

Paging & Page Replacement

Sirak Kaewjamnong

517-312 Operating Systems

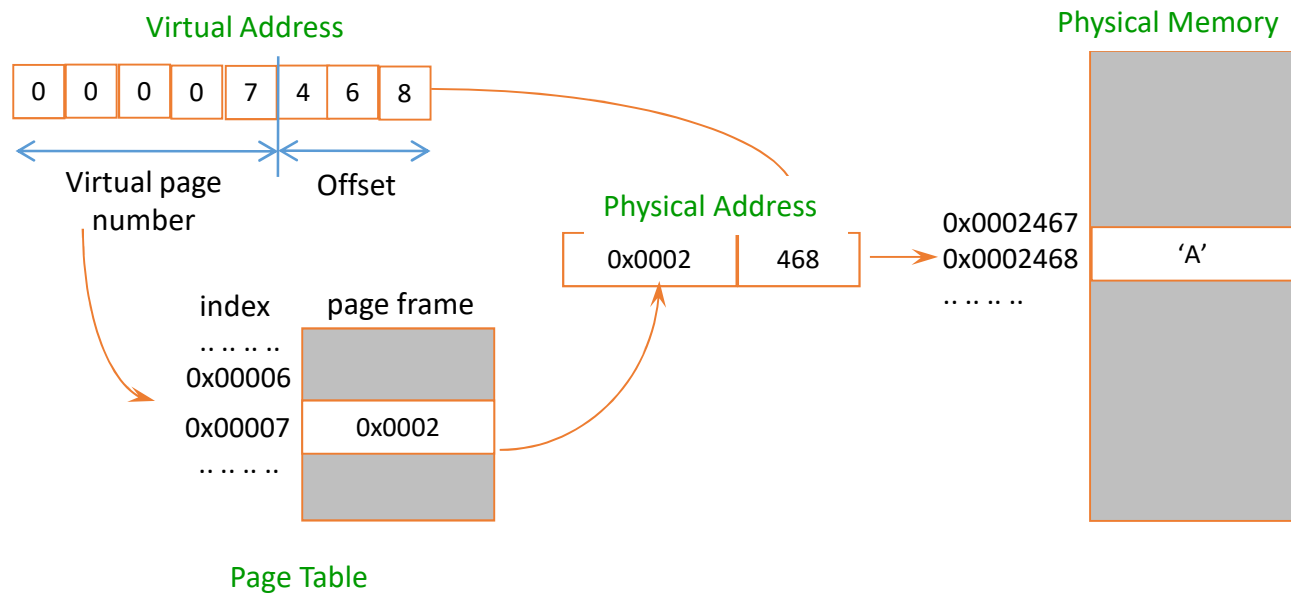
Speeding Up Paging

Major issues faced:

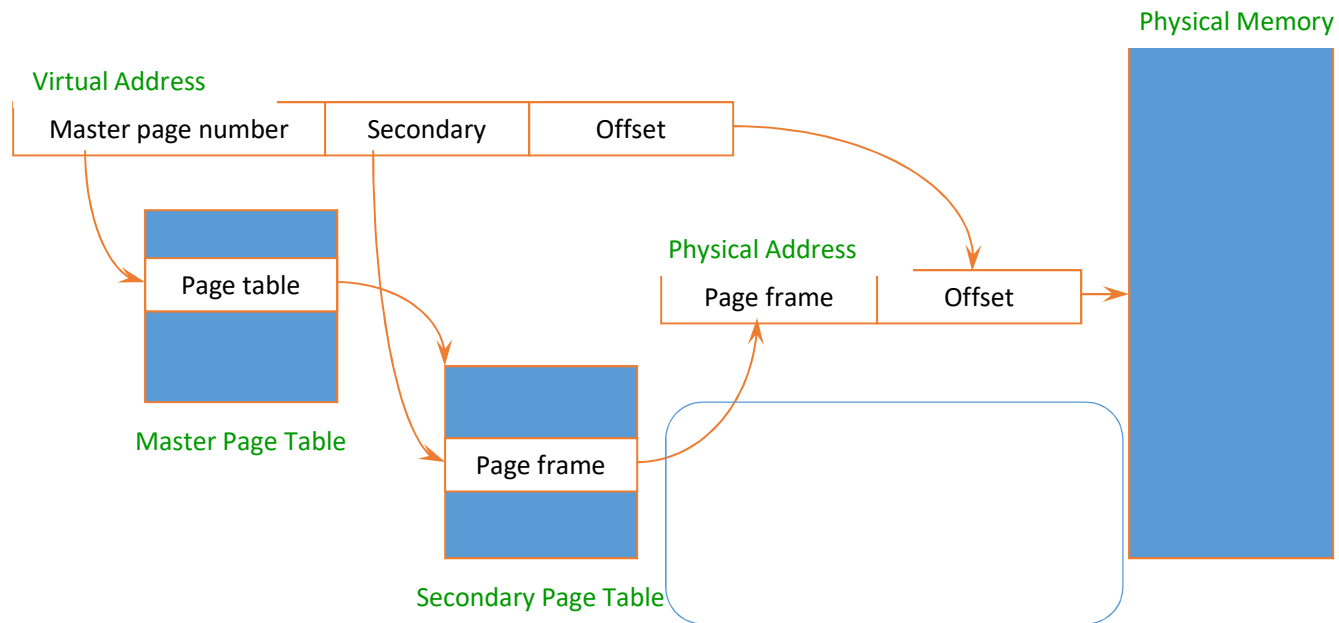
1. The mapping from virtual address to physical address must be fast.
2. If the virtual address space is large, the page table will be large.

