

CSC 369

Operating Systems

Lecture 7: Paged Virtual Memory



University of Toronto, Department of Computer Science



Today

LAST TIME:

- Hardware for dynamic relocation
- Paged address translation
- Simple linear page table & page table entries

TODAY

- Better page table designs – (saving space)
- Efficient translations (TLBs) - (saving time)
- Policies for memory management



Paging Limitations - Space

- Memory required for linear page table can be large (space overhead)
 - Need one PTE per page
 - 32 bit virtual address space w/ 4K pages $\Rightarrow 2^{20}$ PTEs
 - 4 bytes/PTE \Rightarrow 4MB/page table
 - 100 processes \Rightarrow 400MB just for page tables!
 - And modern processors have 64-bit address spaces
 \Rightarrow 16 petabytes per page table!
 - Solution 1: Hierarchical page tables
 - Solution 2: Hashed page tables
 - Solution 3: Inverted page tables



Managing Page Tables

- How can we reduce space overhead?
 - *Observation:* Only need to map the portion of the address space actually being used (tiny fraction of entire address space)
- How do we only map what is being used?
 - Use **3 page tables – one per logical segment** (code, stack, heap)
 - What do we need then?
 - Base register (where does the page table start in physical memory)
 - Bound/limit register (where does the page table end)
 - ⇒ **3 pairs of base/bound registers** (typically in the MMU)
 - Big memory savings for page tables. Why?
 - Is it ideal? See any issues?



Managing Page Tables

- Can still end up with lots of page table waste
 - e.g., if we have a large but sparse heap
- External fragmentation
 - Memory is managed in page-sized units, but page tables now can be of arbitrary size (in multiples of PTEs). So, finding free space for page tables in memory gets complicated
- => Hybrid approach won't work so well in practice
- What else?
 - Use another level of indirection



Multi-level page tables in a nutshell

- Linear page table
 - Page table contains **page table entries (PTEs)**
- Multi-level page tables
 - Split page table into **pages of PTEs**
 - Use a **page table directory** (or “master page table”) to point to different pages of PTEs



Intuition

Linear
page table



Unused space in page table is
like empty sheets.
So, how does this save us space?

Multi-level
page table





Comparison

Linear Page Table

PTBR 201

valid	prot	PFN	
1	rx	12	PFN 201
1	rx	13	
0	-	-	
1	rw	100	
0	-	-	PFN 202
0	-	-	
0	-	-	
0	-	-	
0	-	-	PFN 203
0	-	-	
0	-	-	
0	-	-	
0	-	-	PFN 204
0	-	-	
1	rw	86	
1	rw	15	

A valid bit indicates if a page is allocated. The page table can be divided into 4 pages.

Multi-level Page Table

PDBR 200

valid	PFN	
1	201	PFN 200
0	-	
0	-	
1	204	

The Page Directory

Valid bit in Page Directory == whole page of pages is not allocated.

Notice that some pages have all valid bits zeroes.

valid	prot	PFN	
1	rx	12	PFN 201
1	rx	13	
0	-	-	
1	rw	100	
[Page 1 of PT: Not Allocated]			
[Page 2 of PT: Not Allocated]			
0	-	-	PFN 204
0	-	-	
1	rw	86	
1	rw	15	

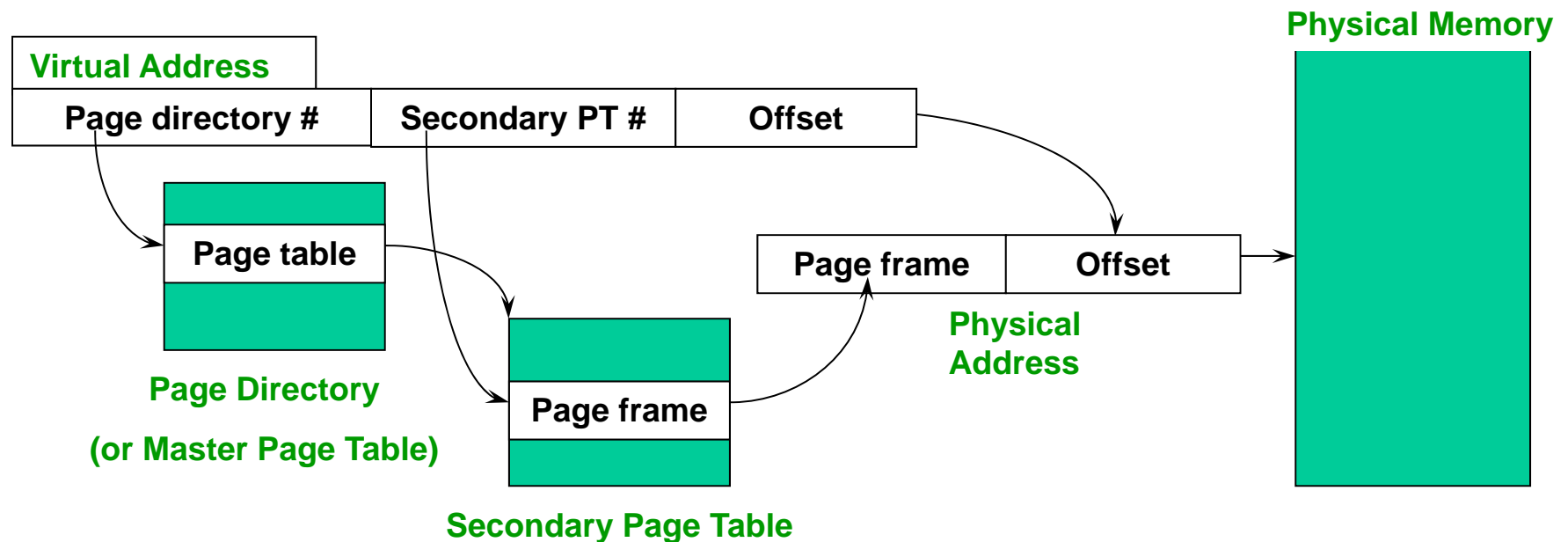
Source: the OSTEP textbook ;)



