

## Paging (cont): Inverted Page Table, Bits in PTE, TLB

ECE595

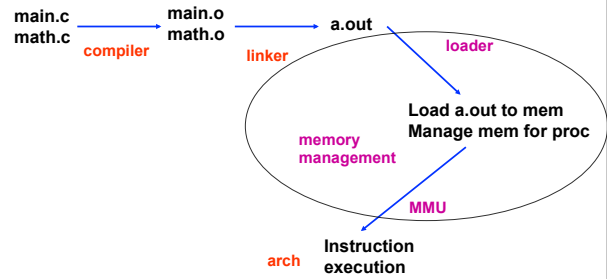
Oct 16

Y. Charlie Hu



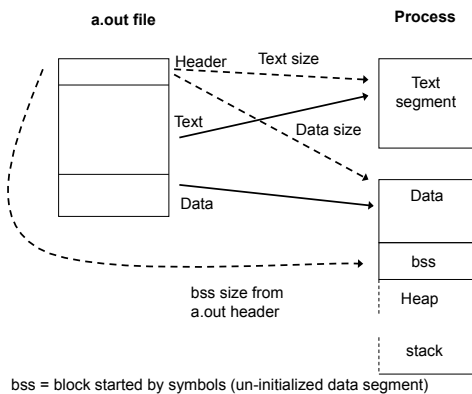
1

## The big picture



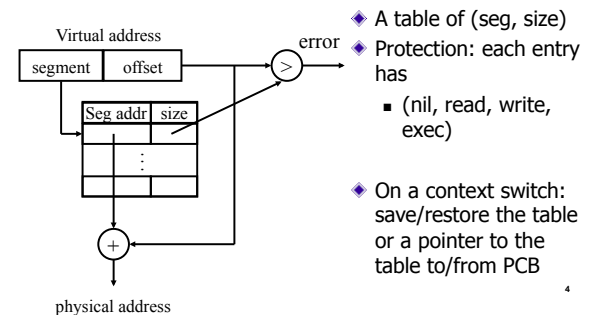
2

## [lec12] Loading



3

## [lec13] Segmentation



4

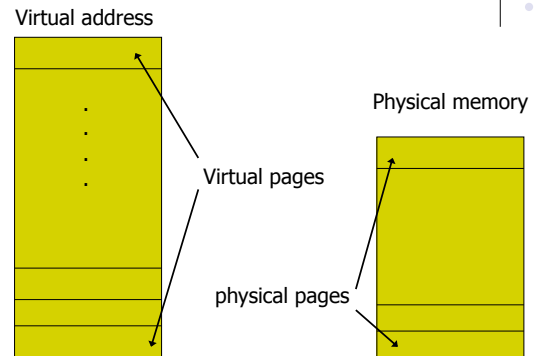
## [lec13] What fundamentally causes external fragmentation?

- Segments of many different sizes
- Each has to be allocated contiguously

- “Million-dollar” question:  
*Physical memory is precious.*  
*Can we limit the waste to a single hole of X bytes?*

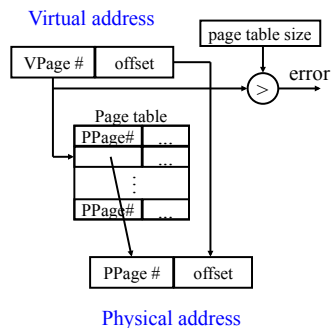
5

## [lec14] Virtual pages / physical pages



6

## [lec14] Paging



- Context switch
  - similar to the segmentation scheme
- Pros:
  - easy allocation, keep a free list
  - easy to swap
  - easy to share

7

## Paging vs. segmentation

- Segmentation:
  - External fragmentation
  - Complicated allocation, swapping
  - + Small segmentation table
- Paging
  - Internal fragmentation
  - + Easy allocation, swapping
  - Large page table

8

## Deep thinking

- In segmentation, why does each segment need to be contiguous in physical memory?
- In segmentation, what to do with heap/stack?
  - What happens when they grow/shrink?
- In paging, do pages belonging to the same “segment” (e.g. heap) need to be contiguous in physical memory?
  - What made this possible?
  - What to do with heap/stack growing/shrinking now?



