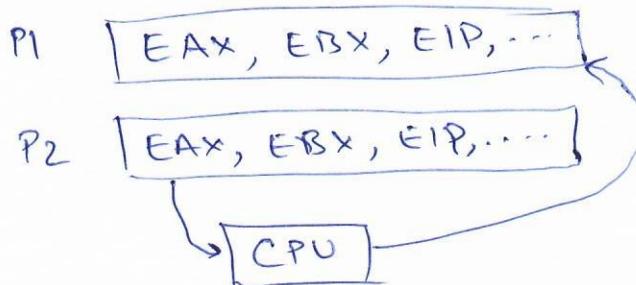


## Context switch

②

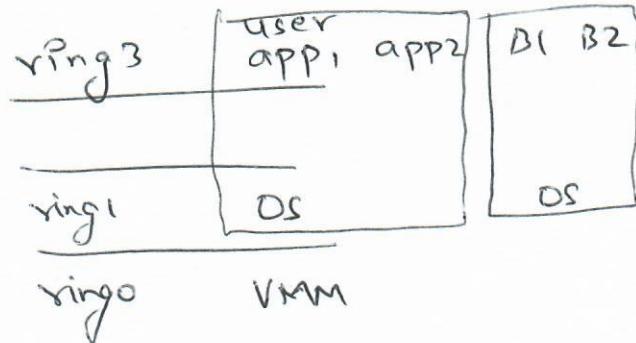
### Process control block

atomic  
Instructions



## Virtualizing OS

②



## Virtualization in OS

②

Let process think it owns  
the hardware

- CPU register - Direct execution
- Address space

\* Isolation \*

## Virtualization in VMM

③

Let <sup>guest</sup> OS think it owns hardware

- CPU registers  
*(Guest) EAX, EBX, .. CR3, Interrupt flags*
- Physical address space
- Interrupt descriptor table
- I/O devices

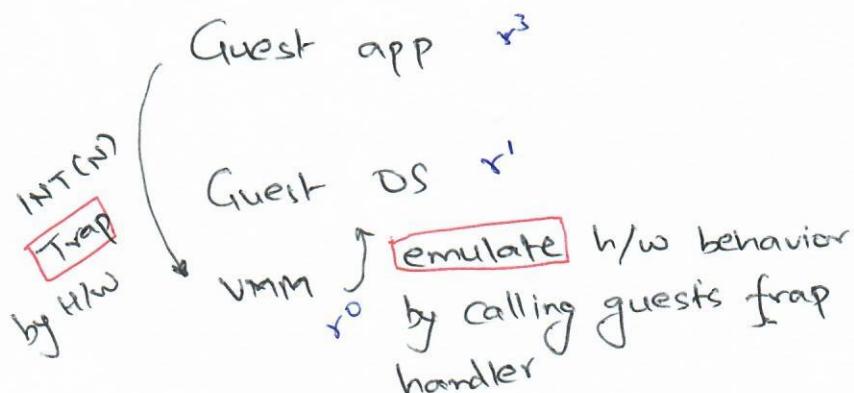
## Trap and emulate hypervisors

②



## Trap and emulate hypervisors

②



## Correctness requirements

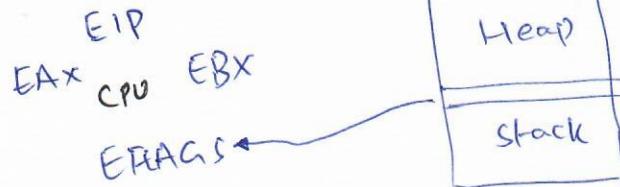
- All sensitive instructions must trap. **Popek Goldberg theorem 1971**
- What happens if change CR3 disable interrupt doesn't trap.

## EFLAGS register

③

- Interrupt enable flag (IF)

not set by **popf** instruction  
in **Ring 2** ↳ Does not trap



## Why forget Popk Goldberg theorem?

①

### Hot areas in 60s - 70s

- expensive mainframe
- multiple operating system (no posix)
- let users run own OS ⇒ no rewrite of apps
- develop new OS

### Not so hot in late 70s - late 90s

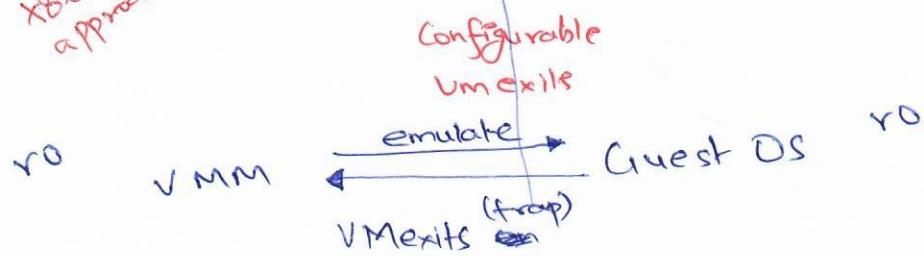
- Personal computing
- My own ~~app~~ computer. One OS