Review: Page Lookups

Example: how do we 'load 0x00007468'?



Review: Two-Level Page Tables



Problems with Efficiency

- If OS is involved in every memory access for PT look up --- too slow (user<->kernel mode switch)
- Every look-up requires multiple memory access
 - Our original page table scheme already doubled the cost of doing memory lookups
 - One lookup into the page table, another to fetch the data
 - Now two-level page tables triple the cost!
 - Two lookups into the page tables, a third to fetch the data
 - And this assumes the page table is in memory

Solution to Efficient Translation

- How can we use paging but also have lookups cost about the same as fetching from memory?
- *Hardware Support*
 - Cache translations in hardware
 - Translation Lookaside Buffer (TLB)
 - TLB managed by Memory Management Unit (MMU)

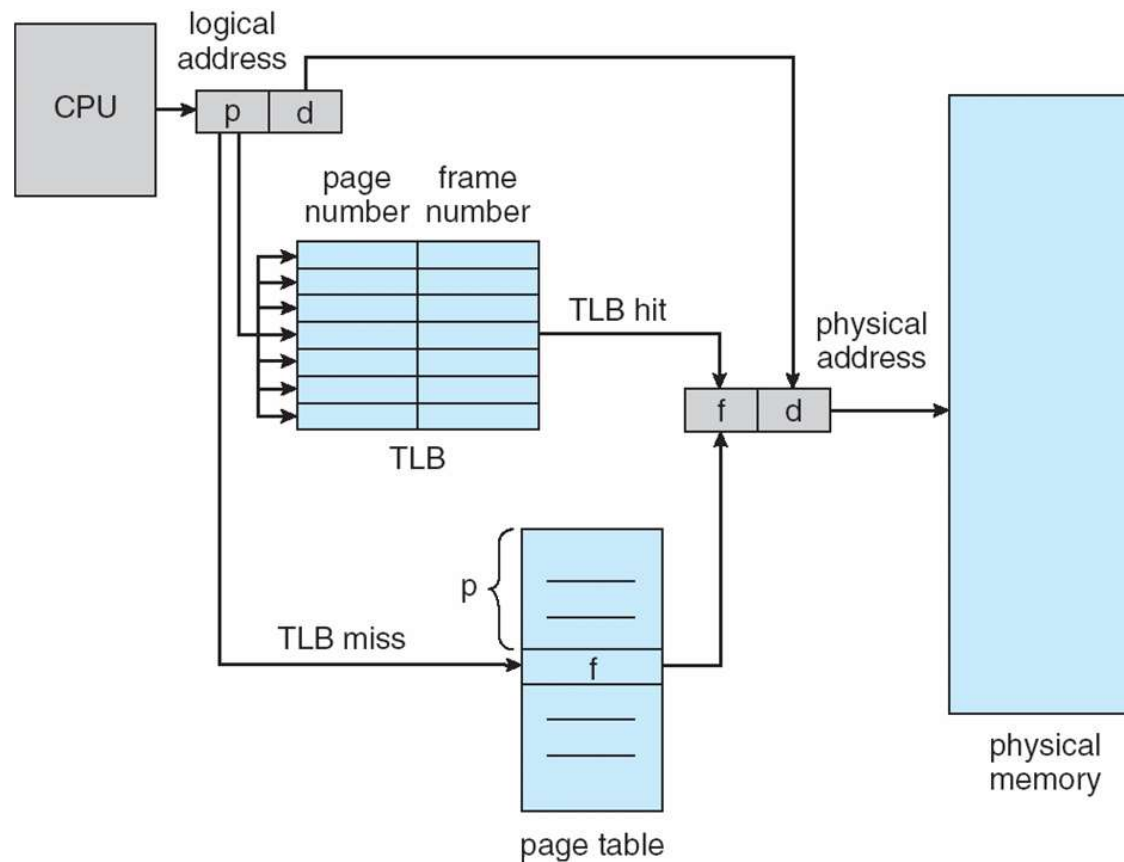
Translation Look-Aside Buffer (TLB)

- Translation Lookaside Buffers
 - *Can be done in a single machine cycle*
- TLB 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
- TLB exploit locality
 - Processes only use a handful of pages at a time
 - 16-48 entries/pages (64-192K)
 - Only need those pages to be “mapped”
 - Hit rates are therefore very important

Translation Look-Aside Buffer

- Some TLBs store **address-space identifiers (ASIDs)** in each TLB entry – uniquely identifies each process to provide address-space protection for that process
 - Otherwise need to flush at every context switch
- TLBs typically small (64 to 1,024 entries)
- On a TLB miss, value is loaded into the TLB for faster access next time
 - Replacement policies must be considered

Paging Hardware With TLB



Effective Access Time

- Hit ratio – percentage of times that a page number is found in the TLB
- An 80% hit ratio means that we find the desired page number in the TLB 80% of the time.
- Suppose that 10 nanoseconds to access memory.
 - If we find the desired page in TLB then a mapped-memory access take 10 ns
 - Otherwise we need two memory access so it is 20 ns
- **Effective Access Time (EAT)**
$$\text{EAT} = 0.80 \times 10 + 0.20 \times 20 = 12 \text{ nanoseconds}$$

implying 20% slowdown in access time
- Consider a more realistic hit ratio of 99%,
$$\text{EAT} = 0.99 \times 10 + 0.01 \times 20 = 10.1\text{ns}$$

implying only 1% slowdown in access time.