## [lec14] How many PTEs do we need?

- Worst case for 32-bit address machine
  - # of processes  $\times 2^{20}$  (if page size is 4096 bytes)
- What about 64-bit address machine?
  - # of processes  $\times 2^{52}$



10

## Page table

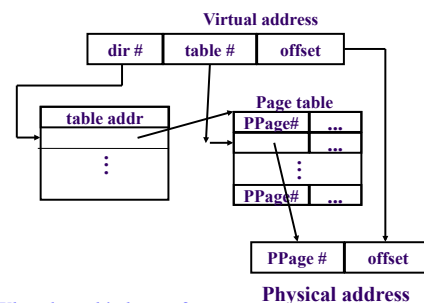
- The page table has to be contiguous in physical mem
  - Potentially large
  - Consecutive pages in mem hard to find
- How can we be flexible?

“All computer science problems can be solved with an extra level of indirection.”



11

## Two-level page tables



What does this buy us?



12

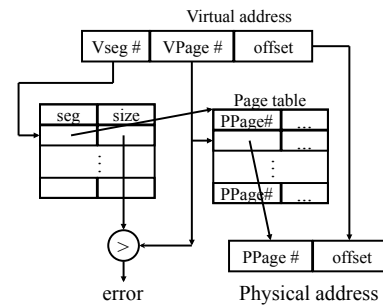
## Multi-level page tables

- 3 Advantages over 1-level page table?

*The power of an extra level of indirection!*

13

## Segmentation with paging



Ex: IBM System 370 (24-bit, 4-bit segment #, 8-bit page #) <sup>14</sup>

## [lec1] Separating Policy from Mechanism

Mechanism – tool to achieve some effect

Policy – decisions on how to use tool  
examples:

- All users treated equally
- All program instances treated equally
- Preferred users treated better

Separation leads to flexibility

15

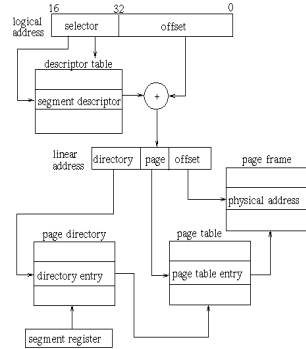
## Segmentation + paging vs. multi-level paging

- *Mechanisms* are similar
- Difference lies in *policy*
  - Segmentation + paging still maintains notion of **segments**
  - Multi-level paging deals **the whole, uniform address space** (like one-level paging)

16

## The Intel Pentium (1993) (pro, II, III, 4) (Ch 8.7, fig 8.22, 8.23)

- Supports both pure segmentation and segmentation with 1-level paging (page size=4M) or 2-level paging (page size=4k)
- CPU generates logical addresses
  - (selector, offset), 16 bits and 32 bits
  - As many as 16K segments
  - Up to 4GB per segment



## Linux on Pentium

- Linux uses 3-level paging
  - For both 32-bit and 64-bit architectures
- On Pentium, degenerates to 2-level paging
  - Middle-level directory has zero bits

18

## Today's topics

- Inverted page table
- Bits in a PTE
- TLB

19

## [lec14] How many PTEs do we need?

- Worst case for 32-bit address machine
  - # of processes  $\times 2^{20}$  (if page size is 4096 bytes)
- What about 64-bit address machine?
  - # of processes  $\times 2^{52}$

*Hmm, but my PC only has 1GB, 256K PTEs should be enough?!*

20

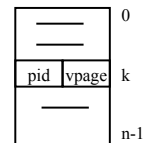
## Inverted Page Table

- Motivation
  - Example: 2 processes, page table has 1M entries, 10 phy pages
- Is there a way to save page table space?