### Not again from late 90s

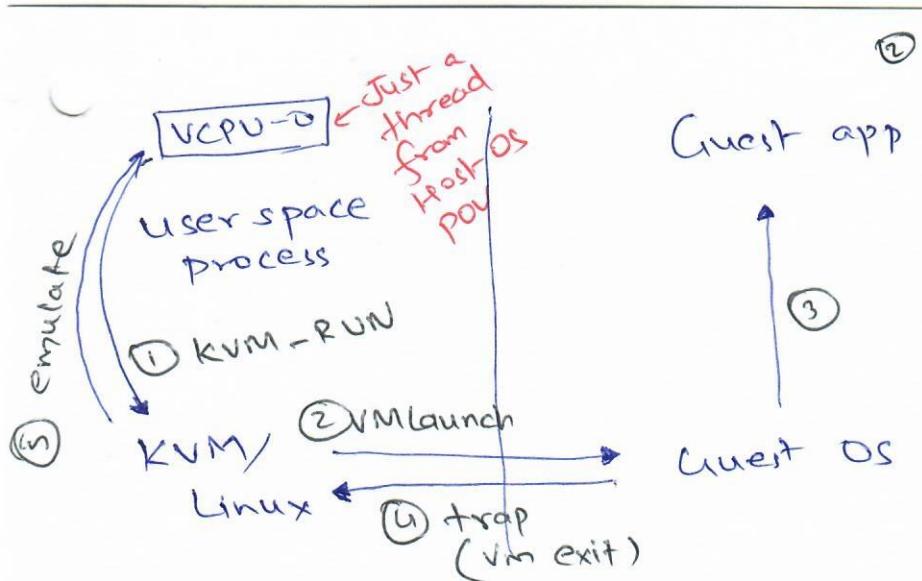
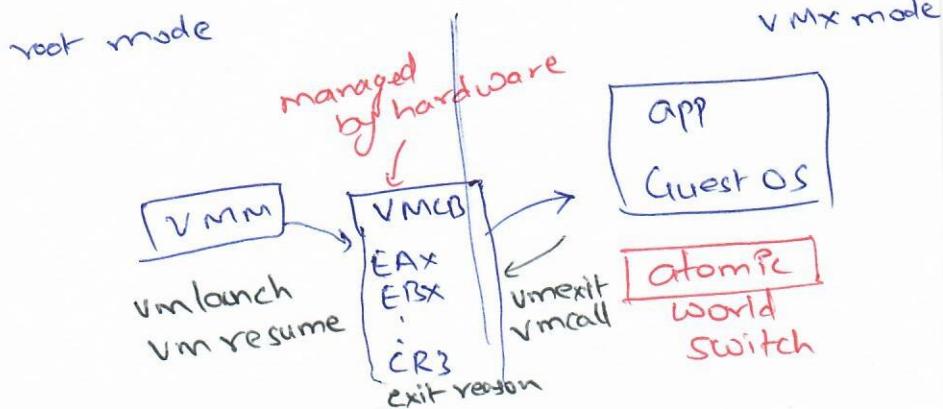
- Intrusion detection
- Fault tolerance
- multiple OS

① Root mode VMX mode/ Non-root

backward compatible  
X86 approach



② Virtual machine control block (VMCB)/ structure (vmcs)



## Direct execution

①

When CPU goes to VMX mode,  
none of host processes are running.

(Similar to when process runs,  
OS does not run)

## KVM?

①

Abstracts over hardware details

Intel  $\leftrightarrow$  ARM  $\leftrightarrow$  AMD

different instructions

Semantics

## Rvm API

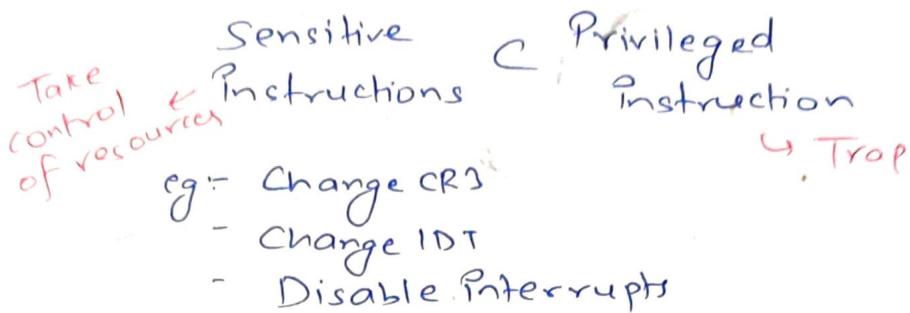
: Poch based : KVM

①

- Create VM
- allocate mem to VM
- read / write virt registers
- Run vcpu
- Inject interrupts

## Last class: Popek Goldberg Theorem

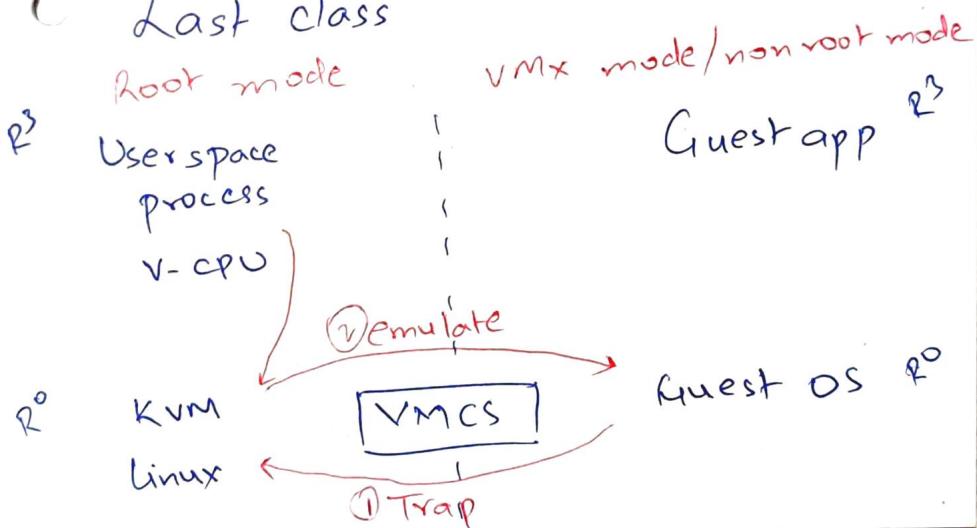
(3)



Problem with popf

## Last class

(3)



Great performance for CPU benchmarks

(2)

• 0-9% slowdown on SPEC benchmark

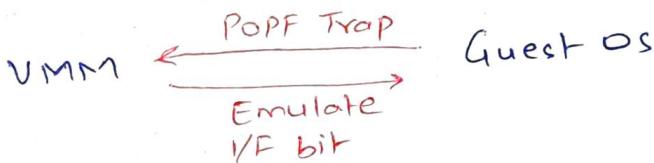
(Figure 2 ASPLoS-06. Comparison of  
sw hw tech for x86 virt)

• Housekeeping overhead Timer Interrupts

③ Exits determine performance for trap-and-emulate hypervisors

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

④ Example: guest running popf frequently



Solution: Avoid exits as much as possible!

⑤ Dealing with popf: Shadow EFLAGS register in VMCB

Do not send  
interrupts if  
shadow eflags.IF

VMM is disabled.

