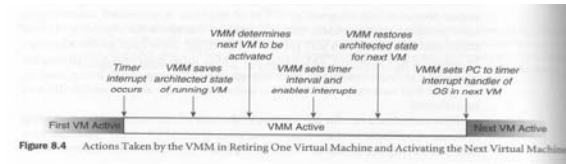


## System VMs

CS680V/RI

## Resource Control

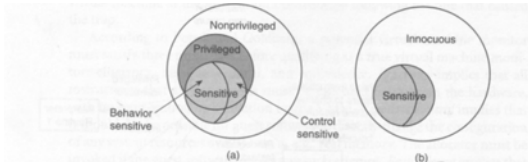
- Issue: How to retain control of resources in the VMM?
- Timer interval control performed by VMM
  - Also, guest OS not allowed to read the timer value
    - Guest OS sees a *virtual interval timer*
- VMM also gains control whenever guest OS executes privileged instructions.



CS680V/RI

## Instruction Types

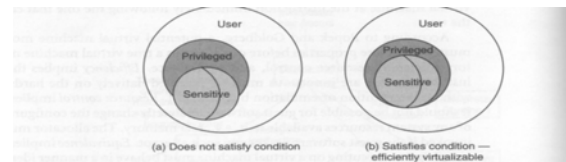
- Non-privileged: Do not cause traps
- Privileged: Cause Traps
- Sensitive: Change/depend upon system state
- Innocuous: Not "sensitive"



CS680V/RI

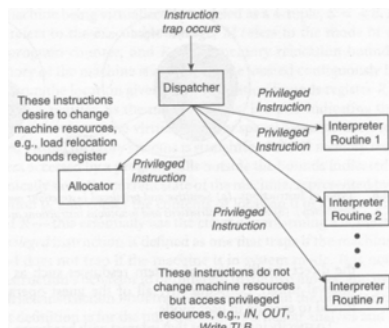
## Conditions of ISA Virtualizability

- Theorem: A computer architecture is *fully virtualizable* if the set of sensitive instructions for that computer is a subset of the set of privileged instructions.



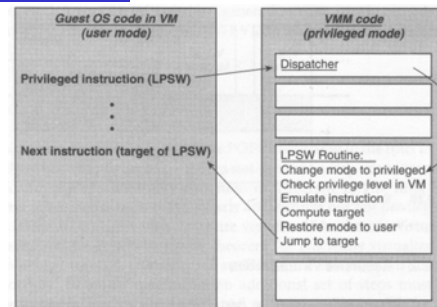
CS680V/RI

## Execution of Privileged Instruction



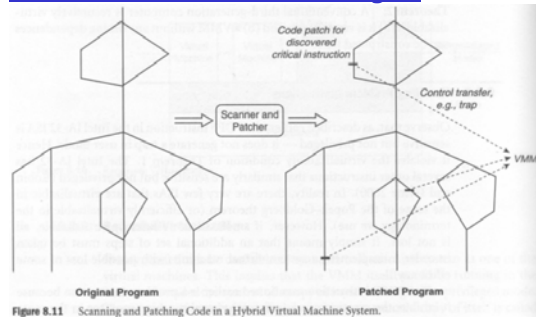
CS680V/RI

## Handling Privileged Instructions in a Guest OS



CS680V/RI

## Handling Problem Instructions - Sensitive But Not Privileged



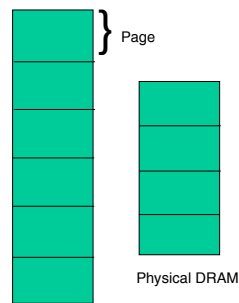
CS680V/BI

## Memory Virtualization

CS680V/BI

## Virtual Memory Basics

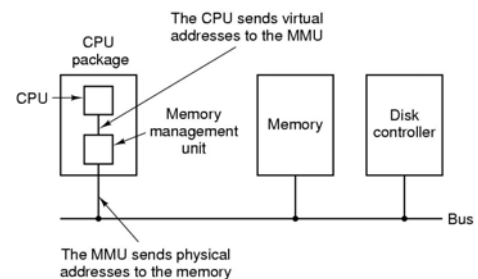
- Sometimes, a single process might require more memory than the available DRAM in the system.
- In such cases swapping is not enough. Rather, we need to break up the memory space of a process into smaller equal-sized pieces (called pages).
- Operating system then decides which pages stay in memory and which get moved to disk.
- **Virtual memory:** means that each process gets an illusion that it has more memory than the physical DRAM in the system.