Back to Virtual memory

- Code needs to be in memory to execute, but entire program rarely used
 - Error code, unusual routines, large data structures
- Entire program code not needed at same time
- Consider ability to execute partially-loaded program
 - Program no longer constrained by limits of physical memory
 - Each program takes less memory while running -> more programs run at the same time
 - Increased CPU utilization and throughput with no increase in response time or turnaround time
 - Less I/O needed to load or swap programs into memory -> each user program runs faster

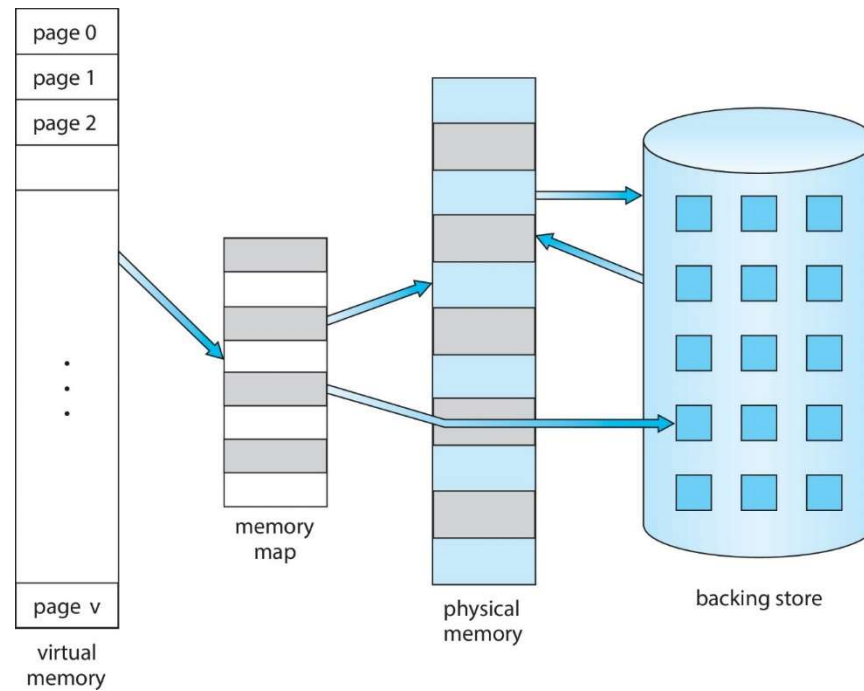
Virtual memory

- **Virtual memory** – separation of user logical memory from physical memory
 - Only part of the program needs to be in memory for execution
 - Logical address space can therefore be much larger than physical address space
 - Allows address spaces to be shared by several processes
 - Allows for more efficient process creation
 - More programs running concurrently
 - Less I/O needed to load or swap processes

Virtual memory (Cont.)

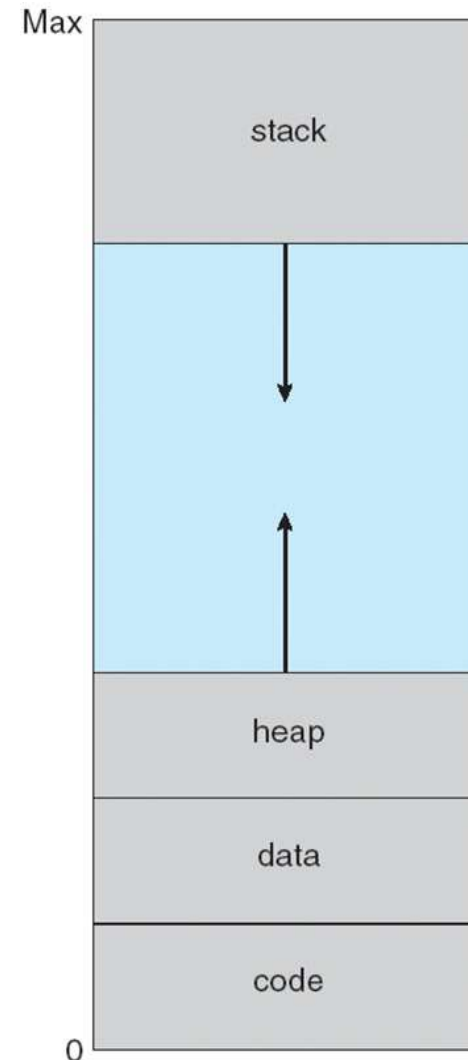
- **Virtual address space** – logical view of how process is stored in memory
 - Usually start at address 0, contiguous addresses until end of space
 - Meanwhile, physical memory organized in page frames
 - MMU must map logical to physical
- Virtual memory can be implemented via:
 - Demand paging
 - Demand segmentation