Two-Level Page Tables

Virtual addresses (VAs) have three parts:

- Page directory (Master page table) number, secondary page number, and offset
- A. Page directory maps VAs to secondary page table
- B. Secondary page table maps page number to physical frame
- C. Offset selects address within physical frame





2-Level Paging Example

- And now .. Some math!
- 32-bit virtual address space
 - 4K pages, 4 bytes/PTE
 - How many bits in offset?
 - $4K = 12$ bits, leaves 20 bits
 - Want master/secondary page tables in 1 page frame each:
 - $4K/4$ bytes = 1K entries. How many bits for master/secondary?
 - Master (1K) = 10, offset = 12, inner = $32 - 10 - 12 = 10$ bits
 - Note: this is why 4K is such a common page size on 32-bit architectures!

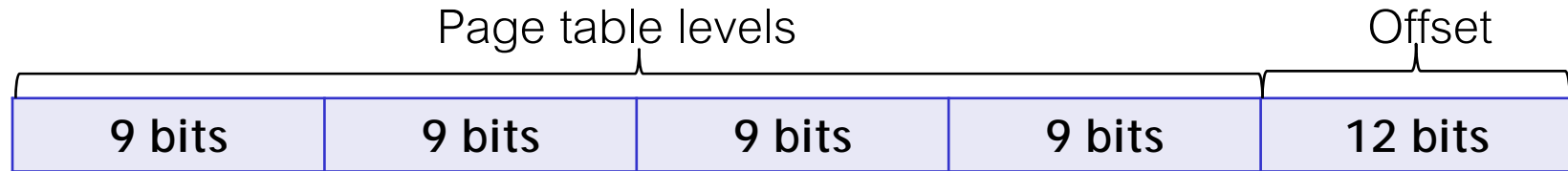


Tradeoff: space vs. time

- Multi-level page table
 - Saves space
 - Adds more levels of indirection when translating addresses
 - How many memory accesses on each translation, compared to linear?
 - Also more complexity in the translation algorithm
- FYI: Linux uses 3-level page tables!
- We'll see later how a TLB speeds up the “time” aspect



64-bit Address Space



- Really only using 48 bits
- Why? No need for more yet, wastes transistors.
- ISA supports 64-bit, but current CPUs only use lower 48-bits.
- Can be extended later to 64-bits without breaking compatibility.



