

CS162 Operating Systems and Systems Programming Lecture 11

Page Allocation and Replacement

February 27, 2012
Anthony D. Joseph and Ion Stoica
<http://inst.eecs.berkeley.edu/~cs162>

Goals for Today

- Page Replacement Policies
 - FIFO, LRU
 - Clock Algorithm
- Working Set/Thrashing

Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Many slides generated from my lecture notes by Kubiawicz.

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.2

Memory Topics (61C and 162)

Topic	61C	162
Protection	HW-based address spaces	HW-based address spaces, SW-based strong typing, SW fault isolation
Virtual Memory	Base & bound and single-level paging approaches	Base & bound, swapping, multiple segments, paging, multi-level, inverted page table approaches. Shared memory and msg passing.
Caching	Mem hierarchy, temporal/spatial locality, 3 sources of cache misses, direct/associative caches, write-through/back policies, access time calc	Mem hierarchy, temporal/spatial locality, five sources of cache misses, direct/associative caches, write-through/back policies, access time calc, context switch implications
TLB	Basic concept	Concept and overlapped with cache
Paging	Overview, LRU algorithm	Detailed steps, analyzing algorithms, approximating LRU, implementing second-chance and nth chance algorithms, working sets

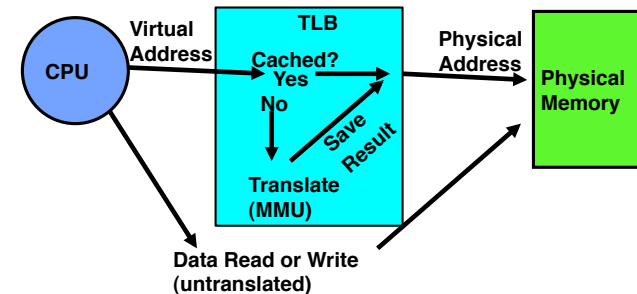
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.3

Review: Caching Applied to Address Translation

- Problem: address translation expensive (especially multi-level)
- Solution: cache address translation (TLB)
 - Instruction accesses spend a lot of time on the same page (since accesses sequential)
 - Stack accesses have definite locality of reference
 - Data accesses have less page locality, but still some...



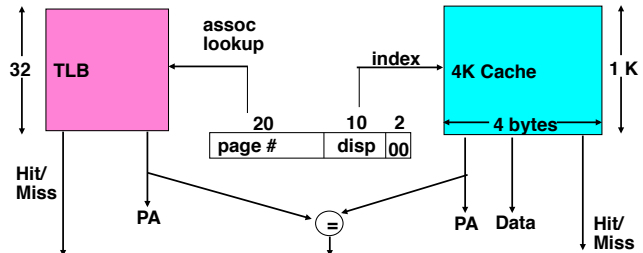
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.4

Overlapping TLB & Cache Access

- Here is how this might work with a 4K cache:



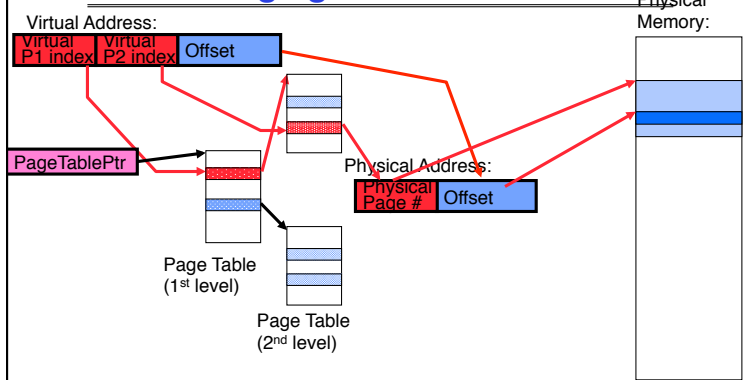
- What if cache size is increased to 8KB?
 - Overlap not complete
 - Need to do something else. See CS152/252
- Another option: Virtual Caches
 - Tags in cache are virtual addresses
 - Translation only happens on cache misses