Guest OS

CPU Real  
eflags

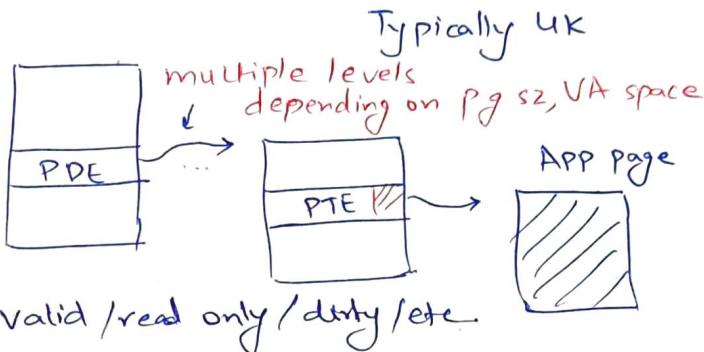
Shadow eflags  
modifies

⇒ POPF no longer traps!

# Virtualize memory! Paging

(3)

CR3 →



## C Metrics of success

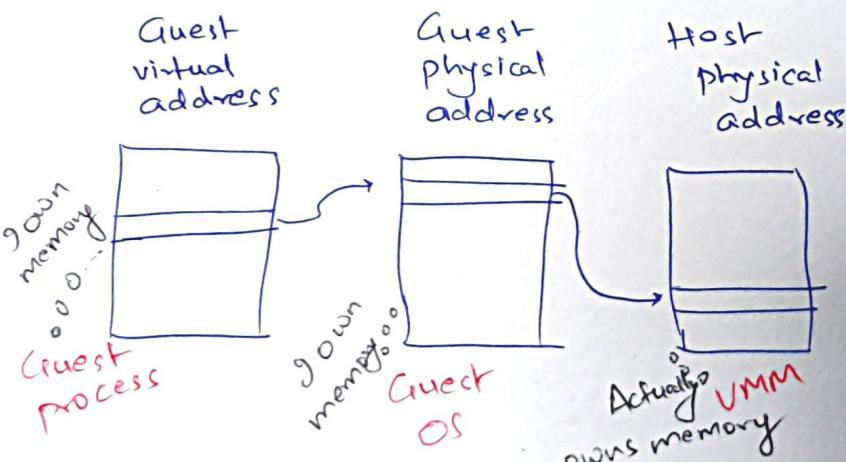
Transparency: Process thinks it <sup>Guest OS, Guest process</sup> owns memory

Isolation: One process can not read <sup>Guest OS, Guest process</sup> another's memory

Performance: Once mapped, address <sup>Guest OS, VMM</sup> translation need not involve OS  
Fast load/stores

## C General idea!

Two levels of address translation



## Idea 1: Interpretive Execution ② (Lab)

$30xx \rightarrow \text{mem}[xx] = AC$   
Virtual address  
array owned  
by Interpreter

✓ Isolation

✗ Performance

✓ Transparency

eg: IBM System 370  
SIE instruction → Start Interpretive Exec

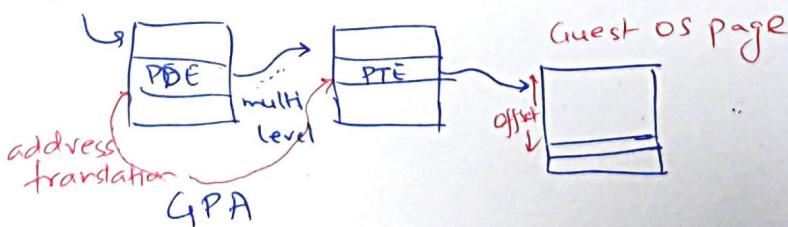
## Idea 2: Reuse paging hardware ② Shadow Page table:

map directly from GVA → HPA

- Most memory ops LD/ST run at native speed

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

CR3 points to a page in host PA



Now guest OS can boot. (3)

- Initially, address translation is disabled on real hardware
- CR0, bit 31 can enable-disable paging.

Guest tries to disable paging

Trap → VMM changes  
← emulate { CR3 to the prev. page table

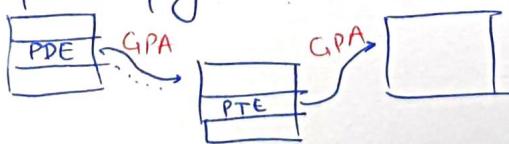
Guest "thinks" it has disabled paging. (3)

Reading CR0.31 shall { trap or return false { shadow register

But address translation is active.  
Addresses treated as GPA

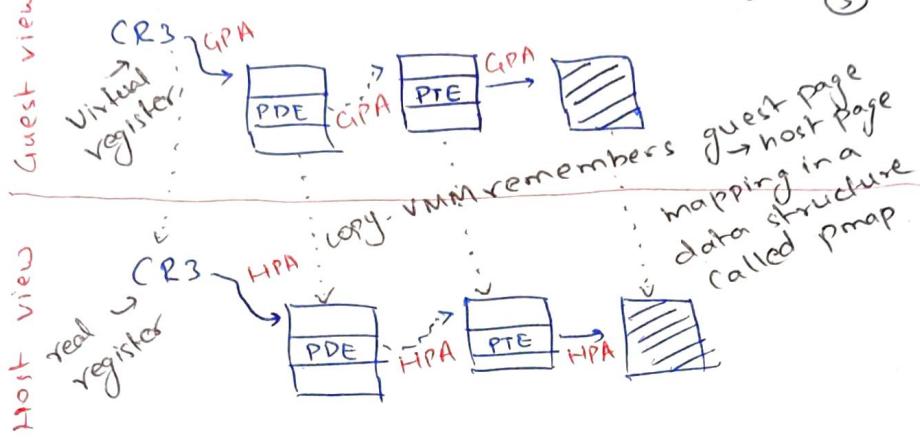
Guest tries to run a process

① Prepares page table (2)