21

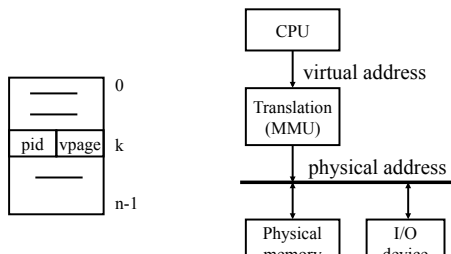
## Ideally,

- One PTE for each physical page frame, disregarding how many processes
  - Assuming rest virtual address not allocated/used



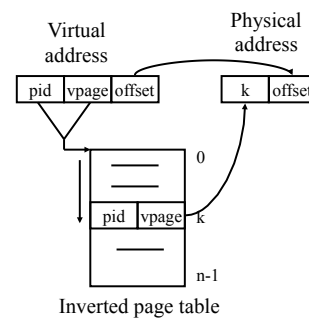
22

## But,



23

## Inverted page tables



- Main idea
  - One PTE for each physical page frame
  - Hash (Vpage, pid) to Page#
- Pros
  - Small page table for large address space
- Cons
  - Lookup is difficult
  - Overhead of managing hash chains, etc
- Ex: 64-bit UltraSPARC<sub>24</sub>, PowerPC

## Deep thinking

- How can two processes share memory under inverted page table?

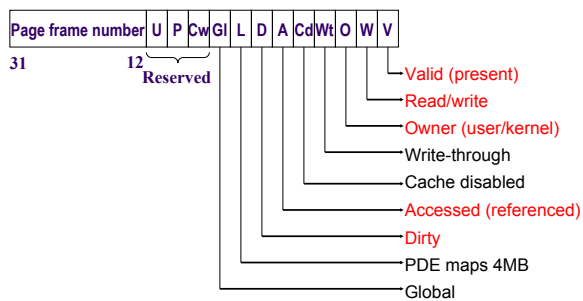
25

## What is happening to heap phy mem allocation before/after malloc()?

- Before malloc()?
- After malloc()?
- Upon first access?
- How to capture the first write to a virtual page?
  - e.g. want to trap into page fault handler

26

## x86 Page Table Entry



27

## What is happening before and after malloc()?

- Before malloc()?
- After malloc()?
- Upon first access?
- How to capture the first write to a virtual page?
  - e.g. want to trap into page fault handler
    - Use valid bit
  - In handler, check if vpage is malloced.
    - If not, segmentation fault
    - Else allocate physical page

28

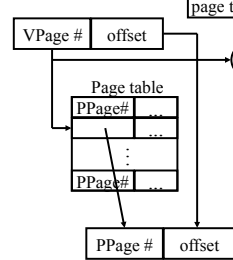
## Today's topics

- Inverted page table
- Bits in a PTE
- TLB

29

## [lec14] Paging implementation – how does it really work?

Virtual address



Physical address

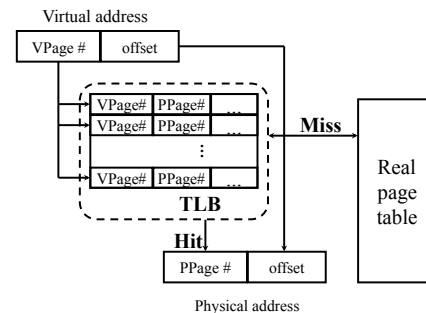
- Where to store page table?
- How to use MMU?
  - Even small page tables too large to load into MMU
  - Page tables kept in mem and MMU only has their base addresses
    - What does MMU have to do?
- Page size?
  - Small page -> big table
    - 32-bit with 4k pages
  - Large page -> small table but large internal fragmentation

## Performance problem with paging

- How many extra memory references to access page tables?
  - One-level page table?
  - Two-level page table (midterm problem)?
- Solution: *reference locality!*
  - In a short period of time, a process is likely accessing only a few pages
  - Instead of storing only page table starting address in MMU,

31

## Translation Look-aside Buffer (TLB)



TLB often fully set-associative → least conflict misses  
Expensive → typically 64 – 1024 entries

32



## Bits in a TLB Entry

VPage#	PPage#	...
VPage#	PPage#	...
		⋮
VPage#	PPage#	...