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.5

Review: Paging & Address Translation

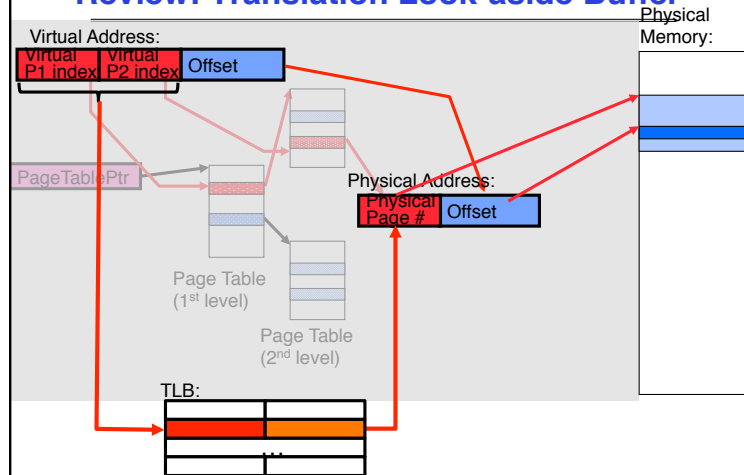


2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.6

Review: Translation Look-aside Buffer

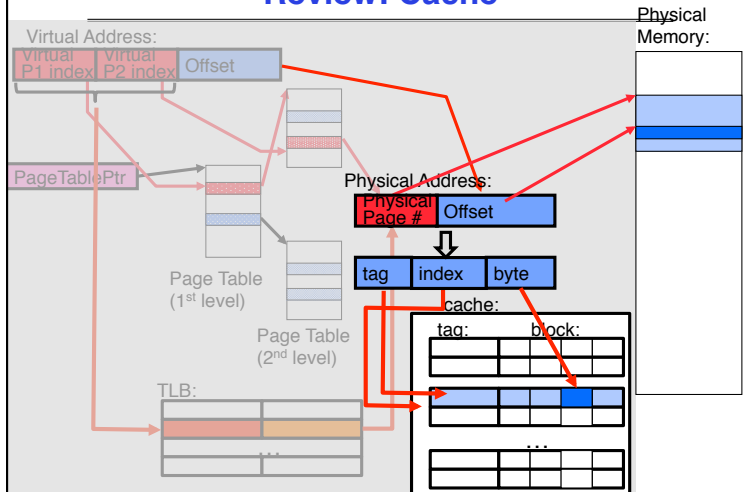


2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.7

Review: Cache



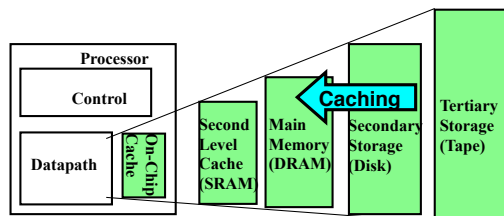
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.8

Demand Paging

- Modern programs require a lot of physical memory
 - Memory per system growing faster than 25%-30%/year
- But they don't use all their memory all of the time
 - 90-10 rule: programs spend 90% of their time in 10% of their code
 - Wasteful to require all of user's code to be in memory
- Solution: use main memory as cache for disk



2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.9

Demand Paging is Caching

- Since Demand Paging is Caching, we must ask:

Question	Choice
What is the block size?	
What is the organization of this cache (i.e., direct-mapped, set-associative, fully-associative)?	
How do we find a page in the cache?	
What is page replacement policy? (i.e., LRU, Random, ...)	
What happens on a miss?	
What happens on a write? (i.e., write-through, write-back)	

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.10

Demand Paging Mechanisms

- PTE helps us implement demand paging
 - Valid \Rightarrow Page in memory, PTE points at physical page
 - Not Valid \Rightarrow Page not in memory; use info in PTE to find it on disk when necessary

- Suppose user references page with invalid PTE?

- Memory Management Unit (MMU) traps to OS
 - » Resulting trap is a "Page Fault"

- What does OS do on a Page Fault?:

- » Choose an old page to replace
- » If old page modified ("D=1"), write contents back to disk
- » Change its PTE and any cached TLB to be invalid
- » Load new page into memory from disk
- » Update page table entry, invalidate TLB for new entry
- » Continue thread from original faulting location

- TLB for new page will be loaded when thread continued!