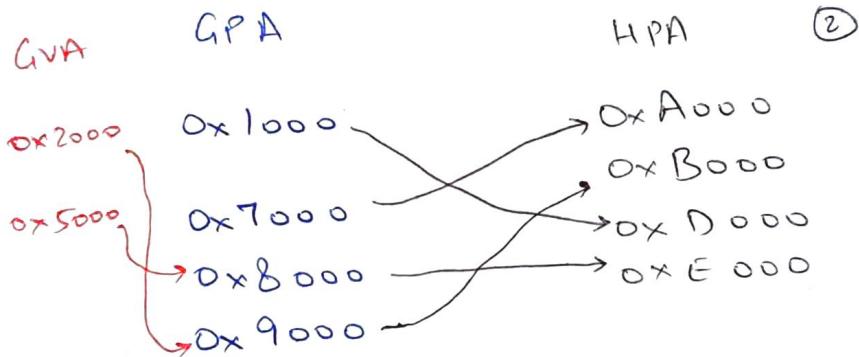


② Changes CR3 ↗ Trap to top-level GPA

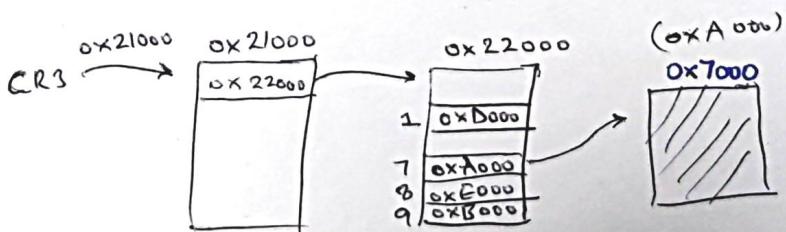
## VMM constructs shadow Page table ③



## Example: memory mapping setup for VM ②

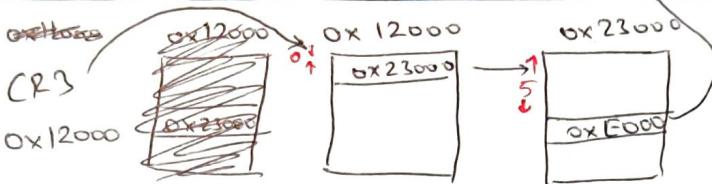
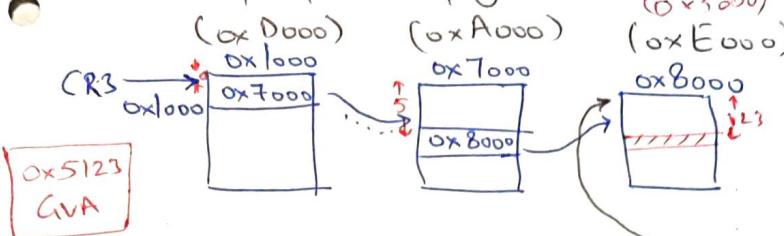


## Example: memory setup for VM



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

Guest prepares page table



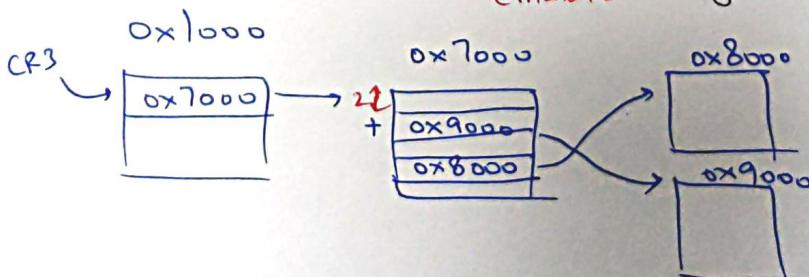
- Transparency: If guest reads CR3, it  
is emulated or kept as shadow CR3

Performance: Directly use MMU once page table is setup

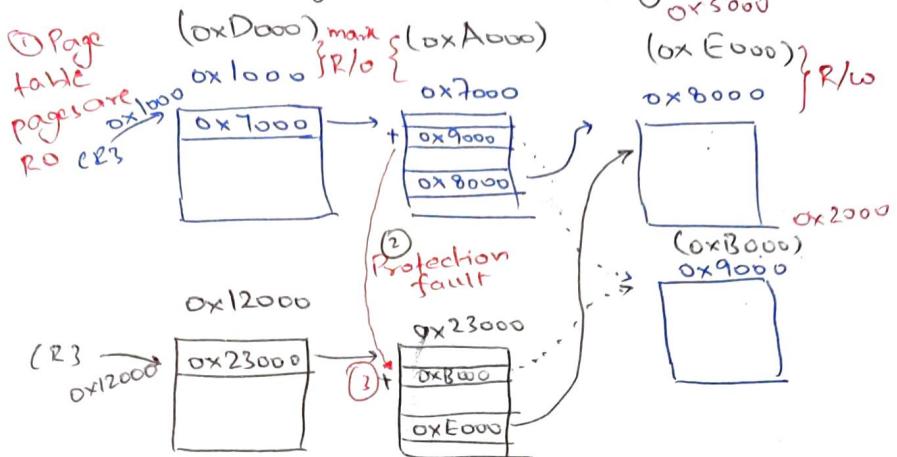
Safety / Isolation: Guest can't directly touch shadow page table. 0x12000, 0x23000 are NOT in guest's address space

- Guest 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.

- 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

→ Page fault

Bring from disk to 0xF000.

Update ALL shadow page tables pointing to 0xE000.

③ Maintain backward mappings in "pmap"

- ③ Move 0xF000 to disk. When ② guest OS tries to update page table.  
→ Page fault  
Bring to 0xC000

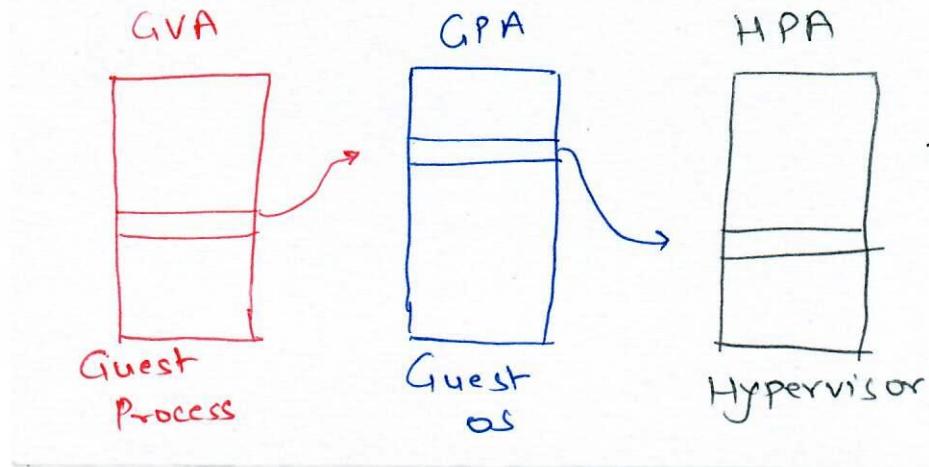
Update GPA  $\rightarrow$  ~~GVA~~ HPA mapping  
 $0x7000 \rightarrow 0xC000$