Virtual Memory That is Larger Than Physical Memory

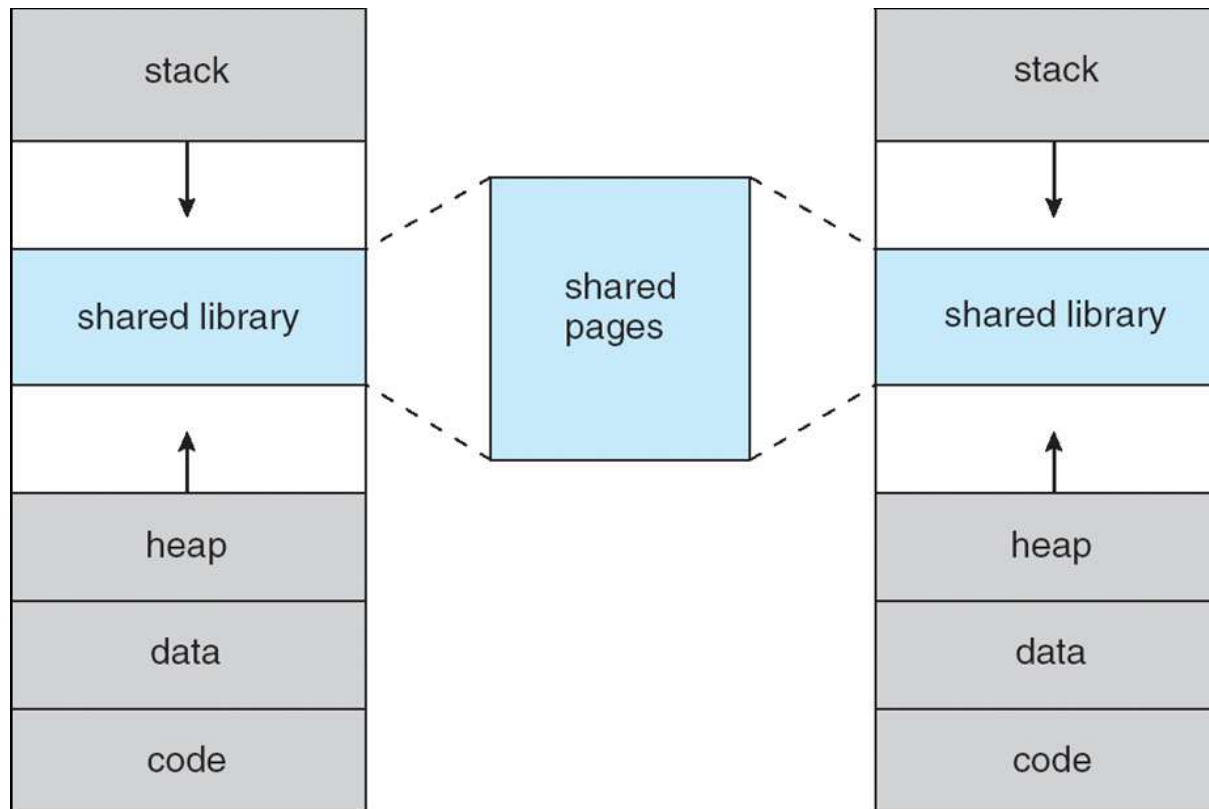


Virtual-address Space

- Usually design logical address space for stack to start at Max logical address and grow “down” while heap grows “up”
 - Maximizes address space use
 - Unused address space between the two is hole
 - No physical memory needed until heap or stack grows to a given new page
- Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc
- System libraries shared via mapping into virtual address space
- Shared memory by mapping pages read-write into virtual address space
- Pages can be shared during `fork()`, speeding process creation



Shared Library Using Virtual Memory

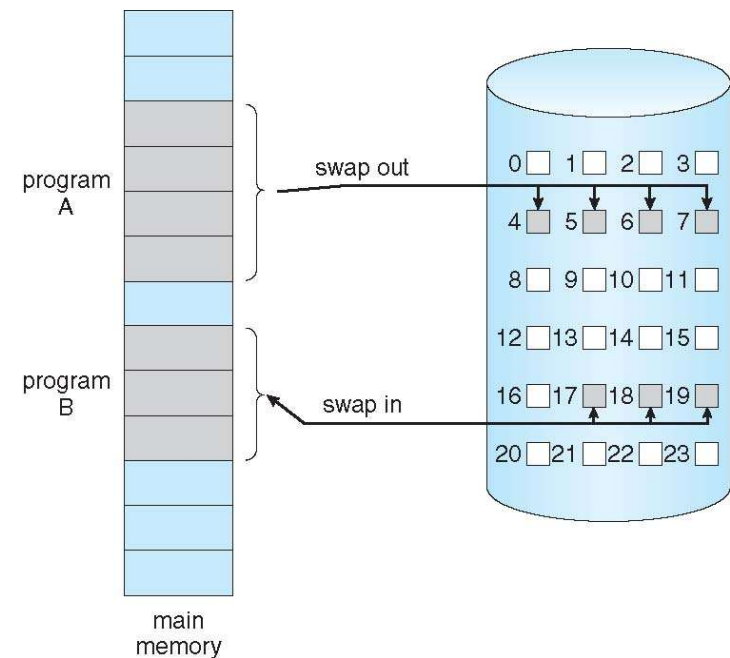


Demand paging

- 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/backing store)
 - The movement of pages between memory and disk is done by the OS, and is transparent to the application

Demand paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed
 - Less I/O needed, no unnecessary I/O
 - Less memory needed
 - Faster response
 - More users
- Similar to paging system with swapping (diagram on right)
- Page is needed \Rightarrow reference to it
 - invalid reference \Rightarrow abort
 - not-in-memory \Rightarrow bring to memory
- **Lazy swapper** – never swaps a page into memory unless page will be needed
 - Swapper that deals with pages is a **pager**