- While pulling pages off disk for one process, OS runs another process from ready queue

- » Suspended process sits on wait queue

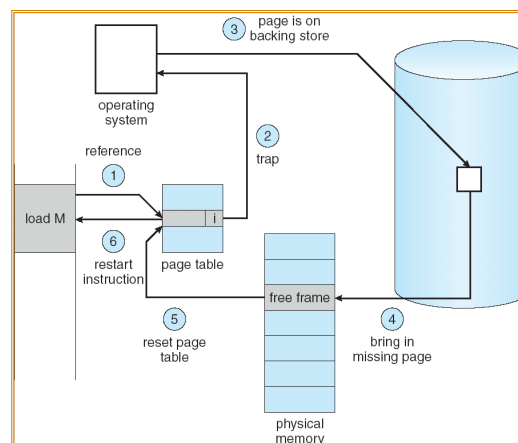
Cache

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.11

Steps in Handling a Page Fault



2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.12

Demand Paging Example

- Since Demand Paging like caching, can compute average access time! ("Effective Access Time")
 - $EAT = \text{Hit Rate} \times \text{Hit Time} + \text{Miss Rate} \times \text{Miss Time}$
- Example:
 - Memory access time = 200 nanoseconds
 - Average page-fault service time = 8 milliseconds
 - Suppose $p = \text{Probability of miss}$, $1-p = \text{Probably of hit}$
 - Then, we can compute EAT as follows:
$$EAT = (1 - p) \times 200\text{ns} + p \times 8 \text{ ms}$$
$$= (1 - p) \times 200\text{ns} + p \times 8,000,000\text{ns}$$
$$= 200\text{ns} + p \times 7,999,800\text{ns}$$
- If one access out of 1,000 causes a page fault, then $EAT = 8.2 \mu\text{s}$:
 - This is a slowdown by a factor of 40!
- What if want slowdown by less than 10%?
 - $200\text{ns} \times 1.1 < EAT \Rightarrow p < 2.5 \times 10^{-6}$
 - This is about 1 page fault in 400,000 !

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.13

What Factors Lead to Misses?

- **Compulsory Misses:**
 - Pages that have never been paged into memory before
 - How might we remove these misses?
 - » Prefetching: loading them into memory before needed
 - » Need to predict future somehow! More later.
- **Capacity Misses:**
 - Not enough memory. Must somehow increase size.
 - Can we do this?
 - » One option: Increase amount of DRAM (not quick fix!)
 - » Another option: If multiple processes in memory: adjust percentage of memory allocated to each one!
- **Conflict Misses:**
 - Technically, conflict misses don't exist in virtual memory, since it is a "fully-associative" cache
- **Policy Misses:**
 - Caused when pages were in memory, but kicked out prematurely because of the replacement policy
 - How to fix? Better replacement policy

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.14

Page Replacement Policies

- Why do we care about Replacement Policy?
 - Replacement is an issue with any cache
 - Particularly important with pages
 - » The cost of being wrong is high: must go to disk
 - » Must keep important pages in memory, not toss them out
- **FIFO (First In, First Out)**
 - Throw out oldest page. Be fair – let every page live in memory for same amount of time.
 - Bad, because throws out heavily used pages instead of infrequently used pages
- **MIN (Minimum):**
 - Replace page that won't be used for the longest time
 - Great, but can't really know future...
 - Makes good comparison case, however
- **RANDOM:**
 - Pick random page for every replacement
 - Typical solution for TLB's. Simple hardware
 - Unpredictable

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.15

Replacement Policies (Con't)

- **LRU (Least Recently Used):**
 - Replace page that hasn't been used for the longest time
 - Programs have locality, so if something not used for a while, unlikely to be used in the near future.
 - Seems like LRU should be a good approximation to MIN.
- How to implement LRU? Use a list!

```
graph LR; Head --> P6[Page 6]; P6 --> P7[Page 7]; P7 --> P1[Page 1]; P1 --> P2[Page 2]; Tail["Tail (LRU)"] --> P2;
```

 - On each use, remove page from list and place at head
 - LRU page is at tail

- Problems with this scheme for paging?
 - List operations complex
 - » Many instructions for each hardware access
- In practice, people **approximate** LRU (more later)

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.16

Example: FIFO

- Suppose we have 3 page frames, 4 virtual pages, and following reference stream:
 - A B C A B D A D B C B
- Consider FIFO Page replacement:

Ref:	A	B	C	A	B	D	A	D	B	C	B
Page:											
1	A					D				C	
2		B					A				
3			C						B		

- FIFO: 7 faults.
- When referencing D, replacing A is bad choice, since need A again right away

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.17

Example: MIN

- Suppose we have the same reference stream:
 - A B C A B D A D B C B
- Consider MIN Page replacement:

Ref:	A	B	C	A	B	D	A	D	B	C	B
Page:											
1	A									C	
2		B									
3			C			D					

- MIN: 5 faults
- Look for page not referenced farthest in future.
- What will LRU do?
 - Same decisions as MIN here, but won't always be true!

