

Virtual Memory (1)

Dr. Jun Zheng
CSE325 Principles of Operating
Systems
11/4/2019



Background

- ❑ 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

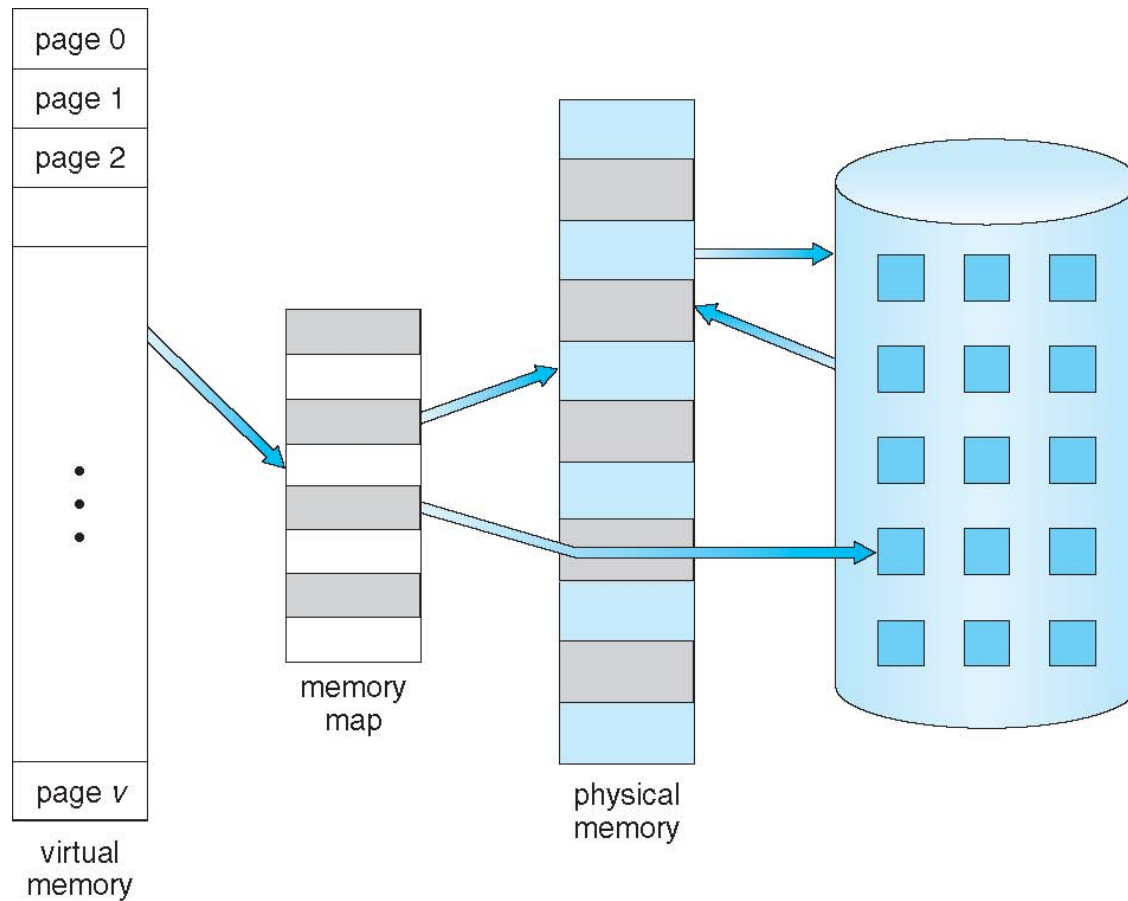
Background (Cont.)

- ❑ **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

Background (Cont.)

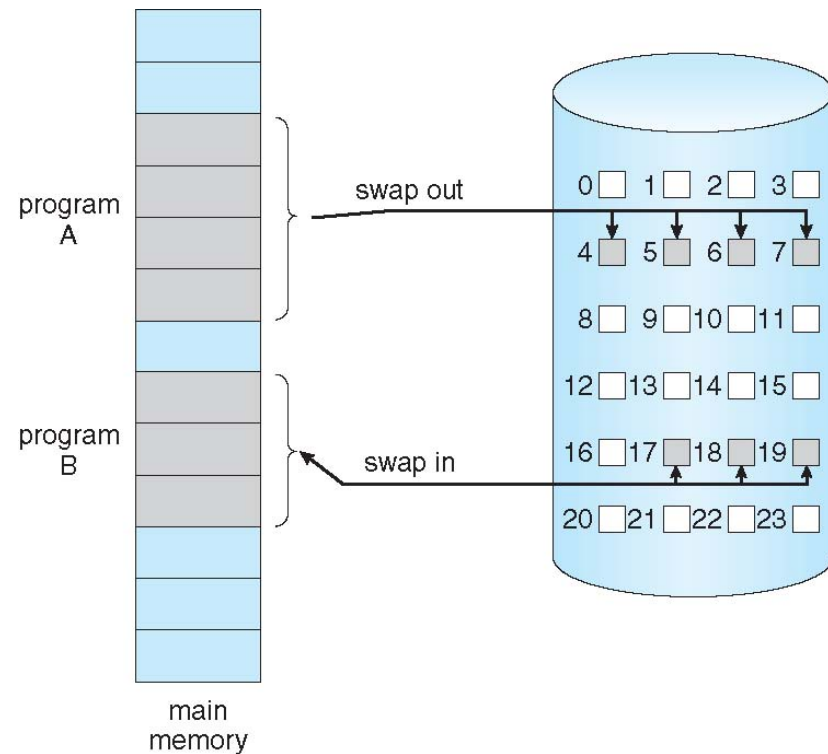
- ❑ **Virtual address space** – logical view of how process is stored in memory
 - ❑ Usually start at address 0, contiguous addresses until end of space
 - ❑ Meanwhile, physical memory organized in page frames
 - ❑ MMU must map logical to physical
- ❑ Virtual memory can be implemented via demand paging: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