Valid-Invalid Bit

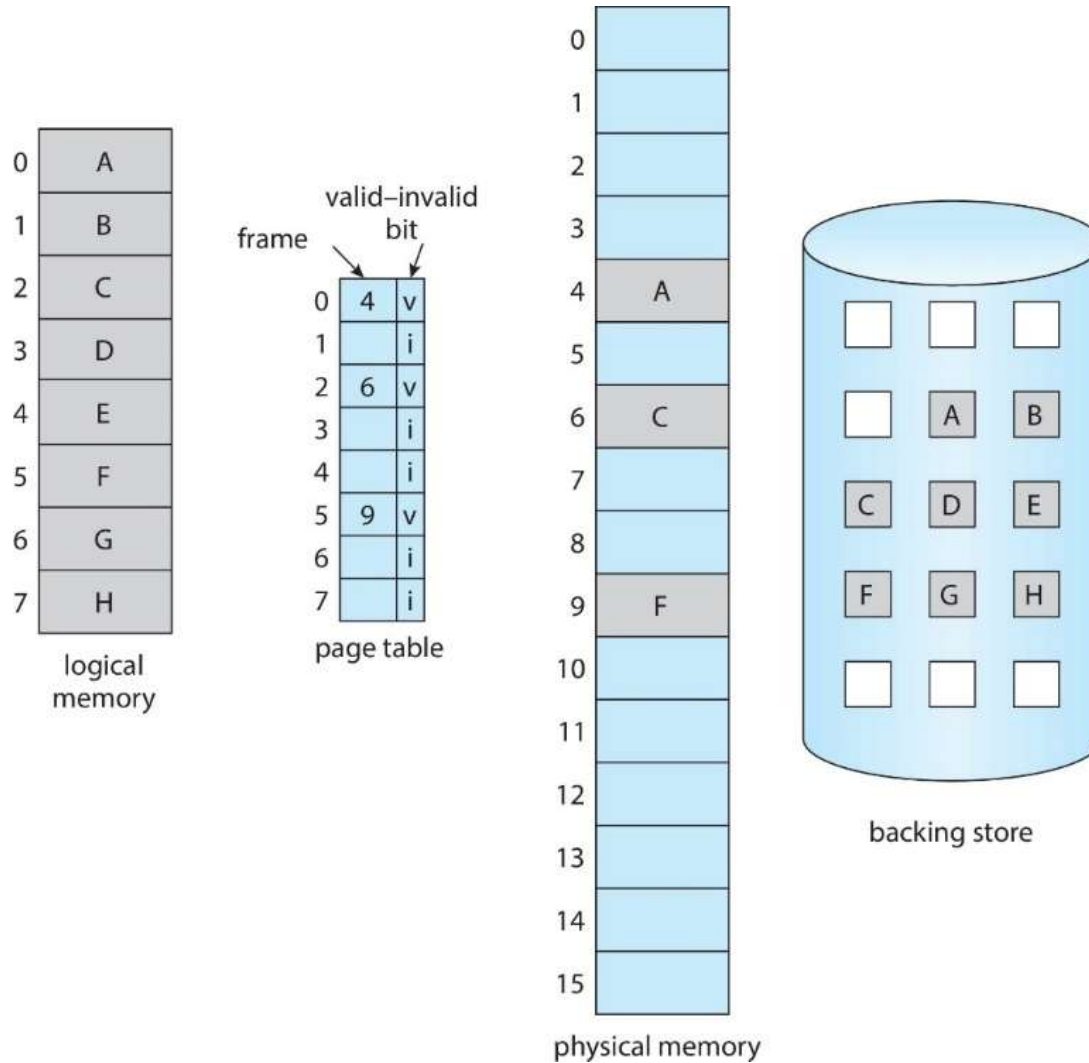
- With each page table entry a valid–invalid bit is associated (**v** \Rightarrow in-memory – **memory resident**, **i** \Rightarrow not-in-memory)
- Initially valid–invalid bit is set to **i** on all entries
- Example of a page table snapshot:

Frame #	valid-invalid bit
	v
	v
	v
	i
...	
	i
	i

page table

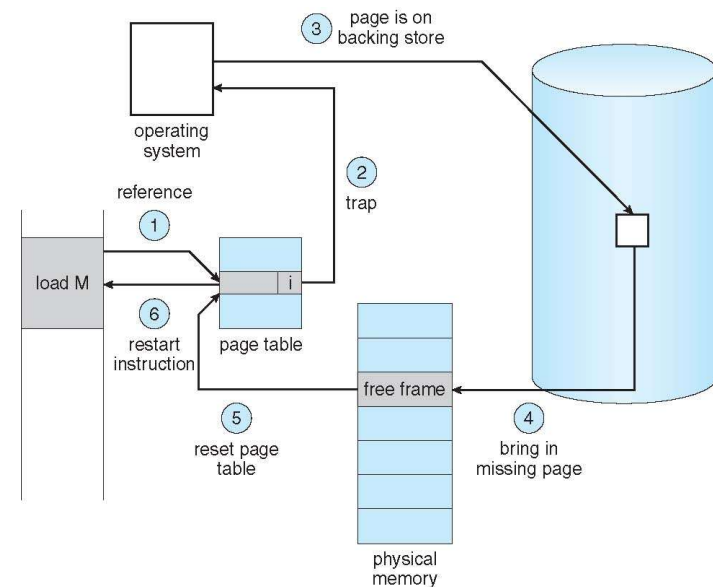
- During MMU address translation, if valid–invalid bit in page table entry is **i** \Rightarrow page fault