640K ought to be enough
for anybody.
-Bill Gates, Microsoft 1981



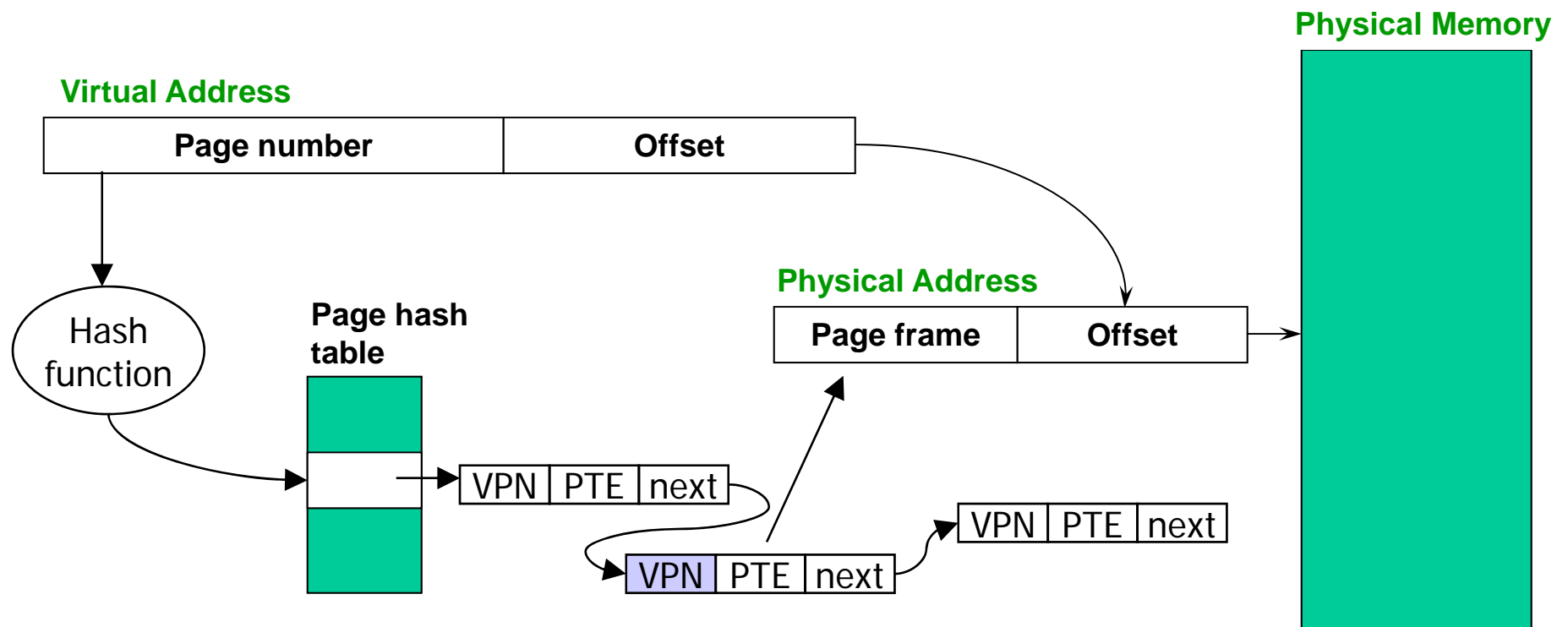
64-bit Address Spaces

- Suppose we just extended the hierarchical page tables with more levels
 - 4K pages \Rightarrow 12 bits for offset \Rightarrow 52 bits for page numbers
 - Assuming 8 bytes for PTE (common on 64-bit architectures)
 - \Rightarrow Maximum $4\text{KB}/8\text{B} = 512$ entries per level
 - \Rightarrow 9 bits for each level $\Rightarrow 52/9$ means 6 levels
 - Too much overhead!
 - Try bigger pages: 16K pages \Rightarrow 14 bits for offset, 50 bits for page numbers
 - \Rightarrow Maximum $16\text{K}/8\text{B} = 2\text{k}$ entries per level (pages are 4X larger)
 - \Rightarrow 11 bits for each level $\Rightarrow 50/11$ means 5 levels
 - Better, but still a lot of memory reads



Solution: Hashed Page Tables

- Hash function maps virtual page number (VPN) to bucket in fixed-size hash table
- Search entries at that bucket for matching VPN





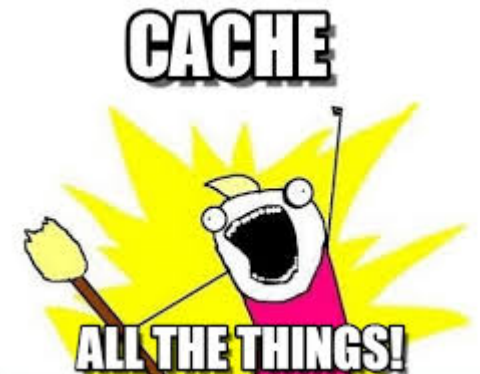
Inverted Page Tables

- So far ... one page table per process
- Instead, keep only one page table for the whole system
 - An entry for each physical page frame
 - Entries record which virtual page # is stored in that frame
 - Need to record process id as well
- Less space, but lookups are slower
 - References use virtual addresses, table is indexed by physical addresses
 - Use hashing to reduce the search time



Paging Limitations - Time

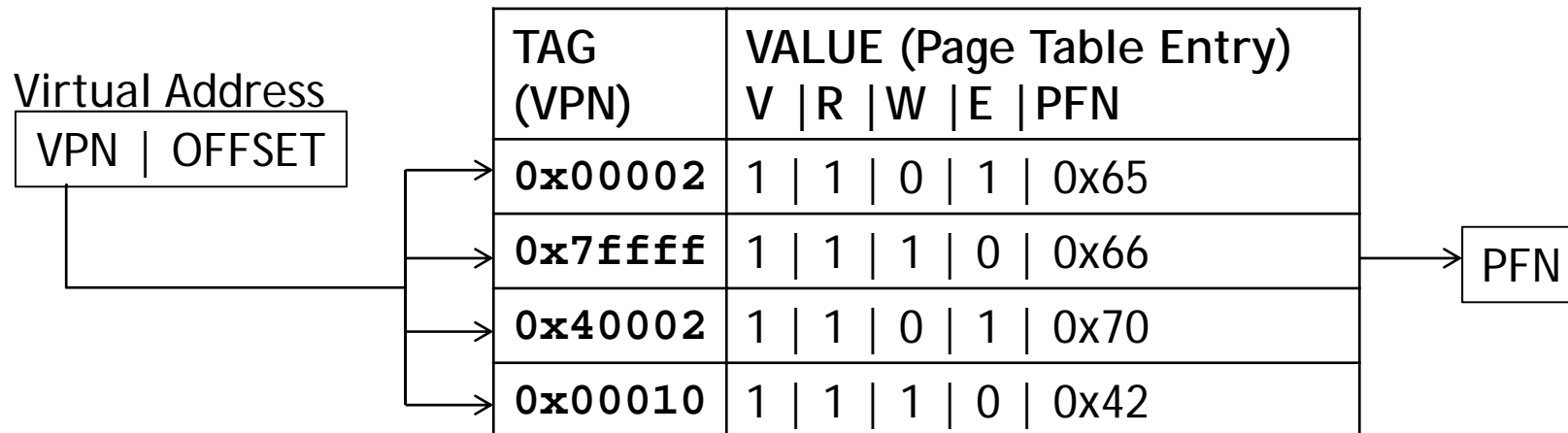
- Memory reference overhead (time)
 - 2 memory reads (references) per address lookup
 - First read page table, then actual memory
 - Hierarchical page tables make it worse
 - One read per level of page table
 - **Solution: use a hardware cache of lookups!**
- Translation Lookaside Buffer (TLB)
 - Small, fully-associative **hardware** cache of recently used translations
 - Part of the Memory Management Unit (MMU)