Process Virtual Address Space

CS680V/BI

## MMU and Virtual Memory

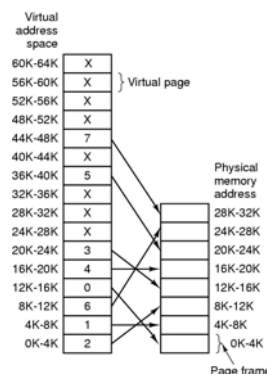


- MMU = Memory Management Unit
- Part of Hardware that accompanies the CPU
- Converts Virtual Addresses to Physical Addresses

CS680V/BI

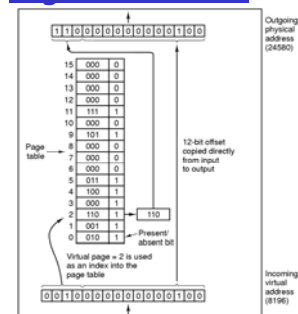
## Paging

The mapping between virtual page numbers and physical page numbers is stored in the page table



CS680V/BI

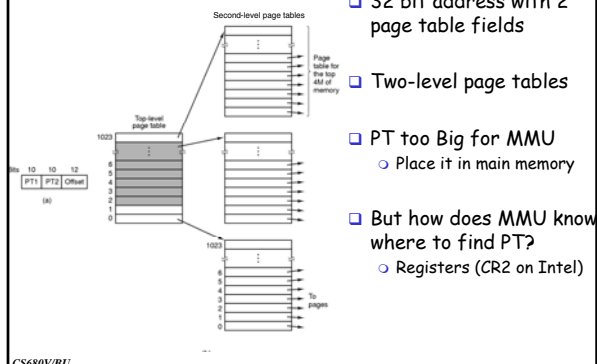
## Page Tables (1)



Internal operation of MMU with 16 4 KB pages

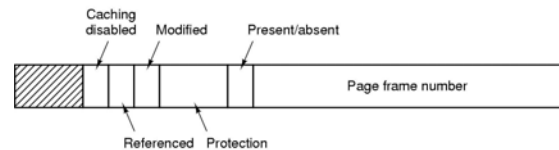
CS680V/BI

## Page Tables (2)



CS680V/RI

## Page Tables (3)



Typical page table entry

CS680V/RI

## TLBs - Translation Lookaside Buffers

Valid	Virtual page	Modified	Protection	Page frame
1	140	1	RW	31
1	20	0	R X	38
1	130	1	RW	29
1	129	1	RW	62
1	19	0	R X	50
1	21	0	R X	45
1	860	1	RW	14
1	861	1	RW	75

- TLB is a small cache that speeds up the translation of virtual addresses to physical addresses.
- TLB is part of the MMU hardware (comes with CPU)
- It is not a Data Cache or Instruction Cache. Those are separate.
- TLB simply caches translations from virtual page number to physical page number so that the MMU don't have to access page-table in memory too often.
- On x86 architecture, TLB has to be "flushed" upon every context switch because there is no field in TLB to identify the process context.

CS680V/RI

## Impact of Page Size on Page tables

Small page size

- Advantages
  - less internal fragmentation
  - better fit for various data structures, code sections
- Disadvantages
  - If a process needs more pages, then it has larger page table

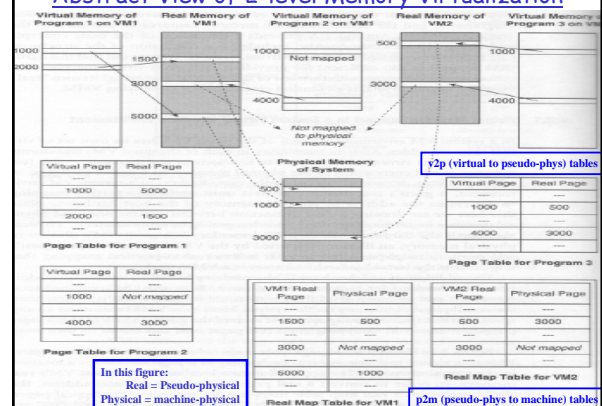
CS680V/RI

## Memory Virtualization for VMs

- Not the same as virtual memory for processes
  - But very similar
  - Just another layer of indirection
- Address Translation process
  - Level 1: Guest VM page table
    - Virtual address → pseudo-physical address
  - Level 2: VMM Page table
    - Pseudo-physical address → machine-physical address
- Sum of real memory of all guest VMs **could be** greater than physical memory of the system
  - Thus VMM may have its own swap space, independent of what guest VMs have.