Page Table When Some Pages Are Not in Main Memory



Steps in Handling Page Fault

1. If there is a reference to a page, first reference to that page will trap to operating system
 - Page fault
2. Operating system looks at another table to decide:
 - Invalid reference \Rightarrow abort
 - Just not in memory
3. Find free frame
4. Swap page into frame via scheduled disk operation
5. Reset tables to indicate page now in memory
Set validation bit = **v**
6. Restart the instruction that caused the page fault

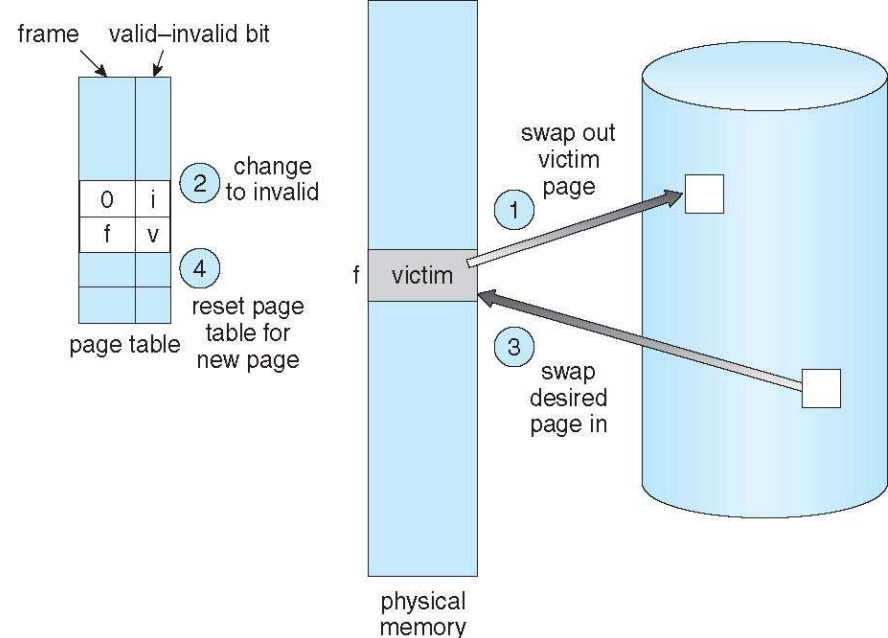


Page Replacement

- Prevent **over-allocation** of memory by modifying page-fault service routine to include page replacement
- Page replacement completes separation between logical memory and physical memory – large virtual memory can be provided on a smaller physical memory
- Use **modify (dirty) bit** to reduce overhead of page transfers – only modified pages are written to disk
- 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

Basic Page Replacement

1. Find the location of the desired page on disk
2. Find a free frame:
 - If there is a free frame, use it
 - If there is no free frame, use a page replacement algorithm to select a **victim frame**
 - Write victim frame to disk if dirty
3. Bring the desired page into the (newly) free frame; update the page and frame tables
4. Continue the process by restarting the instruction that caused the trap



Page Replacement Strategies

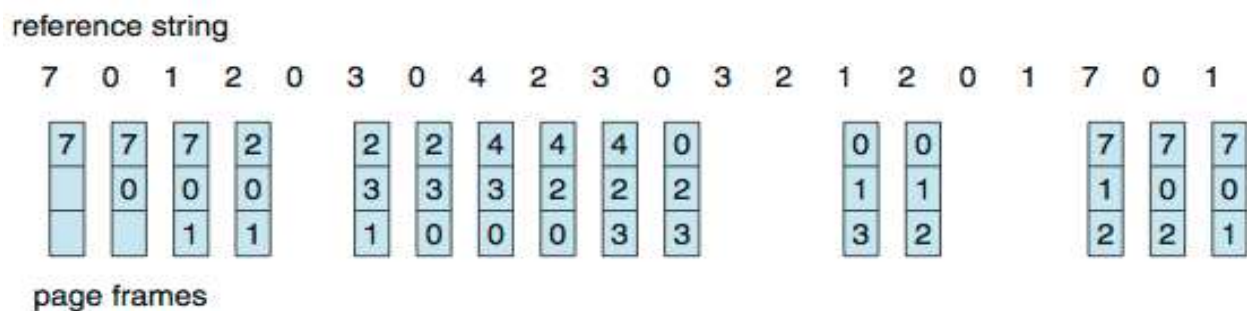
- The Principle of Optimality
 - Replace the page that will not be used again the farthest time in the future
- Random replacement
 - Choose a page randomly
- FIFO – First In First Out
 - Replace the page that has been in memory the longest
- LRU – Least Recently Used
 - Replace the page that has not been used for the longest time
- NRU – Not Recently Used
 - An approximation to LRU

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**)

First-In-First-Out (FIFO) Algorithm

- Reference string: **7,0,1,2,0,3,0,4,2,3,0,3,2,1,2,0,1,7,0,1**
- 3 frames (3 pages can be in memory at a time per process)



15 page faults

- Can vary by reference string: consider 1,2,3,4,1,2,5,1,2,3,4,5
 - Adding more frames can cause more page faults!
 - Belady's Anomaly
- How to track ages of pages?
 - Just use a FIFO queue