Translation Lookaside Buffer (TLB)

- Translates virtual page #s into PTEs (not physical addresses!)
 - Can be done in a single machine cycle
- TLBs are implemented in hardware
 - Fully associative cache (all entries looked up in parallel)
 - Cache tags are virtual page numbers
 - Cache values are PTEs (entries from page tables)
 - With PTE + offset, can directly calculate physical address





TLBs Exploit Locality

- Processes only use a handful of pages at a time
 - Only need those pages to be “mapped”
 - 16-64 entry TLB, so 16-64 pages can be mapped (64-256K)
- Hit rates are very important
 - typically >99% of translations are *hits*
- But we still need to deal with misses
 - Called a *TLB miss* or *TLB fault*



Managing TLBs

- Who places translations into the TLB (loads the TLB)?
 - **Hardware-loaded** (Memory Management Unit, e.g. Intel Pentium)
 - Knows where page tables are in main memory (PTBR)
 - OS maintains tables in memory, HW accesses them directly
 - Tables have to be in HW-defined format (inflexible)
 - **Software-loaded** TLB (OS, e.g. MIPS R2000)
 - TLB faults to the OS, OS finds appropriate PTE, loads it in TLB
 - Must be fast (but still 20-200 cycles)
 - CPU ISA has instructions for manipulating TLB
 - Add entry, read entry, invalidate entry, invalidate all entries (flush)
 - Tables can be in any format convenient for OS (flexible)

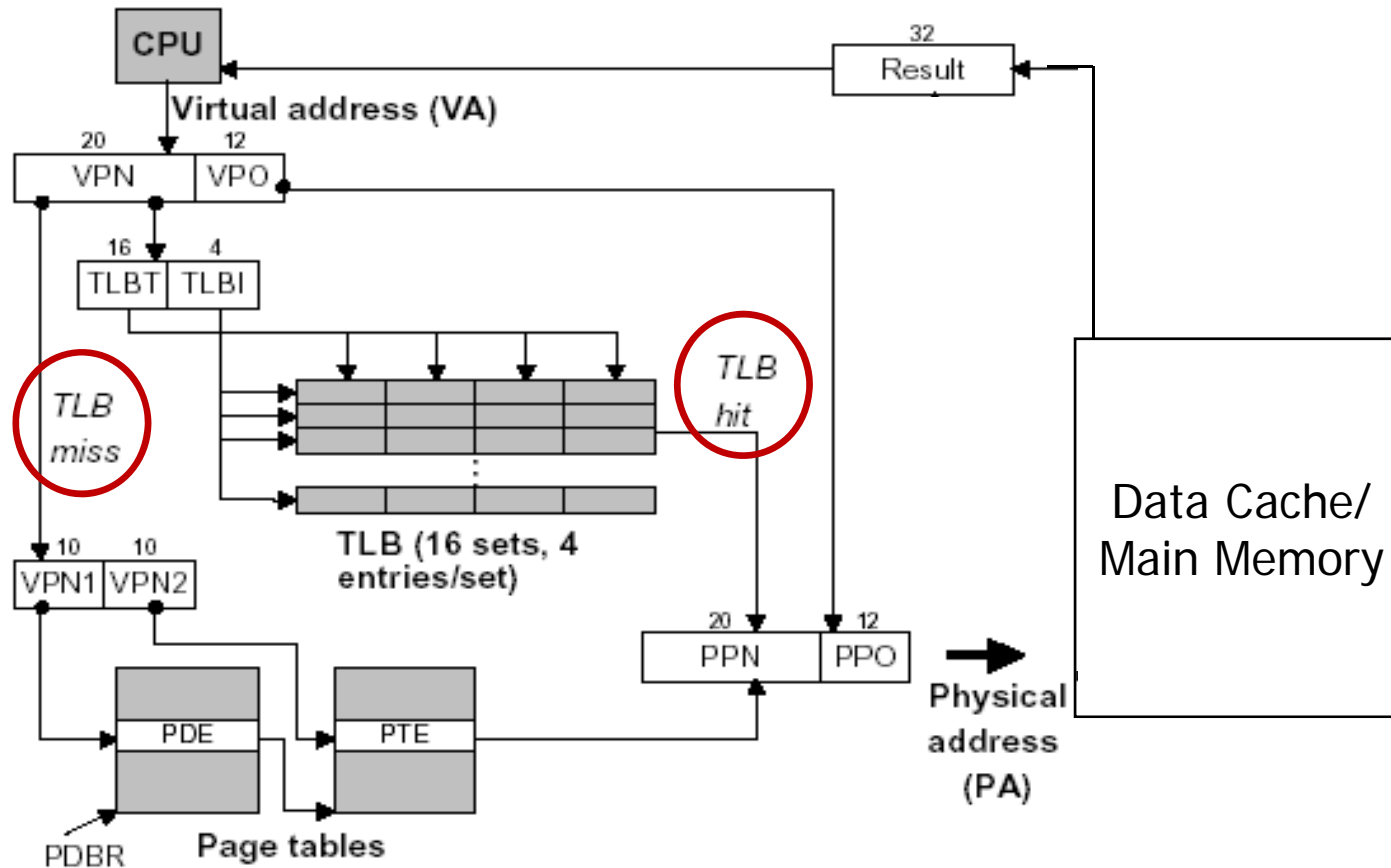


Managing TLBs (2)