CS680V/RI

## Abstract View of 2-level Memory Virtualization



## Two Architectures for implementing the abstract view

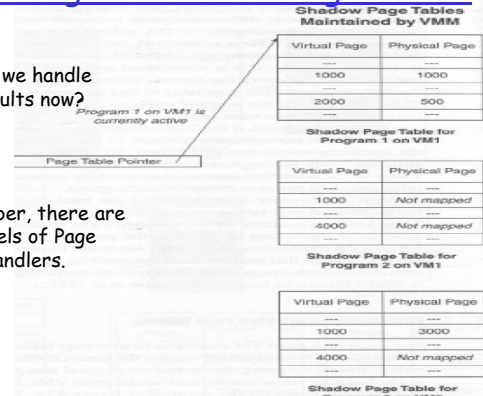
- ❑ Architected Page Tables
  - Page table interface defined by ISA and understood by memory translation hardware
  - E.g. x86 architecture
- ❑ Architected TLBs
  - TLB interface defined by ISA and understood by memory translation hardware
  - E.g. alpha architecture

CS680V/RI

## Virtualizing Architected Page Tables

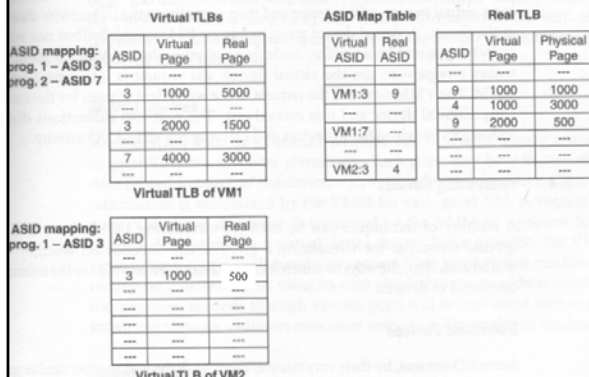
- ❑ How do we handle page-faults now?

- ❑ Remember, there are two-levels of Page Fault handlers.



CS680V/RI

## Virtualizing Architected TLBs



## I/O Virtualization

- ❑ A number of different types of I/O devices
- ❑ Construct a virtual version of the device
- ❑ I/O activity directed at the device is intercepted by VMM and converted to equivalent request for underlying physical device.

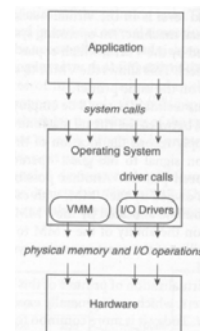
CS680V/RI

## Device Types

- ❑ Dedicated Device
  - Display, keyboard, mouse etc.
  - VMM routes, but does not interpret the I/O instructions
- ❑ Partitioned Devices
  - E.g. A hard disk can host several virtual disks
- ❑ Shared Devices
  - E.g. network adapter
- ❑ Spooled Devices,
  - E.g. printers
  - Shared, but at a higher granularity than shared devices
- ❑ Nonexistent Physical Devices
  - E.g. network adapter to communicate only among VMs
  - E.g. Virtual disk hosted on a non-disk physical device

CS680V/RI

## Different I/O interfaces



CS680V/RI

### Virtualizing at the I/O Operations Level

- ❑ Directly intercept I/O instructions below the device driver.
- ❑ Execute the effect of the instruction on physical device
- ❑ Difficult to reverse engineer the combination of individual I/O operations that device drivers intends to do.

CS680V/BI

### Virtualizing at the Device Driver Level

- ❑ VMM has knowledge of the guest OS internals
- ❑ It inserts virtual drivers for each virtual device in the guest OS.
- ❑ Typically used in hosted VMs.

CS680V/BI

### Virtualizing at the System Call Level

- ❑ VMM can intercept system calls for more efficient virtualization.
- ❑ However, VMM needs to know all the calls in ABI (application binary interface).
- ❑ Plus need to account for interaction between ABI and other parts of the OS.

CS680V/BI