

CS 202 Advanced Operating Systems

Virtual Memory (cont'd)

Elephant(s) in the room



- Problem 1: Translation is slow!
 - Many memory accesses for each memory access
 - Caches are useless!



"Unfortunately, there's another elephant in the room."

Speeding up Translation with a TLB



- Page table entries (PTEs) are cached in L1 like any other memory word
 - > PTEs may be evicted by other data references
 - > PTE hit still requires a small L1 delay
- Solution: Translation Lookaside Buffer (TLB)
 - Small hardware cache in MMU
 - Maps virtual page numbers to physical page numbers
 - Contains complete page table entries for small number of pages

TLB Hit





A TLB hit eliminates a memory access

TLB Miss





A TLB miss incurs an additional memory access (the PTE) Fortunately, TLB misses are rare. Why?

Reloading the TLB



- > If the TLB does not have mapping, two possibilities:
 - 1. MMU loads PTE from page table in memory
 - > Hardware managed TLB, OS not involved in this step
 - OS has already set up the page tables so that the hardware can access it directly
 - 2. Trap to the OS
 - > Software managed TLB, OS intervenes at this point
 - > OS does lookup in page table, loads PTE into TLB
 - > OS returns from exception, TLB continues
- > A machine will only support one method or the other
- > At this point, there is a PTE for the address in the TLB

Page Faults



- > PTE can indicate a protection fault
 - Read/write/execute operation not permitted on page
 - Invalid virtual page not allocated, or page not in physical memory
- > TLB traps to the OS (software takes over)
 - R/W/E OS usually will send fault back up to process, or might be playing games (e.g., copy on write, mapped files)
 - Invalid
 - Virtual page not allocated in address space
 - > OS sends fault to process (e.g., segmentation fault)
 - Page not in physical memory
 - > OS allocates frame, reads from disk, maps PTE to physical frame

Multi-Level Page Tables

- > Suppose:
 - > 4KB (2¹²) page size, 48-bit address space, 8-byte PTE
- > Problem:
 - Would need a 512 GB page table!
 - > $2^{48} * 2^{-12} * 2^3 = 2^{39}$ bytes
- Common solution:
 - Multi-level page tables
 - > Example: 2-level page table
 - Level 1 table: each PTE points to a page table (always memory resident)
 - Level 2 table: each PTE points to a page (paged in and out like any other data)





A Two-Level Page Table Hierarchy







Page Replacement

Mapped Files



- Mapped files enable processes to do file I/O using loads and stores
 - Instead of "open, read into buffer, operate on buffer, …"
- > Bind a file to a virtual memory region (mmap() in Unix)
 - > PTEs map virtual addresses to physical frames holding file data
 - Virtual address base + N refers to offset N in file
- Initially, all pages mapped to file are invalid
 - OS reads a page from file when invalid page is accessed
 - > OS writes a page to file when evicted, or region unmapped
 - > If page is not dirty (has not been written to), no write needed
 - > Another use of the dirty bit in PTE

Demand Paging (OS)



- We use demand paging (similar to other caches):
 - Pages loaded from disk when referenced
 - Pages may be evicted to disk when memory is full
 - Page faults trigger paging operations
- > What is the alternative to demand paging?
 - Some kind of prefetching
- Lazy vs. aggressive policies in systems

Demand Paging (Process)



- > Demand paging when a process first starts up
- > When a process is created, it has
 - A brand-new page table with all valid bits off
 - > No pages in memory
- > When the process starts executing
 - > Instructions fault on code and data pages
 - Faulting stops when all necessary code and data pages are in memory
 - > Only code and data needed by a process needs to be loaded
 - > This, of course, changes over time...

Page replacement policy



- What we discussed so far (page faults, swap, page table structures, etc...) is mechanisms
- Page replacement policy: determine which page to remove when we need a victim
- > Does it matter?
 - > Yes! Page faults are super expensive
 - Getting the number down, can improve the performance of the system significantly

Evicting the Best Page



- > Goal is to reduce the page fault rate
- > The best page to evict is the one never touched again
 - Will never fault on it
- Never is a long time, so picking the page closest to "never" is the next best thing
 - Evicting the page that won't be used for the longest period minimizes the number of page faults
 - > Proved by Belady
- We're now going to survey various replacement algorithms, starting with Belady's

Belady's Algorithm

UCR

- > Belady's algorithm
 - Idea: Replace the page that will not be used for the longest time in the future
 - > Optimal? How would you show?
 - Problem: Have to predict the future
- > Why is Belady's useful then?
 - > Use it as a yardstick/upper bound
 - Compare implementations of page replacement algorithms with the optimal to gauge room for improvement
 - > If optimal is not much better, then algorithm is pretty good
 - > What's a good lower bound?
 - > Random replacement is often the lower bound