- OS ensures that **TLB and page tables are consistent**
 - When it changes the protection bits of a PTE, it needs to update the copy of the PTE in the TLB (if one exists)
 - Typically, just **invalidate the entry**
- Reload TLB on a **process context switch**
 - Invalidate all entries ("**flush the TLB**")
 - Why?
 - What is one way to fix it?
- When the TLB misses and a new PTE has to be loaded, **a cached PTE must be evicted**
 - Choosing PTE to evict is called the **TLB replacement policy**
 - Implemented in hardware (for hardware loaded TLB), often simple (e.g., LRU)



Example: Pentium Address Translation





Paged Virtual Memory

- We've mentioned before that pages can be moved between memory and disk
 - This process is called **demand paging**
- OS uses main memory as a page cache of all the data allocated by processes in the system
 - Initially, pages are allocated from memory
 - When memory fills up, allocating a page in memory requires some other page to be evicted from memory
 - Evicted pages go to disk (where? **the swap file**)
 - The movement of pages between memory and disk is done by the OS, and is transparent to the application



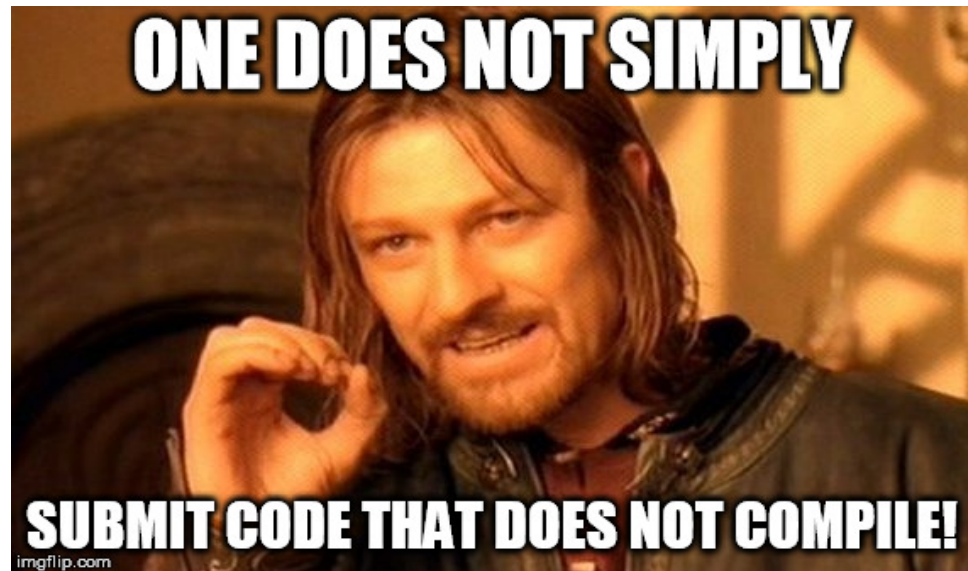
Page Faults

- What happens when a process accesses a page that has been evicted?
 1. When it evicts a page, the OS sets the PTE as invalid and stores the location of the page in the swap file in the PTE
 2. When a process accesses the page, the invalid PTE will cause a trap (page fault)
 3. The trap will run the OS page fault handler
 4. Handler uses the invalid PTE to locate page in swap file
 5. Reads page into a physical frame, updates PTE to point to it (what else?)
 6. Resumes the process
- But where does it put it? Each time we evict something else => slow!
 - OS usually keeps a pool of free pages around so that allocations do not always cause evictions



Assignment 1

- Common mistakes:
 - No synchronization on reads to shared data => must synchronize any access to shared data, *both* reads and writes!
 - Did not clear the list when demonitoring all, or when deintercepting the system call => on next intercept or monitor, residual stuff was in the list
 - Did not restore the intercepted system calls when exiting the module
 - Added starter_code/ directory, no INFO.txt file, etc.. => Read the submission instructions carefully!
 - Code did not compile!





Big picture

SO FAR:

- Optimizations:
 - 2-level page table (saving space)
 - Inverted page tables
 - Efficient translations (TLBs) - (saving time)

NEXT

- Policies for memory management



Remember: Locality

- All paging schemes depend on locality
 - Processes reference pages in localized patterns
- Temporal locality
 - Locations referenced recently likely to be referenced again.
- Spatial locality
 - Locations near recently referenced locations are likely to be referenced soon.
- Although the cost of paging is high, if it is infrequent enough it is acceptable
 - Processes usually exhibit both kinds of locality during their execution, making paging possible
 - All caching strategies depend on locality to be effective



Policy Decisions

- Page tables, MMU, TLB, etc. are *mechanisms* that make virtual memory possible
- Now, we'll look at *policies* for virtual memory management:
 - Fetch Policy – *when* to fetch a page
 - Demand paging vs. Prepaging
 - Placement Policy – *where* to put the page
 - Are some physical pages preferable to others?
 - Replacement Policy – *what* page to evict to make room?
 - Lots and lots of possible algorithms!