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.18

When will LRU perform badly?

- Consider the following: A B C D A B C D A B C D
- LRU Performs as follows (same as FIFO here):

Ref:	A	B	C	D	A	B	C	D	A	B	C	D
Page:												
1	A			D			C			B		
2		B			A			D			C	
3			C			B			A			D

– Every reference is a page fault!

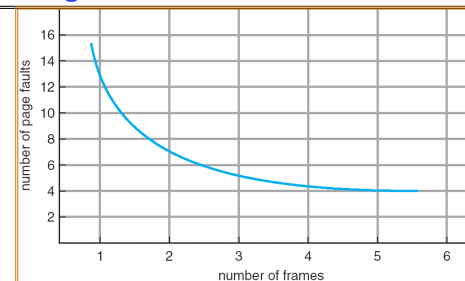
- MIN Does much better:

Ref:	A	B	C	D	A	B	C	D	A	B	C	D
Page:												
1	A									B		
2		B					C					
3			C	D								

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

Graph of Page Faults Versus The Number of Frames



- One desirable property: When you add memory the miss rate goes down
 - Does this always happen?
 - Seems like it should, right?
- No: Belady's anomaly
 - Certain replacement algorithms (FIFO) don't have this obvious property!

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.20

Adding Memory Doesn't Always Help Fault Rate

- Does adding memory reduce number of page faults?
 - Yes for LRU and MIN
 - Not necessarily for FIFO! (Called Belady's anomaly)

Page:	A	B	C	D	A	B	E	A	B	C	D	E
1	A			D			E					
2		B			A					C		
3			C			B					D	

Page:	A	B	C	D	A	B	E	A	B	C	D	E
1	A						E				D	
2		B						A				E
3			C						B			
4				D						C		

- After adding memory:
 - With FIFO, contents can be completely different
 - In contrast, with LRU or MIN, contents of memory with X pages are a subset of contents with X+1 Page

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.21

Administrivia

- All deadlines are on lectures page: <http://inst.eecs.berkeley.edu/~cs162/sp12/lectures.html>
 - Project 2 Design Doc due Thursday 3/1 at 11:59PM
 - Midterm next Wednesday 3/7 at 5-6:30PM in **10 Evans**
 - Closed-book, 1 double-sided page of handwritten notes
 - Covers lectures/readings #1-12 (Wed 3/1) and project one
 - Midterm review session: Monday 3/5 7-9PM in 141 McCone
- Please fill the anonymous course survey at <https://www.surveymonkey.com/s/DZ5Y8XM>
 - We'll try to make changes *this* semester based on your feedback
- New lecture schedule (topic order) posted:
 - I/O, File Systems, Key-Value Stores, Networking, Fault-Tolerant Computing, Security, Capstones

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.22

5min Break

Implementing LRU & Second Chance

- Perfect:
 - Timestamp page on each reference
 - Keep list of pages ordered by time of reference
 - Too expensive to implement in reality for many reasons
- Second Chance Algorithm:**
 - Approximate LRU
 - Replace **an** old page, not **the oldest** page
 - FIFO with "use" bit
- Details
 - A "use" bit per physical page
 - On page fault check page at head of queue
 - If use bit=1 → clear bit, and move page at tail (give the page second chance!)
 - If use bit=0 → replace page
 - Moving pages to tail still complex

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.24

Clock Algorithm

- **Clock Algorithm:** more efficient implementation of second chance algorithm
 - Arrange physical pages in circle with single clock hand
- Details:
 - On page fault:
 - » Check use bit: 1→used recently; clear and leave it alone
 - » 0→selected candidate for replacement
 - » Advance clock hand (not real time)
 - Will always find a page or loop forever?
- What if hand moving slowly?
 - Good sign or bad sign?
 - » Not many page faults and/or find page quickly
- What if hand is moving quickly?
 - Lots of page faults and/or lots of reference bits set



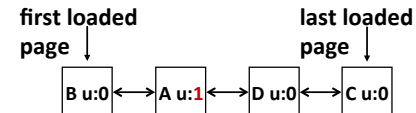
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.25

Second Chance Illustration

- Max page table size 4
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives



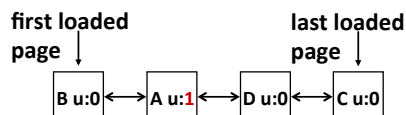
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.26