TLB

- Common (necessary) bits
  - Virtual page number: match with the virtual address
  - PTE
- Optional (useful) bits
  - ASIDs -- Address-space identifiers (process tags)

33

## Miss handling: Hardware-controlled TLB

- On a TLB hit, MMU checks the valid bit
  - If valid, perform address translation
  - If invalid (e.g. page not in memory), MMU generates a page fault
    - OS performs fault handling
    - Restart the faulting instruction
- On a TLB miss
  - MMU parses page table and loads PTE into TLB
    - Needs to replace if TLB is full
    - PT layout is **fixed**
  - Same as hit ...

34

## Miss handling: Software-controlled TLB

- On a TLB hit, MMU checks the valid bit
  - If valid, perform address translation
  - If invalid (e.g. page not in memory), MMU generates a page fault
    - If page not valid, OS performs page fault handling
    - Restart the faulting instruction
- On a TLB miss, HW raises exception, **traps to the OS**
  - OS parses page table and loads PTE into TLB
    - Needs to replace if TLB is full
    - PT layout is **flexible**
  - Same as in a hit...

35

## Hardware vs. software controlled

- Hardware approach
  - Efficient -- TLB misses handled by hardware
  - OS intervention is required only in case of page fault
  - Page structure prescribed by MMU hardware -- rigid
- Software approach
  - Less efficient -- TLB misses are handled by software
  - MMU hardware very simple, permitting larger, faster TLB
  - OS designer has complete flexibility in choice of MM data structure
    - e.g. 2-level page table, inverted page table)

36

## Deep thinking

- Without TLB, how MMU finds PTE is fixed
- With TLB, it can be flexible, e.g. software-controlled is possible
- What enables this?

37

## More TLB Issues

- Which TLB entry should be replaced?
  - Random
  - LRU
- What happens when changing a page table entry (e.g. because of swapping, change read/write permission)?
  - change the entry in memory
  - flush (eg. invalidate) the TLB entry
    - INGLPG on x86

38

## What happens to TLB in a process context switch?

- During a process context switch, cached translations can not be used by the next process
  - Invalidate all entries during a context switch
    - Lots of TLB misses afterwards
- Tag each entry with an ASID
  - Add a HW register that contains the process id of the current executing process
  - TLB hits if an entry's process id matches that reg

39

## Cache vs. TLB

- Similarities:
  - Both cache a part of the physical memory
- Differences:
  - Associativity
    - TLB is usually fully associative
    - Cache can be direct mapped
  - Consistency
    - TLB does not deal with consistency with memory
    - TLB can be controlled by software

40

## More on consistency Issues

- Snoopy cache protocols can maintain consistency with DRAM, even in the presence of DMA
- No hardware maintains consistency between DRAM and TLBs:
  - OS needs to flush related TLBs whenever changing a page table entry in memory
- On multiprocessors, when you modify a page table entry, you need to do “*TLB shoot-down*” to flush all related TLB entries on all processors

41

## Midterm topics

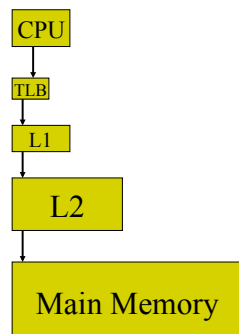
- Processes
- Threads
- Synchronization (wait-free sync)
- IPC with messages
- CPU scheduling
- Deadlocks
- Memory management (including today lec)
- Extra-level of indirection!

42

## Memory Hierarchy Revisited

What does this imply about L1 addresses?

Where do we hope requests get satisfied?

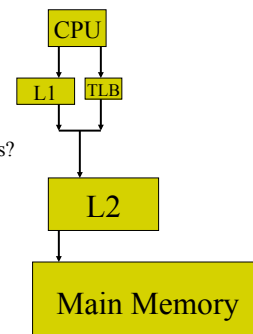


43

## Memory Hierarchy Re-Revisited

What does this imply about L1 addresses?

Any speed benefits?  
Any drawbacks?



44

## Reading

- Chapter 8