Demand Paging (OS)

- Demand paging from the OS perspective:
 - Pages are evicted to disk when memory is full
 - Pages loaded from disk when referenced again
 - References to evicted pages cause a TLB miss
 - PTE was invalid, causes fault
 - OS allocates a page frame, reads page from disk
 - When I/O completes, the OS fills in PTE, marks it valid, and resumes faulting process
- Dirty vs. clean pages
 - Actually, only dirty pages (modified) need to be written to disk
 - Clean pages do not – but you need to know where on disk to read them from again

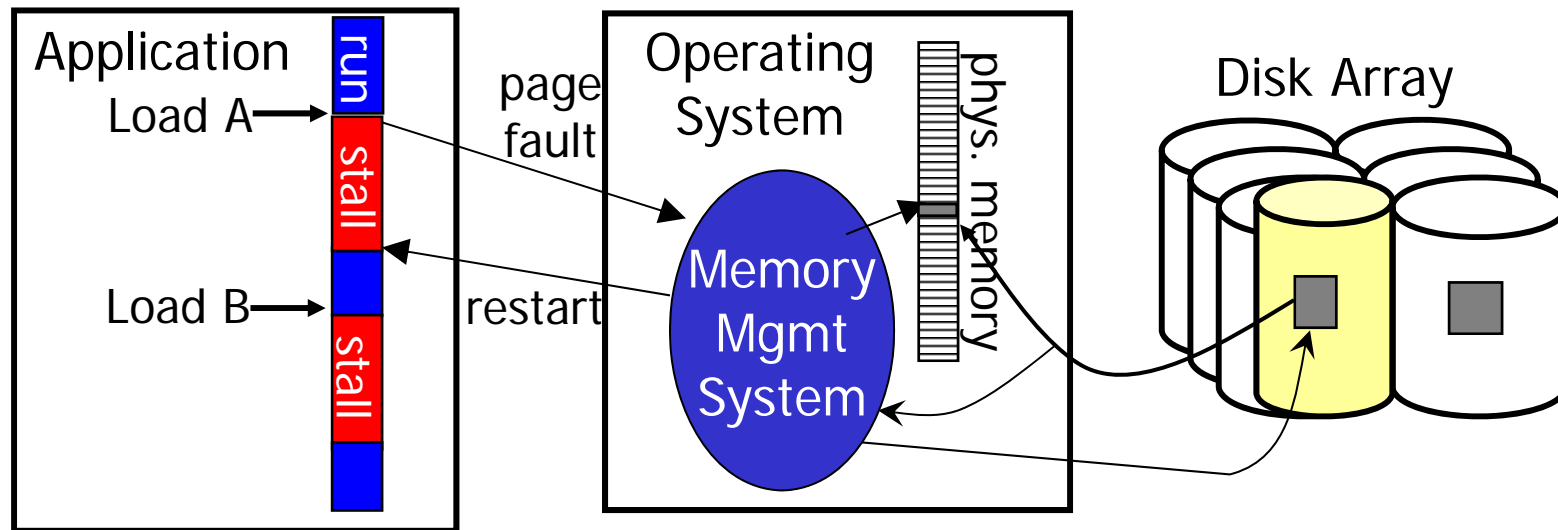


Demand Paging (Process)

- Demand paging is also used 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...



Costs of Demand Paging

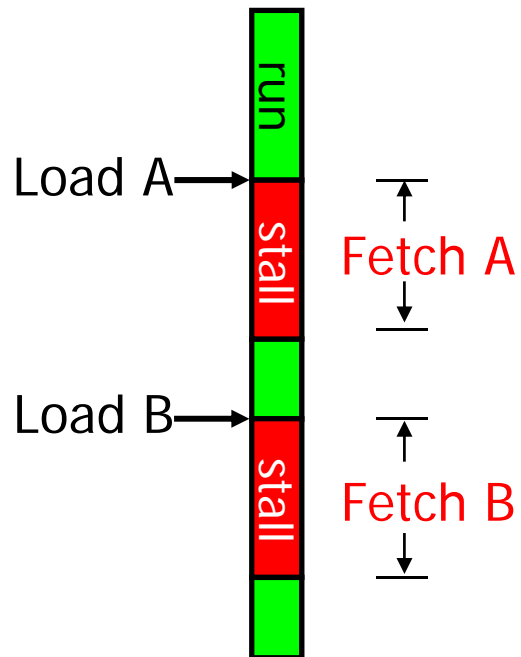


- Timing: Disk read is initiated *when the process needs the page*
- Request size: Process can only page fault on one page at a time, disk sees single page-sized read
- What alternative do we have?