- ④ Summary:
  - 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  $\rightarrow$  mapping backward

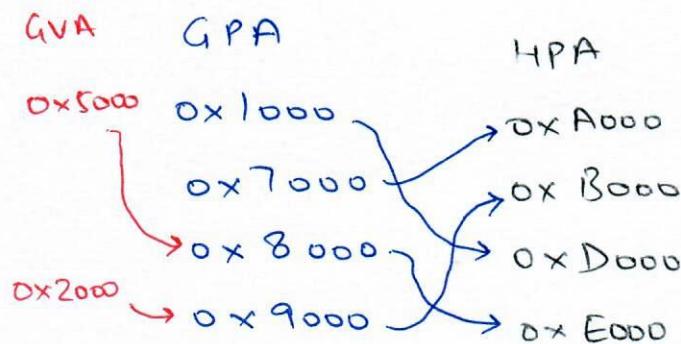
O Transparency ✓ Isolation ✓ ②

- 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

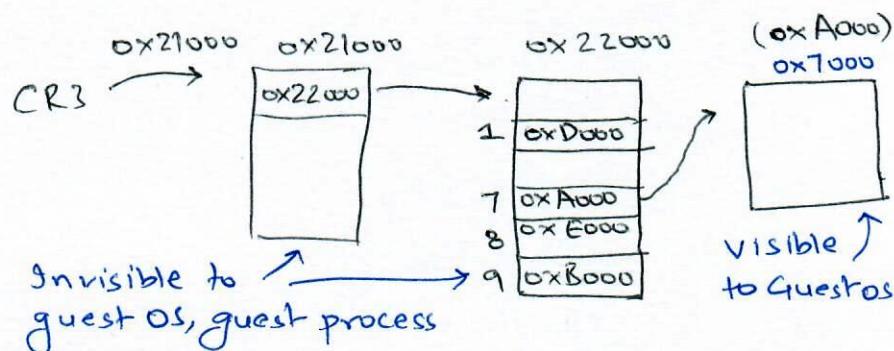
## Memory virtualization in VMM ②

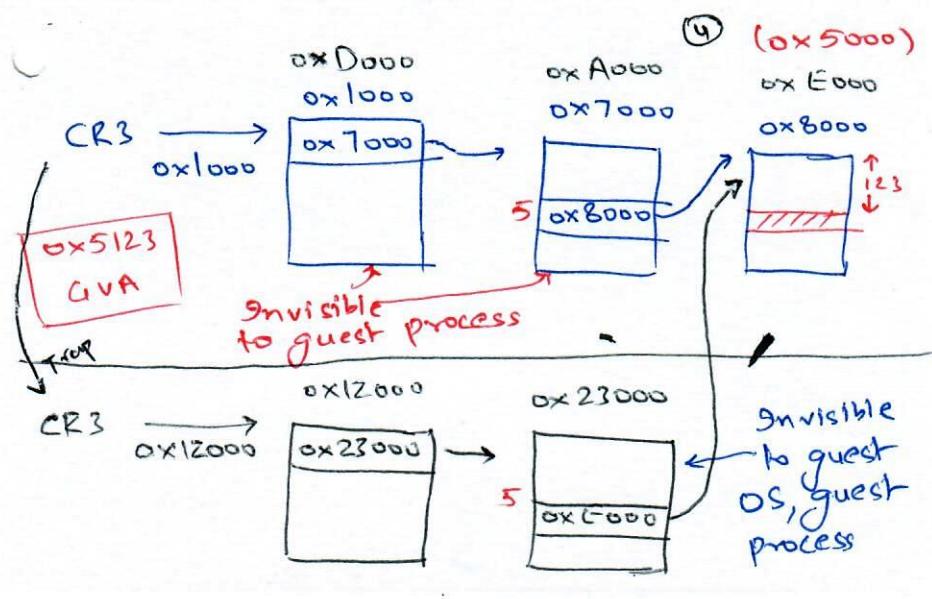


### Example: memory mapping for VM ②

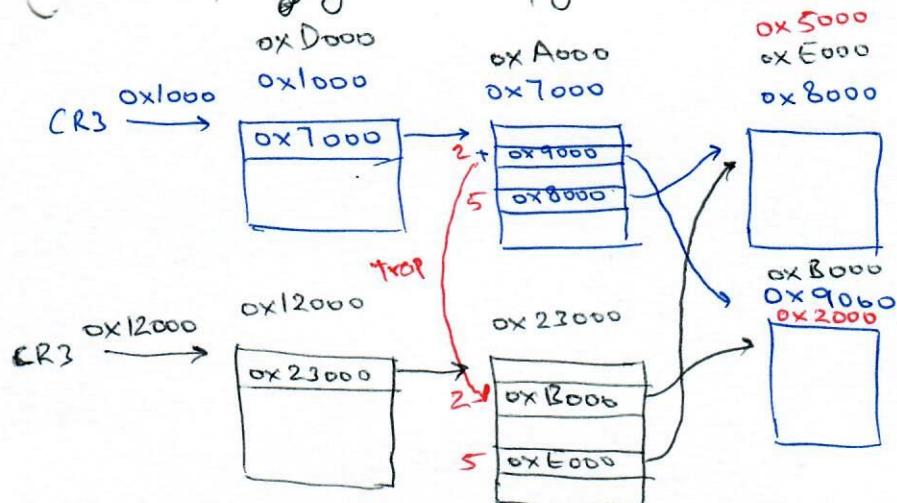


When guest boots, CR0.31 is disabled (virtual memory is disabled)  
 $\text{GPA} \rightarrow \text{HPA}$





### Maintaining shadow page tables ②



### Process fork

(fork.c) ③

```

int main() {
    printf("Hello");
    int rc = fork();
    int x = 5;
    if (rc < 0) { exit(1); /* fork failed */ }
    if (rc > 0) { printf("parent"); x = 1; }
    if (rc == 0) { printf("child"); x = 2; }
    printf("x: %d", x);
}