Second Chance Illustration

- Max page table size 4
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives
 - Page F arrives



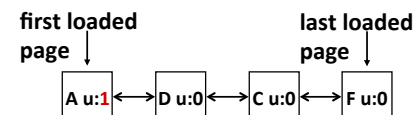
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.27

Second Chance Illustration

- Max page table size 4
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives
 - Page F arrives



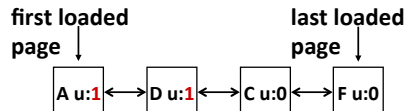
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.28

Second Chance Illustration

- Max page table size 4
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives
 - Page F arrives
 - Access page D
 - Page E arrives



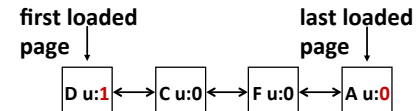
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.29

Second Chance Illustration

- Max page table size 4
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives
 - Page F arrives
 - Access page D
 - Page E arrives



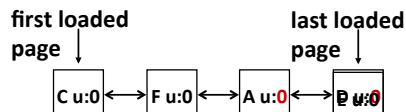
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.30

Second Chance Illustration

- Max page table size 4
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives
 - Page F arrives
 - Access page D
 - Page E arrives



2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.31

Clock Replacement Illustration

- Max page table size 4
- Invariant: point at oldest page
 - Page B arrives



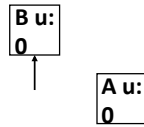
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.32

Clock Replacement Illustration

- Max page table size 4
- Invariant: point at oldest page
 - Page B arrives
 - Page A arrives
 - Access page A



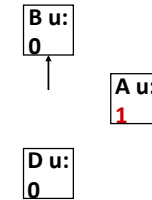
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.33

Clock Replacement Illustration

- Max page table size 4
- Invariant: point at oldest page
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives



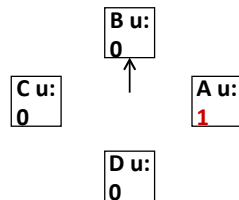
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.34

Clock Replacement Illustration

- Max page table size 4
- Invariant: point at oldest page
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives



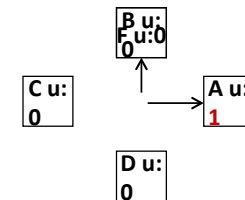
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.35

Clock Replacement Illustration

- Max page table size 4
- Invariant: point at oldest page
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives
 - Page F arrives



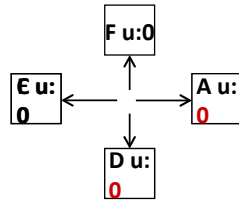
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.36

Clock Replacement Illustration

- Max page table size 4
- Invariant: point at oldest page
 - Page B arrives
 - Page A arrives
 - Access page A
 - Page D arrives
 - Page C arrives
 - Page F arrives
 - Access page D
 - Page E arrives



2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.37

Nth Chance version of Clock Algorithm

- **Nth chance algorithm:** Give page N chances
 - OS keeps counter per page: # sweeps
 - On page fault, OS checks use bit:
 - » 1 ⇒ clear use and also clear counter (used in last sweep)
 - » 0 ⇒ increment counter; if count=N, replace page
 - Means that clock hand has to sweep by N times without page being used before page is replaced
- How do we pick N?
 - Why pick large N? Better approx to LRU
 - » If N ~ 1K, really good approximation
 - Why pick small N? More efficient
 - » Otherwise might have to look a long way to find free page
- What about dirty pages?
 - Takes extra overhead to replace a dirty page, so give dirty pages an extra chance before replacing?
 - Common approach:
 - » Clean pages, use N=1
 - » Dirty pages, use N=2 (and write back to disk when N=1)

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.38

Clock Algorithms: Details

- Which bits of a PTE entry are useful to us?
 - **Use:** Set when page is referenced; cleared by clock algorithm
 - **Modified:** set when page is modified, cleared when page written to disk
 - **Valid:** ok for program to reference this page
 - **Read-only:** ok for program to read page, but not modify
 - » For example for catching modifications to code pages!
- Do we really need hardware-supported “modified” bit?
 - No. Can emulate it (BSD Unix) using read-only bit
 - » Initially, mark all pages as read-only, even data pages
 - » On write, trap to OS. OS sets software “modified” bit, and marks page as read-write.
 - » Whenever page comes back in from disk, mark read-only

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.39

Clock Algorithms Details (cont'd)