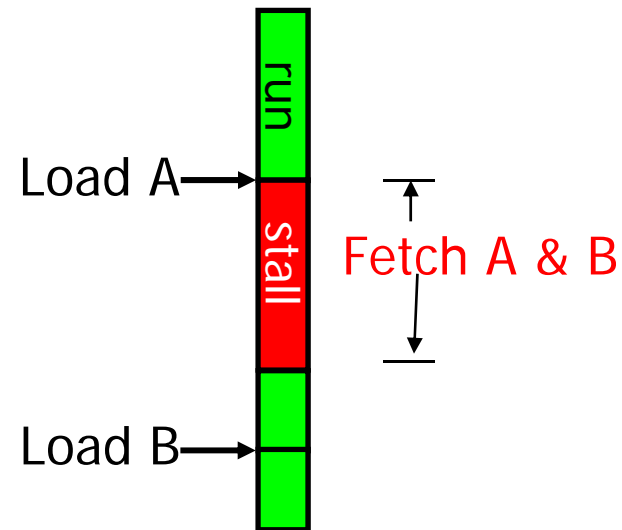


Prepaging (aka Prefetching)

Without Prepaging



With Prepaging



- Predict future page use at time of current fault
 - On what should we base the prediction? What if it's wrong?



Policies



- Fetch Policy – *when* to fetch a page
 - Demand paging vs. Prepaging
- Placement Policy – *where* to put the page
 - Are some physical pages preferable to others?
- Replacement Policy – *what* page to evict to make room?
 - Lots and lots of possible algorithms!



Placement Policy

- In paging systems, memory management hardware can translate any virtual-to-physical mapping equally well
- Why would we prefer some mappings over others?
 - NUMA (non-uniform memory access) multiprocessors
 - Any processor can access entire memory, but local memory is faster
 - FYI: check it out on CDF, run: `numactl --hardware`
 - Cache performance
 - Choose physical pages to minimize cache conflicts
 - These are active research areas!
- Placement policy will have smaller effect than fetch and replacement policies



Policies

- ✓ • Fetch Policy – *when* to fetch a page
 - Demand paging vs. Prepaging
- ✓ • Placement Policy – *where* to put the page
 - Are some physical pages preferable to others?
- Replacement Policy – *what* page to evict to make room?
 - Lots and lots of possible algorithms!



Page Replacement Policy

- When a page fault occurs, the OS loads the faulted page from disk into a page frame of memory
- There may be no free frames available for use
- When this happens, the OS must **replace** a page for each page faulted in
 - It must evict a page (called the **victim**) to free up a frame
 - If the victim has been modified since the last page-out, it must be written back to disk on eviction
 - *Modify* or *dirty* bit in PTE
- The **page replacement algorithm** determines how a victim is chosen



Evicting the Best Page

- Goal of replacement algorithm: to reduce the fault rate by selecting the best victim page to remove
- Replacement algorithms are evaluated on a *reference string* by counting the number of page faults
- The best page to evict is the one never used 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 of time minimizes the number of page faults
 - Proved by Belady, 1966



Belady's Algorithm

- Known as the *optimal* page replacement algorithm because it has the lowest fault rate for any page reference string (aka OPT or MIN)
 - Idea: Replace the page that will not be used for the longest period of time
 - Problem: Have to know the future perfectly
- Why is Belady's useful then? Use it as a yardstick
 - 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
 - If optimal is much better, then algorithm could use some work
 - Random replacement is often the lower bound





Modelling Belady's Algorithm

- Page address list: 2,3,2,1,5,2,4,5,3,2,5,2

Two types:

Cold misses:
first access to
a page
(unavoidable)

**Capacity
misses:**
caused by
replacement
due to
limited size
of memory

2	2	2	2	2	2
	3	3	3	3	3
			1	5	5
4	4	4	2	2	2
3	3	3	3	3	3
5	5	5	5	5	5



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**)





Modelling FIFO

- Page Address List: 0,1,2,3,0,1,4,0,1,2,3,4

	0	1	2	3	0	1	4	0	1	2	3	4	3 frames, 9 faults
Youngest	0	1	2	3	0	1	4	4	4	2	3	3	
		0	1	2	3	0	1	1	1	4	2	2	
Oldest			0	1	2	3	0	0	0	1	4	4	

Anomaly	0	1	2	3	0	1	4	0	1	2	3	4	4 frames, 10 faults
Youngest	0	1	2	3	3	3	4	0	1	2	3	4	
		0	1	2	2	2	3	4	0	1	2	3	
			0	1	1	1	2	3	4	0	1	2	
Oldest				0	0	0	1	2	3	4	0	1	



Conditions for anomaly exist at this point.



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?





Implementing Exact LRU

- Option 1:
 - Time stamp every reference
 - Evict page with oldest time stamp
 - Problems?
 - Need to make PTE large enough to hold meaningful time stamp (may double size of page tables, TLBs)
 - Need to examine every page on eviction to find one with oldest time stamp
- Option 2:
 - Keep pages in a stack. On reference, move the page to the top of the stack. On eviction, replace page at bottom.
 - Problems:
 - Need costly software operation to manipulate stack on EVERY memory reference!