```

## Observations

②

- Child does not print "Hello"
- Parent prints  $x \rightarrow 6$   
Child prints  $x \rightarrow 10$

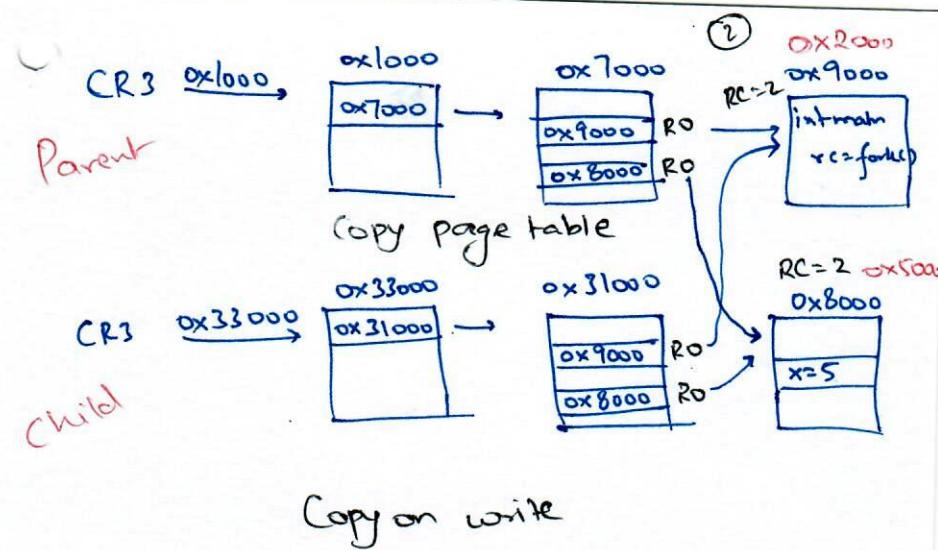
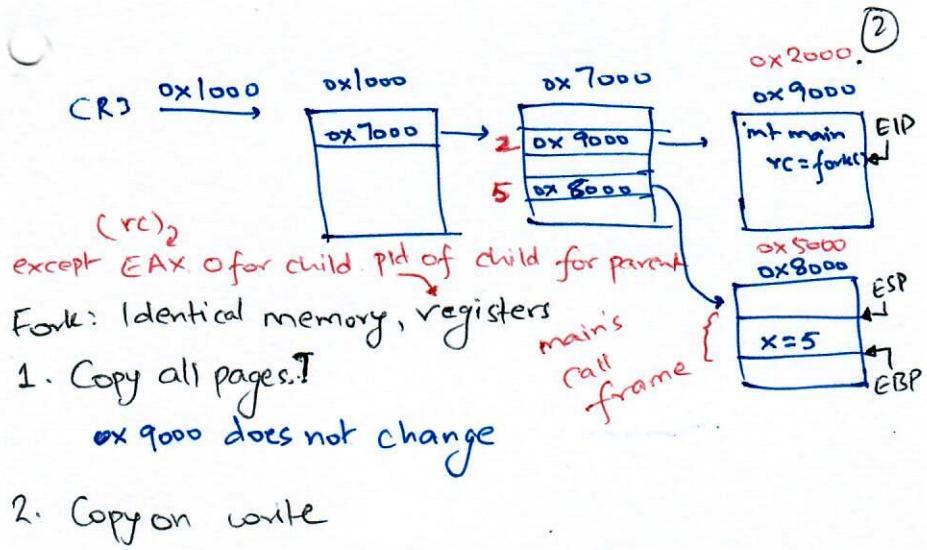
other things are also passed along-

- open file descriptor
- namespaces, priority, etc.

## Why?

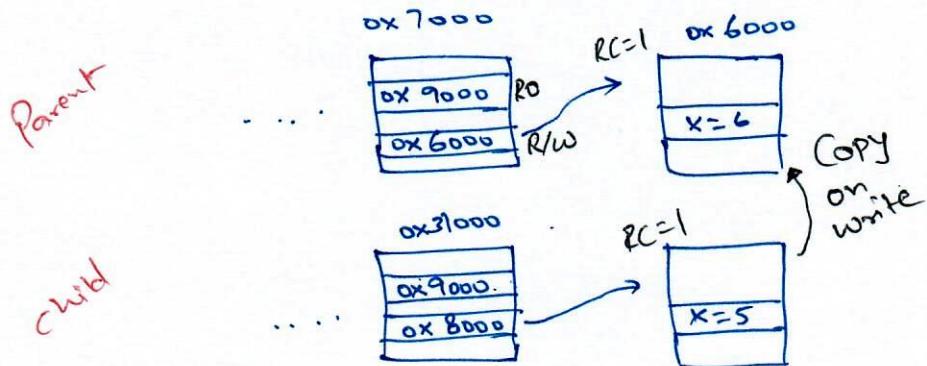
①

- \* shell.c → gap between fork and exec for output redirection
- \* In-memory storage (e.g. Redis) snapshot
  - Fork → write DB to disk
  - parent process continues to serve requests



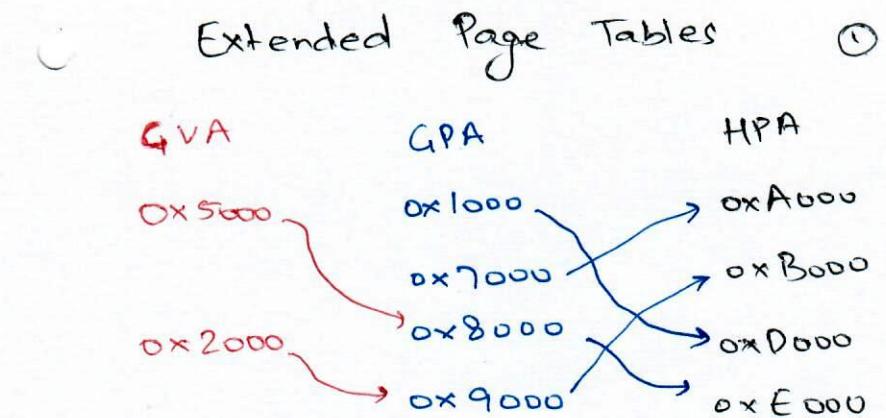
① Parent does  $x += 1 \Rightarrow ST 6 \rightarrow x$