Belady's Anomaly in FIFO

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

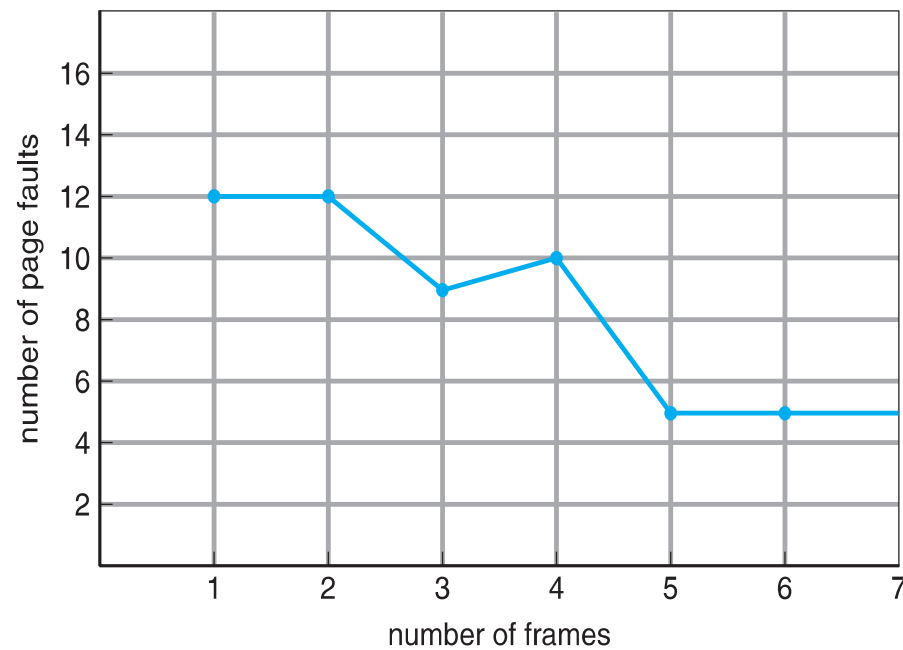
1	1	1	2	3	4	1	1	1	2	5	5
	2	2	3	4	1	2	2	2	5	3	3
		3	4	1	2	5	5	5	3	4	4
PF	PF	PF	PF	PF	PF	PF	X	X	PF	PF	X

- 3 frames and 9 page faults

1	1	1	1	1	1	2	3	4	5	1	2
	2	2	2	2	2	3	4	5	1	2	3
		3	3	3	3	4	5	1	2	3	4
			4	4	4	5	1	2	3	4	5
PF	PF	PF	PF	X	X	PF	PF	PF	PF	PF	PF

- 4 frames and 10 page faults

FIFO Illustrating Belady's Anomaly



Least Recently Used (LRU) Algorithm

- Use past knowledge rather than future
- Replace page that has not been used in the most amount of time
- Associate time of last use with each page

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

7	7	7	2																
	0	0	0																
		1	1																

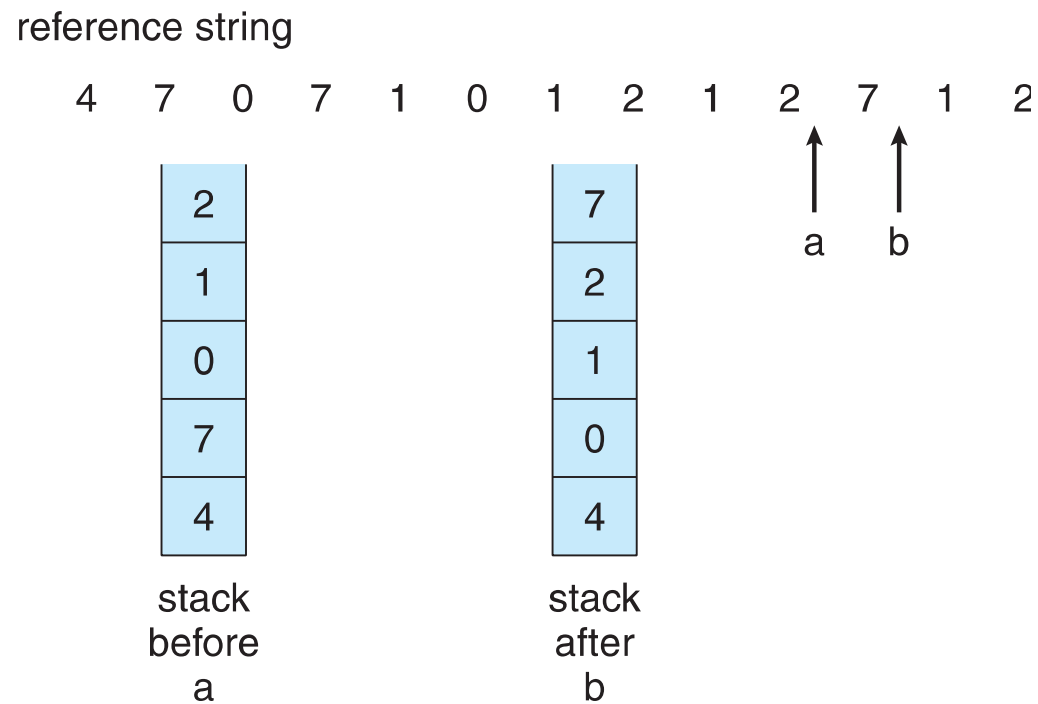
page frames

- 12 faults – better than FIFO Generally good algorithm and frequently used
- But how to implement?

LRU Algorithm (Cont.)