Virtual Memory That is Larger Than Physical Memory



Demand Paging

- ❑ 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
- ❑ **Lazy swapper** – never swaps a page into memory unless page will be needed
 - ❑ Swapper that deals with pages is a **pager**



Basic Concepts

- ❑ With swapping, pager guesses which pages will be used before the process is swapped out again
 - ❑ Instead of swapping in a whole process, pager brings in only those pages into memory
- ❑ How to determine that set of pages?
 - ❑ Need new MMU functionality to implement demand paging
- ❑ If pages needed are already **memory resident**
 - ❑ No difference from non demand-paging
- ❑ If page needed and not memory resident
 - ❑ Need to detect and load the page into memory from storage
 - ❑ Without changing program behavior
 - ❑ Without programmer needing to change code

Valid-Invalid Bit

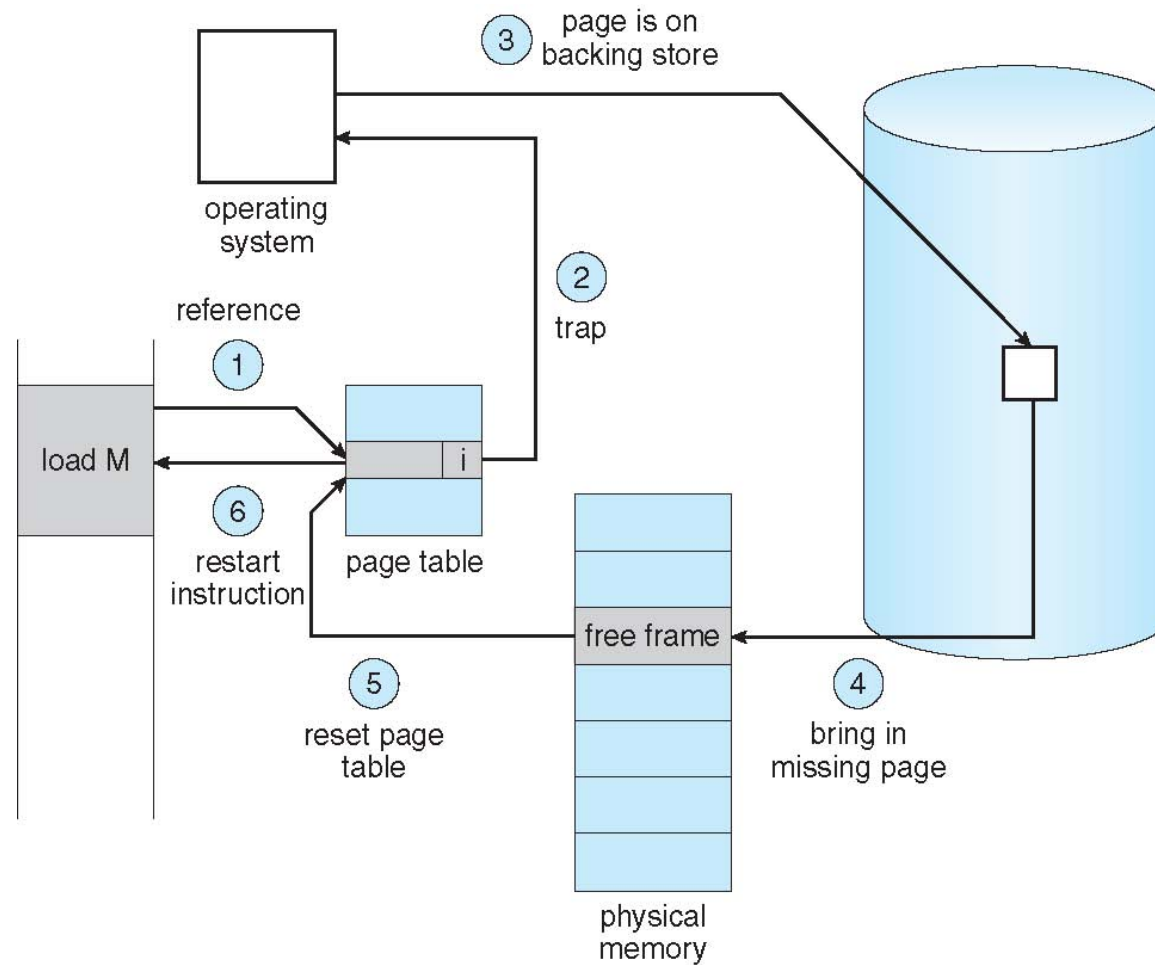
- ❑ With each page table entry a valid–invalid bit is associated (**v** \Rightarrow in-memory – **memory resident**, **i** \Rightarrow not-in-memory)
- ❑ Initially valid–invalid bit is set to **i** on all entries
- ❑ Example of a page table snapshot:

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**

Steps in Handling a Page Fault



Aspects of Demand Paging

- ❑ Extreme case – start process with *no* pages in memory
 - ❑ OS sets instruction pointer to first instruction of process, non-memory-resident -> page fault
 - ❑ And for every other process pages on first access
 - ❑ **Pure demand paging**
- ❑ Actually, a given instruction could access multiple pages -> multiple page faults
 - ❑ Consider fetch and decode of instruction which adds 2 numbers from memory and stores result back to memory
 - ❑ Pain decreased because of **locality of reference**
- ❑ Hardware support needed for demand paging
 - ❑ Page table with valid / invalid bit
 - ❑ Secondary memory (swap device with **swap space**)
 - ❑ Instruction restart