Modelling Exact LRU

- Page Address List: 0,1,2,3,0,1,4,0,1,2,3,4

3 Frames	0	1	2	3	0	1	4	0	1	2	3	4	10 faults
MRU page	0	1	2	3	0	1	4	0	1	2	3	4	
		0	1	2	3	0	1	4	0	1	2	3	
LRU page			0	1	2	3	0	1	4	0	1	2	

4 Frames	0	1	2	3	0	1	4	0	1	2	3	4	8 faults
MRU page	0	1	2	3	0	1	4	0	1	2	3	4	
		0	1	2	3	0	1	4	0	1	2	3	
			0	1	2	3	0	1	4	0	1	2	
LRU page				0	1	2	3	3	3	4	0	1	

- It's a "stack algorithm" => no Belady's Anomaly



Approximating LRU

- Exact LRU is too costly to implement, so approximate it
- LRU approximations use the PTE *reference* bit
- Basic Idea:
 - Initially, all R bits are zero; as processes execute, bits are set to 1 for pages that are used
 - Periodically examine the R bits – we do not know *order* of use, but we know pages that were (or were not) used
- Additional-Reference-Bits Algorithm
 - Keep a counter for each page
 - At regular intervals, for every page do:
 - Shift R bit into high bit of counter register
 - Shift other bits to the right
 - Pages with “larger” counters were used more recently





Second Chance Algorithm

- FIFO, but inspect reference bit when page is selected
 - If ref bit is 0, replace the page
 - If ref bit is 1, clear ref bit, reset arrival time of page to current time
 - Pages that are used often enough to keep reference bits set will not be replaced
- Can combine with Modify bit to create 4 classes of pages
 - Called “Not Recently Used” in text



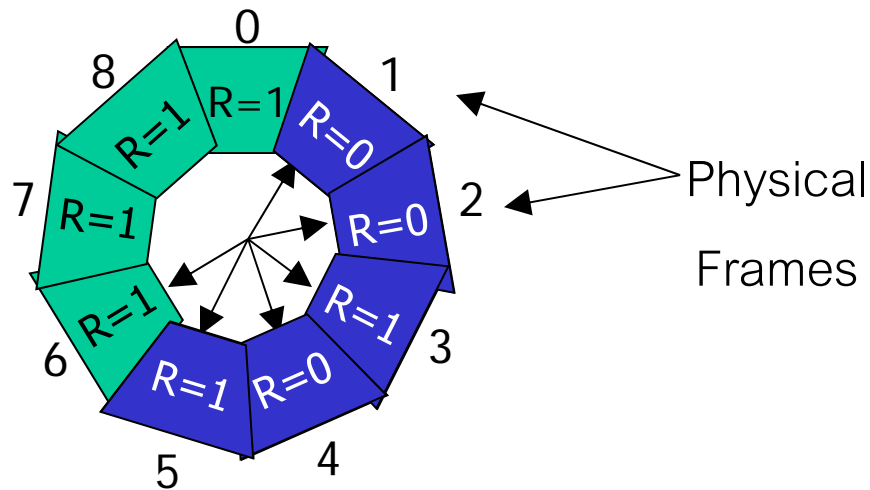
Implementing Second Chance (clock)

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 (aka use 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



Modelling Clock



- 1st page fault:
 - Advance hand to frame 4, use frame 3
- 2nd page fault (assume none of these pages are referenced)
 - Advance hand to frame 6, use frame 5



Counting-based Replacement

- Count number of uses of a page
- Least-Frequently-Used (LFU)
 - Replace the page used least often
 - Pages that are heavily used at one time tend to stick around even when not needed anymore
 - Newly allocated pages haven't had a chance to be used much
- Most-Frequently-Used (MFU)
 - Favours new pages
- Neither is common, both are poor approximations of OPT



Page Buffering

- Preceding discussion assumed the replacement algorithm is run and a victim page selected *when a new page needs to be brought in*
- Most of these algorithms are too costly to run on every page fault
 - Maintain a pool of free pages
 - Run replacement algorithm when pool becomes too small (“low water mark”), free enough pages to at once replenish pool (“high water mark”)
 - On page fault, grab a frame from the free list
 - Frames on free list still hold previous contents, can be “rescued” if virtual page is referenced before reallocation



Assignment 3

- Implementation of virtual memory system
- A few quick mentions:
 - Make sure to work on a Linux machine (CDF preferably!)
 - Don't use the Virtual Machine from A1! If you must, keep in mind that your code should work on CDF in the end.
 - Read the handout carefully and look into the starter code
 - Make sure to test your code before submission!
 - Code that does not compile will get a zero, no excuses!
 - Read the **submission instructions carefully!**
 - Best assignment(s) might receive a **symbolic prize!**



Assignment 3

- Quick tips:
 - Do the online exercise!
 - In the exercise, we generate some memory traces of accesses done by various programs
 - Next, go through the trace of VAs, translate each of them using a 2-level page table and demand paging.
 - Second part: Implement page replacement policies
 - Analysis on memory access patterns
 - Once again, we WILL detect plagiarism! Don't do it!
- More translation exercises to follow
 - Should know by now bitwise operations, conversions from hex to binary, etc.



Announcements

- Extra virtual mem translation exercises will be posted ...
 - Practice during reading week..
 - Online only, not “doable” during lecture
 - No need to submit on MarkUs, submit directly on course webpage (the last answer will be marked)
 - Submit individually, may consult with classmates though!
 - The point is to learn how to do translations and apply replacement policies
 - Important exercises to clarify your understanding of translations and page replacement
 - Make sure to revise bitwise operations, translations binary-to-hex and vice-versa, etc..



Announcements

- Midterm
 - Next week, during Thursday lecture timeslot
 - Evening section – first hour of lecture, followed by lecture as usual..
 - Covers up to Virtual Memory (excluding)
 - Write in pen, cannot regrade tests written in pencil!
- Location and logistics:
 - Check website!
 - Check Piazza announcements!
- Next tutorial will be transformed into a lecture
 - Come to the regular lecture room, not your tutorial room!