First-In First-Out (FIFO)



- > FIFO is an obvious algorithm and simple to implement
 - > Maintain a list of pages in order in which they were paged in
 - > On replacement, evict the one brought in longest time ago
- > Why might this be good?
 - > Maybe the one brought in the longest ago is not being used
- > Why might this be bad?
 - Then again, maybe it's not
 - > We don't have any info to say one way or the other
- > FIFO suffers from "Belady's Anomaly"
 - The fault rate might actually increase when the algorithm is given more memory (very bad)

Least Recently Used (LRU)



- LRU uses reference information to make a more informed replacement decision
 - Idea: We can't predict the future, but we can make a guess based upon past experience
 - On replacement, evict the page that has not been used for the longest time in the past (Belady's: future)
 - When does LRU do well? When does LRU do poorly?
- Implementation
 - To be perfect, need to time stamp every reference (or maintain a stack) – much too costly
 - So we need to approximate it

Approximating LRU



- > LRU approximations use the PTE reference bit
 - > Keep a counter for each page
 - > At regular intervals, for every page do:
 - If ref bit = 0, increment counter
 - If ref bit = 1, zero the counter
 - > Zero the reference bit
 - The counter will contain the number of intervals since the last reference to the page
 - The page with the largest counter is the least recently used
- > Some architectures don't have a reference bit
 - Can simulate reference bit using the valid bit to induce faults

LRU Clock (Not Recently Used)

- Not Recently Used (NRU) Used by Unix
 - Replace page that is "old enough"
 - > Arrange all of physical page frames in a big circle (clock)
 - > A clock hand is used to select a good LRU candidate
 - > Sweep through the pages in circular order like a clock
 - > If the ref bit is off, it hasn't been used recently
 - > What is the minimum "age" if ref bit is off?
 - > If the ref bit is on, turn it off and go to next page
 - > Arm moves quickly when pages are needed
 - Low overhead when plenty of memory
 - > If memory is large, "accuracy" of information degrades
 - > What does it degrade to?
 - > One fix: use two hands (leading erase hand, trailing select hand)

LRU Clock





Example: gcc Page Replace



Example: Belady's Anomaly



0

Fixed vs. Variable Space



- In a multiprogramming system, we need a way to allocate memory to competing processes
- > Problem: How to determine how much memory to give to each process?
 - Fixed space algorithms
 - > Each process is given a limit of pages it can use
 - > When it reaches the limit, it replaces from its own pages
 - Local replacement
 - > Some processes may do well while others suffer
 - Variable space algorithms
 - > Process' set of pages grows and shrinks dynamically
 - Global replacement
 - > One process can ruin it for the rest

Working Set Model



- A working set of a process is used to model the dynamic locality of its memory usage
 - > Defined by Peter Denning in 60s
- > Definition
 - WS(t,w) = {set of pages P, such that every page in P was referenced in the time interval (t, t-w)}
 - t time, w working set window (measured in page refs)
- A page is in the working set (WS) only if it was referenced in the last w references

Working Set Size



- The working set size is the number of pages in the working set
 - > The number of pages referenced in the interval (t, t-w)
- > The working set size changes with program locality
 - > During periods of poor locality, you reference more pages
 - Within that period, the working set size is larger
- Intuitively, want the working set to be the set of pages a process needs in memory to prevent heavy faulting
 - Each process has a parameter w that determines a working set with few faults
 - > Denning: Don't run a process unless working set is in memory

Example: gcc Working Set





Working Set Problems



- > Problems
 - > How do we determine w?
 - > How do we know when the working set changes?
- Too hard to answer
 - So, working set is not used in practice as a page replacement algorithm
- > However, it is still used as an abstraction
 - > The intuition is still valid
 - When people ask, "How much memory does Firefox need?", they are in effect asking for the size of Firefox's working set

Thrashing



- Page replacement algorithms avoid thrashing
 - When most of the time is spent by the OS in paging data back and forth from disk
 - No time spent doing useful work (making progress)
 - > In this situation, the system is overcommitted
 - No idea which pages should be in memory to reduce faults
 - Could just be that there isn't enough physical memory for all of the processes in the system
 - > Ex: Running Windows95 with 4 MB of memory...
 - Possible solutions
 - Swapping write out all pages of a process
 - > Buy more memory

Summary



- > Page replacement algorithms
 - Belady's optimal replacement (minimum # of faults)
 - > FIFO replace page loaded furthest in past
 - LRU replace page referenced furthest in past
 - > Approximate using PTE reference bit
 - LRU Clock replace page that is "old enough"
 - Working Set keep the set of pages in memory that has minimal fault rate (the "working set")
 - Page Fault Frequency grow/shrink page set as a function of fault rate
- Multiprogramming
 - > Should a process replace its own page, or that of another?