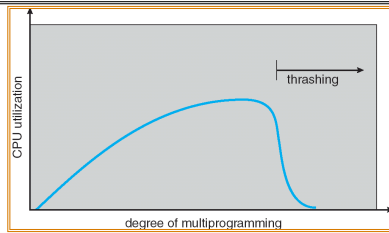
- Do we really need a hardware-supported “use” bit?
 - No. Can emulate it using “invalid” bit:
 - » Mark all pages as invalid, even if in memory
 - » On read to invalid page, trap to OS
 - » OS sets use bit, and marks page read-only
 - When clock hand passes by, reset use bit and mark page as invalid again

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.40

Thrashing



- If a process does not have “enough” pages, the page-fault rate is very high. This leads to:
 - low CPU utilization
 - operating system spends most of its time swapping to disk
- **Thrashing** = a process is busy swapping pages in and out
- Questions:
 - How do we detect Thrashing?
 - What is best response to Thrashing?

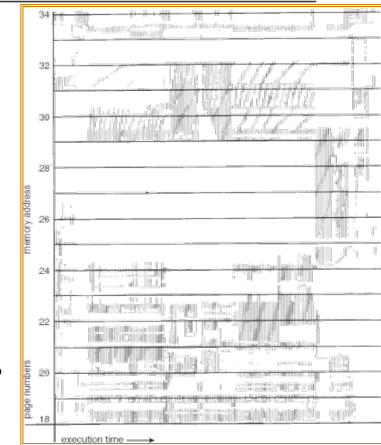
2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.41

Locality In A Memory-Reference Pattern

- Program Memory Access Patterns have temporal and spatial locality
 - Group of Pages accessed along a given time slice called the “Working Set”
 - Working Set defines minimum number of pages needed for process to behave well
- Not enough memory for Working Set \Rightarrow Thrashing
 - Better to swap out process?

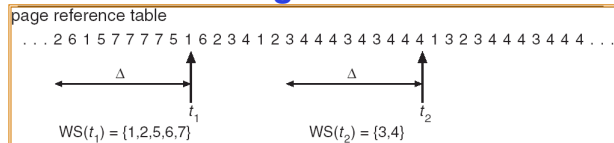


2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.42

Working-Set Model



- Δ = working-set window = fixed number of page references
 - Example: 10,000 instructions
- WS_i (working set of Process P_i) = total set of pages referenced in the most recent Δ (varies in time)
 - if Δ too small will not encompass entire locality
 - if Δ too large will encompass several localities
 - if $\Delta = \infty \Rightarrow$ will encompass entire program
- $D = \sum |WS_i|$ = total demand frames
- if $D > \text{memory} \Rightarrow$ Thrashing
 - Policy: if $D > \text{memory}$, then suspend/swap out processes
 - This can improve overall system behavior by a lot!

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.43

What about Compulsory Misses?

- Recall that compulsory misses are misses that occur the first time that a page is seen
 - Pages that are touched for the first time
 - Pages that are touched after process is swapped out/swapped back in
- **Clustering:**
 - On a page-fault, bring in multiple pages “around” the faulting page
 - Since efficiency of disk reads increases with sequential reads, makes sense to read several sequential pages
- **Working Set Tracking:**
 - Use algorithm to try to track working set of application
 - When swapping process back in, swap in working set

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.44

Summary (1/2)

- Demand Paging:
 - Treat memory as cache on disk
 - Cache miss \Rightarrow find free page, get page from disk
- Transparent Level of Indirection
 - User program is unaware of activities of OS behind scenes
 - Data can be moved without affecting application correctness
- Replacement policies
 - FIFO: Place pages on queue, replace page at end
 - » Fair but can eject in-use pages, suffers from Belady's anomaly
 - MIN: Replace page that will be used farthest in future
 - » Benchmark for comparisons, can't implement in practice
 - LRU: Replace page used farthest in past
 - » For efficiency, use approximation

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.45

Summary (2/2)

- Clock Algorithm: Approximation to LRU
 - Arrange all pages in circular list
 - Sweep through them, marking as not "in use"
 - If page not "in use" for one pass, then can replace
- Second-Chance List algorithm: Yet another approx LRU
 - Divide pages into two groups, one of which is truly LRU and managed on page faults
- Working Set:
 - Set of pages touched by a process recently
- Thrashing: a process is busy swapping pages in and out
 - Process will thrash if working set doesn't fit in memory
 - Need to swap out a process

2/27/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

11.46