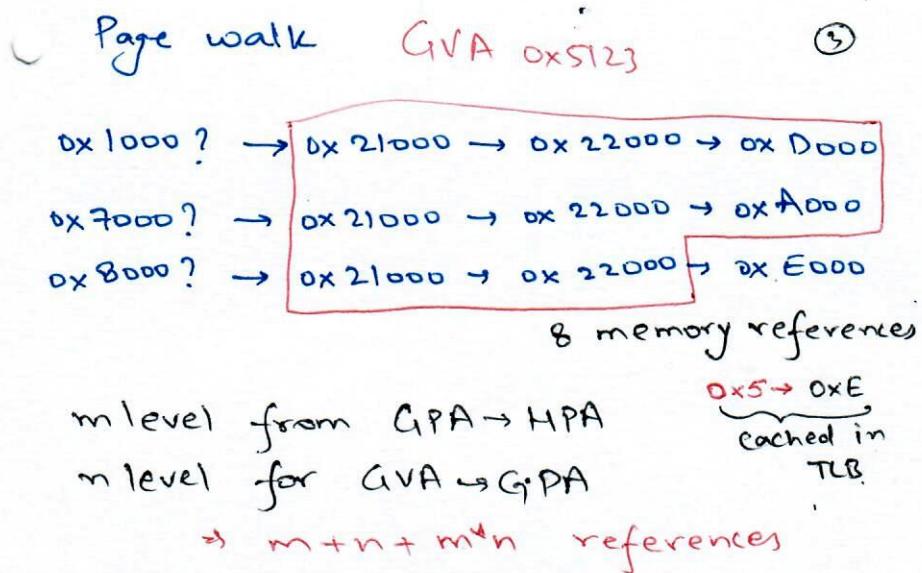
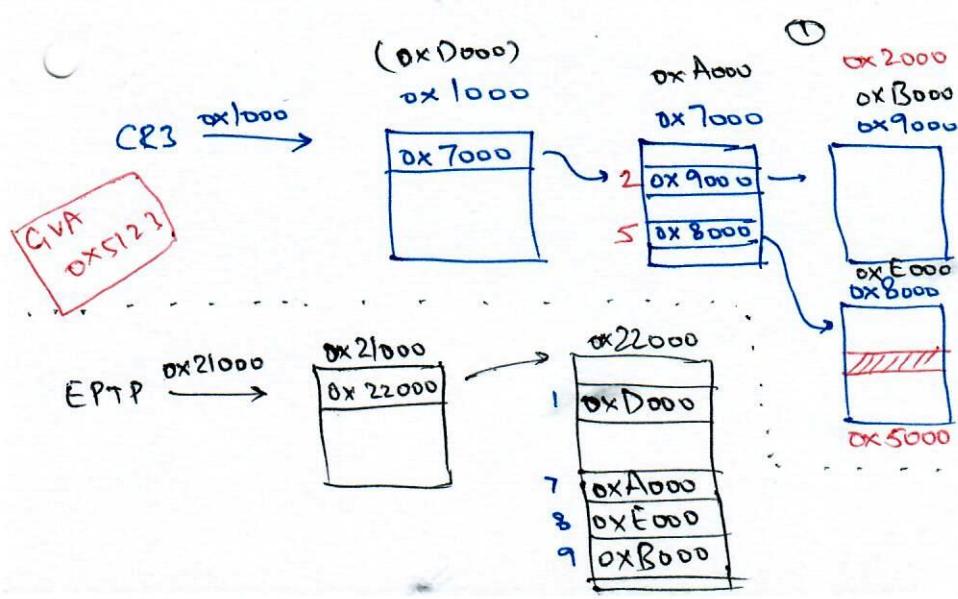
② Protection fault: Page is read-only



- 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 → 4.4 x slowdown
  - PTE modification
    - native → single cycle store
    - shadow PT → 12733 cycles





- Hardware support for 2-D address translations
- Changing CR3 need not trap ⇒ fast ctx sw
- Modifying PTEs need not trap ⇒ fast fork

(4)

Table 1:- Instruction and data translations

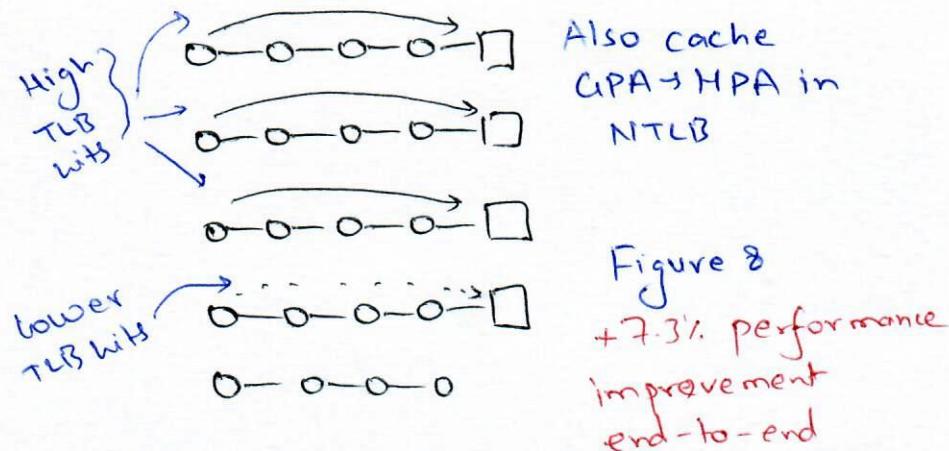
$3.9x - 4.57x$  slowdown

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  
(ASPLOS'08) virtualized systems

Two main ideas: 1) Nested TLB ②



2) Large pages ②

- Increases TLB reach → Also helps in shadow page tables
- Reduces m and n
- 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