- Counter implementation
 - Every page entry has a counter; every time page is referenced through this entry, copy the clock into the counter
 - When a page needs to be changed, look at the counters to find smallest value
 - Search through table needed
- Stack implementation
 - Keep a stack of page numbers in a double link form:
 - Page referenced:
 - move it to the top
 - requires 6 pointers to be changed
 - But each update more expensive
 - No search for replacement
- LRU is **stack algorithms** that don't have Belady's Anomaly

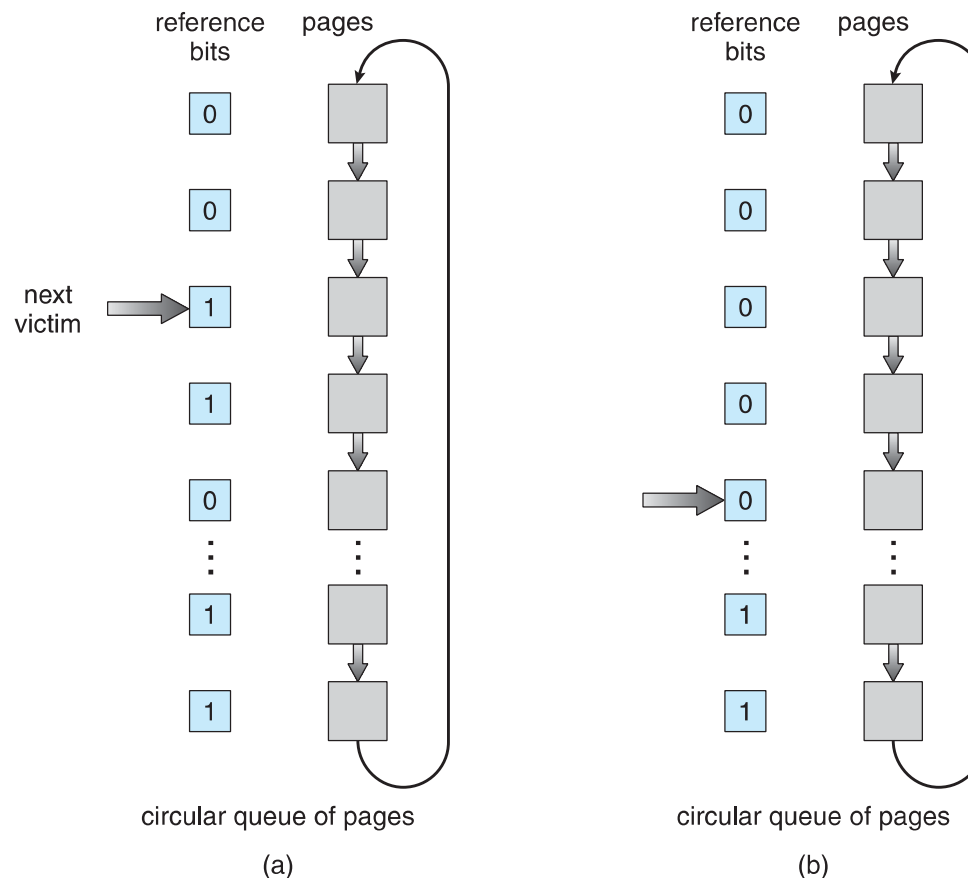
Use Of A Stack to Record Most Recent Page References



LRU Approximation Algorithms

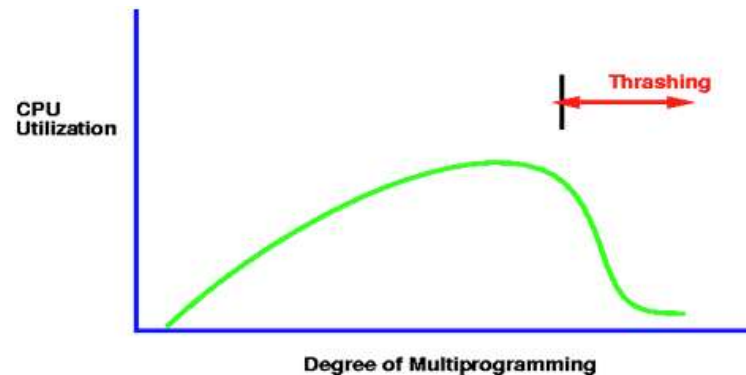
- LRU needs special hardware and still slow
- NRU: Evict a page that is **NOT** recently used;
- LRU: evict a page that is **LEAST** recently used
- **Reference bit**
 - With each page associate a bit, initially = 0
 - When page is referenced bit set to 1
 - Replace any with reference bit = 0 (if one exists)
 - We do not know the order, however
- **Second-chance algorithm**
 - Generally FIFO, plus hardware-provided reference bit
 - **Clock** replacement
 - If page to be replaced has
 - Reference bit = 0 -> replace it
 - reference bit = 1 then:
 - set reference bit 0, leave page in memory
 - replace next page, subject to same rules

Second-Chance (clock) Page-Replacement Algorithm



Thrashing and CPU utilization

- As the page fault rate goes up, processes get suspended on page queues for the disk
- The system may try to optimize performance by starting new jobs
 - But is it always good?
- Starting new jobs will reduce the number of page frames available to each process, increasing the page fault requests
- System throughput plunges



Sources and References

- <https://www.geeksforgeeks.org/operating-system-beladys-anomaly/>
- Andrew S. Tanenbaum, Herbert Bos, Modern Operating 4th edition, Pearson, 2015
- Presentations by Ding Yuan, ECE Dept., University of Toronto
- Operating System Concepts 10th book official slides by Abraham Silberschatz, Greg Gagne, Peter B. Galvin