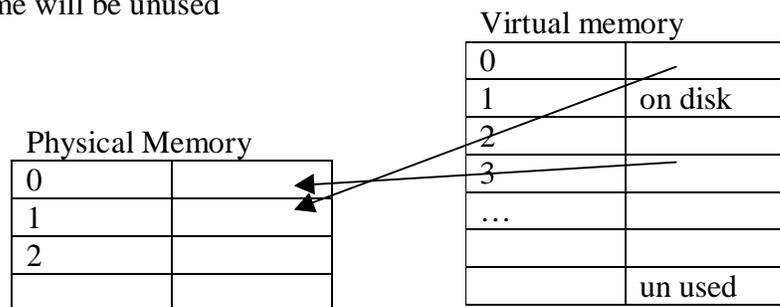


Lecture 29

July 18

Virtual Memory

- Modern PCs can address 2^{32} bytes = 4GB of memory
 - o If we allow some “memory” to be stored on disk, there will be more space to use
- The larger “virtual memory could have parts in memory & other parts on disk”
 - o All of this should be invisible to the program
 - o When a program does a memory read, it doesn’t know if the data is even in memory
 - o If its not in memory, it must be moved from disk into memory
 - Not necessarily in the same location every time
 - o Every byte in the address space is given a “virtual address”
 - $0-2^{32}-1$ in a 32 bit arch.
 - Programs will use the virtual address
 - This will be translated to a “physical address”
 - If its in memory
- The virtual address space is broken into “pages”
 - o E.g. with 4 KB pages
 - $4\text{ GB} - 2^{20} \sim 1\text{ million pages}$
 - o some of the pages will be in memory at a given time
 - o some will be on the disk
 - o some will be unused



- the space on disk that holds the VM pages is called the “swap space”
 - o this is managed by the OS
 - o the CPU doesn't do disk access directly
 - o the amount & location of the swap space is determined by the OS
 - o the transfer memory <-> disk & keeping track of pages on disk is done by the OS
 - o every time memory is accessed there must be a virtual → physical address translation
 - can't be done by the OS (too slow)
 - need hardware support for this
 - the “memory Management Unit” (MMU)
 - o the MMU's job is to manage CPU <-> communication
 - virtual → physical translation
 - if the virtual address is in RAM, read & give the data back
 - if not, send a “page fault” to the OS
 - the OS is activated & must swap the required page into RAM
 - might have to move something else out to make room
 - once it's in memory, the MMU can retrieve the data for the CPU