①

### Reclamation

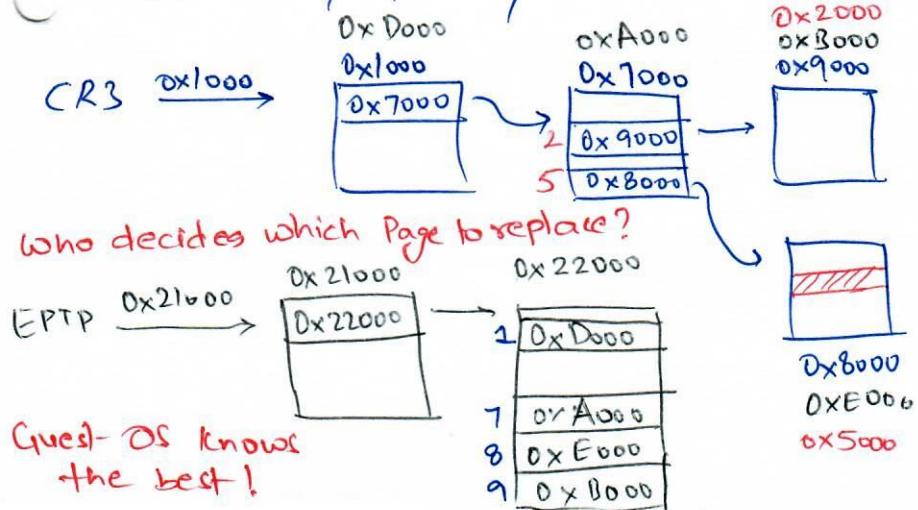
- Ballooning

## Sharing memory

- Content based page sharing

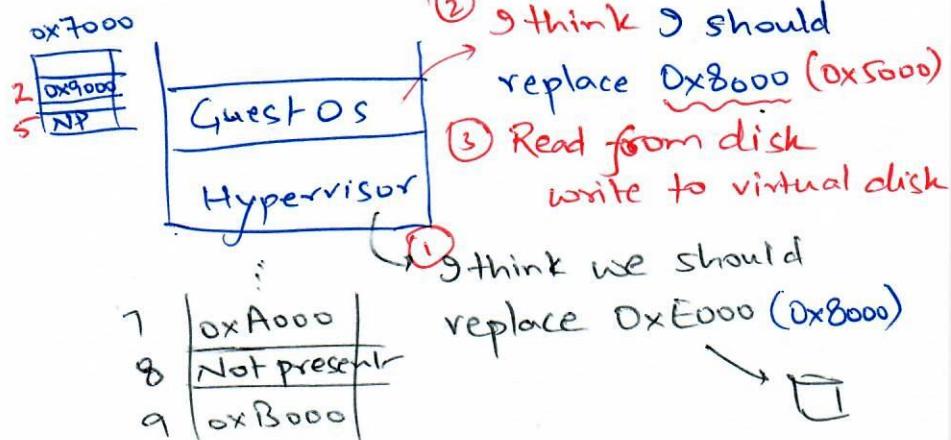
## Issues w/ transparent reclamation

②



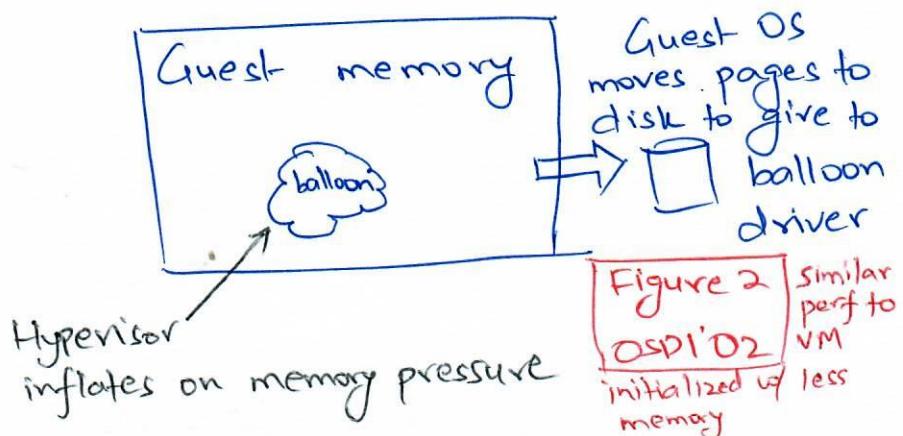
## Double paging problem

③



## 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

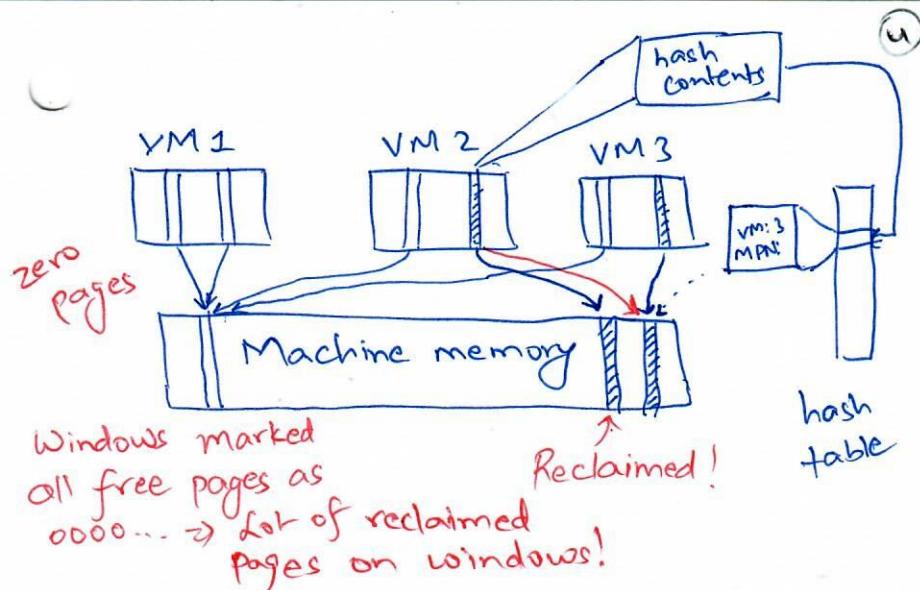
## Content-based page sharing

(4)

- Scan a random page periodically
- Hash ↗ opt: R/o page marked by OS
- Match ⇒ fully match page
- Match ⇒ Mark (w. Reclaim)

⇒ Transparent to Guests!

Figure 5 (OSD1|O2) 7-32: reclaimed memory



### Interesting trivia

Linux KSM (kernel Same page merging)  
/ .. shared memory

- Similar hash-based implementation [Nov '08]
- VMWare Patent!
- Red-black tree based implementation [Apr '09]