

Memory Management

Huginn oc Muninn fljúga hverjan dag Jörmungrund yfir; óumk ek of Hugin at hann aftr né komit, þó sjámk meir um Munin.

O'er Mithgarth Hugin and Munin both Each day set forth to fly; For Hugin I fear lest he come not home, But for Munin my care is more.

— *Grímnismál* (Poetic Edda)

Memory management

- ✦ Basic memory management
- ✦ Swapping
- ✦ Virtual memory
- ✦ Page replacement algorithms
- ✦ Modeling page replacement algorithms
- ✦ Design issues for paging systems
- ✦ Implementation issues
- ✦ Segmentation

In an ideal world...

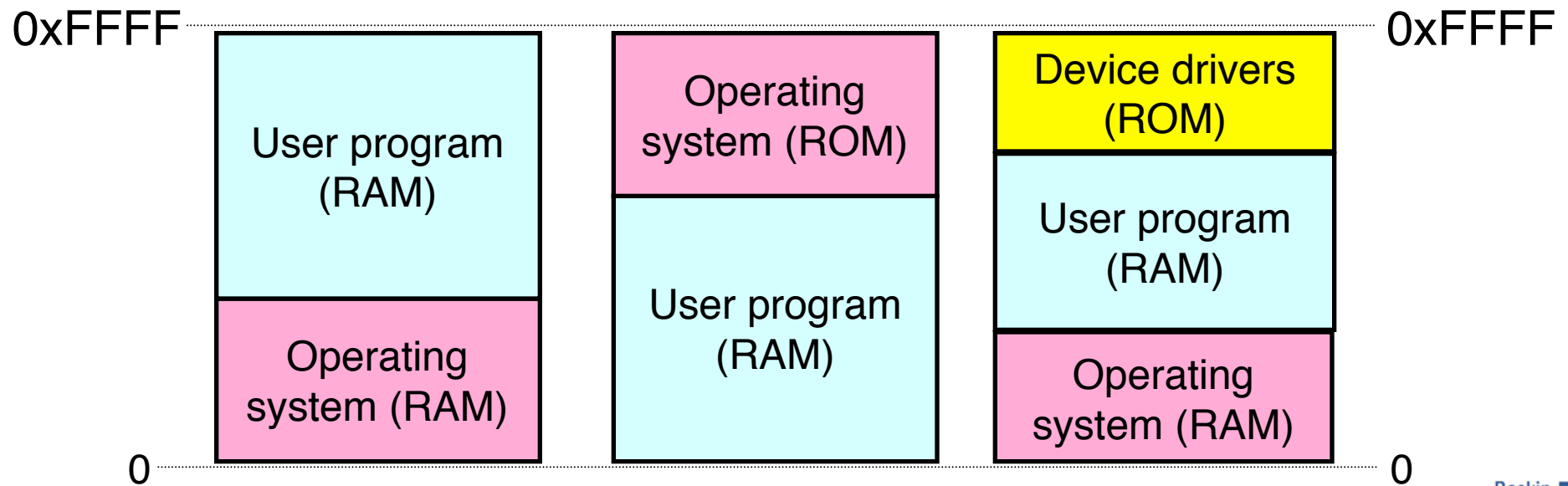
- ✦ The ideal world has memory that is
 - Very large
 - Very fast
 - Non-volatile (doesn't go away when power is turned off)
- ✦ The real world has memory that is:
 - Very large
 - Very fast
 - Affordable!
 - ➡ Pick any two...
- ✦ Memory management goal: make the real world look as much like the ideal world as possible

Memory hierarchy

- ✦ What is the memory hierarchy?
 - Different levels of memory
 - Some are small & fast
 - Others are large & slow
- ✦ What levels are usually included?
 - Cache: small amount of fast, expensive memory
 - L1 (level 1) cache: usually on the CPU chip
 - L2: may be on or off chip
 - L3 cache: off-chip, made of SRAM
 - Main memory: medium-speed, medium price memory (DRAM)
 - Disk: many gigabytes of slow, cheap, non-volatile storage
- ✦ Memory manager handles the memory hierarchy

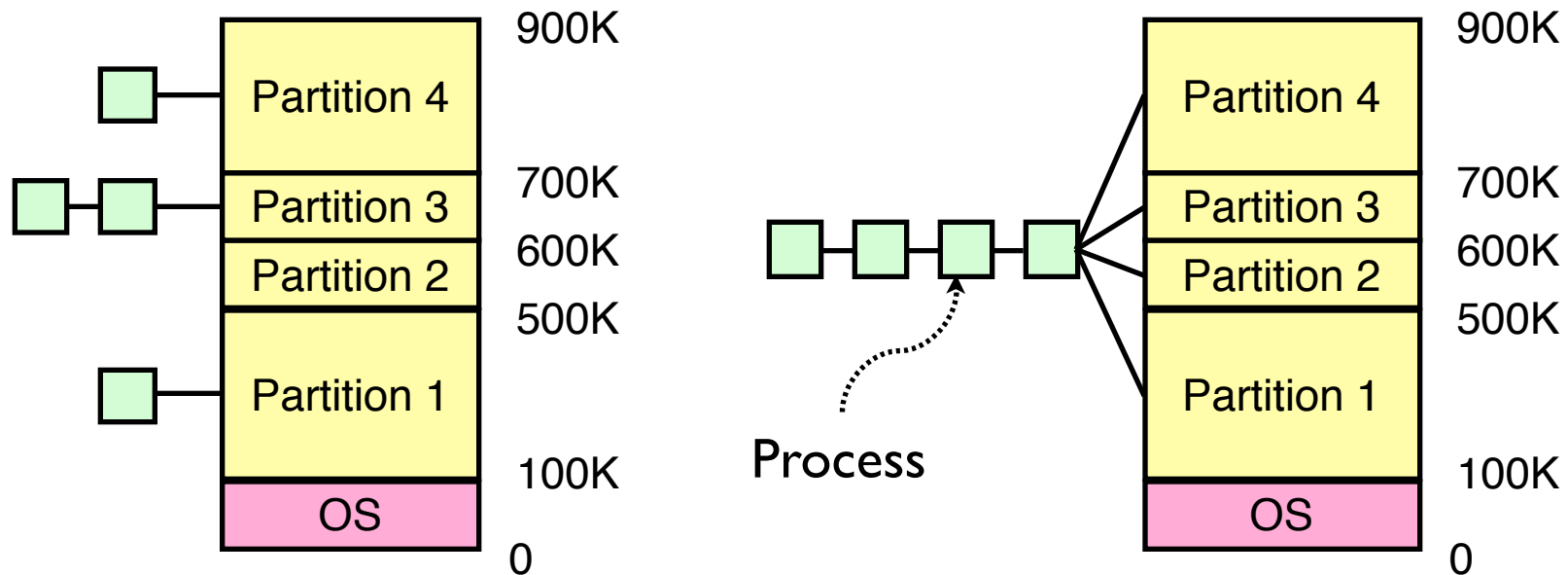
Basic memory management

- ◆ Components include
 - Operating system (perhaps with device drivers)
 - Single process
- ◆ Goal: lay these out in memory
 - Memory protection may not be an issue (only one program)
 - Flexibility may still be useful (allow OS changes, etc.)
- ◆ No swapping or paging



Fixed partitions: multiple programs

- ✦ Fixed memory partitions
 - Divide memory into fixed spaces
 - Assign a process to a space when it's free
- ✦ Mechanisms
 - Separate input queues for each partition
 - Single input queue: better ability to optimize CPU usage

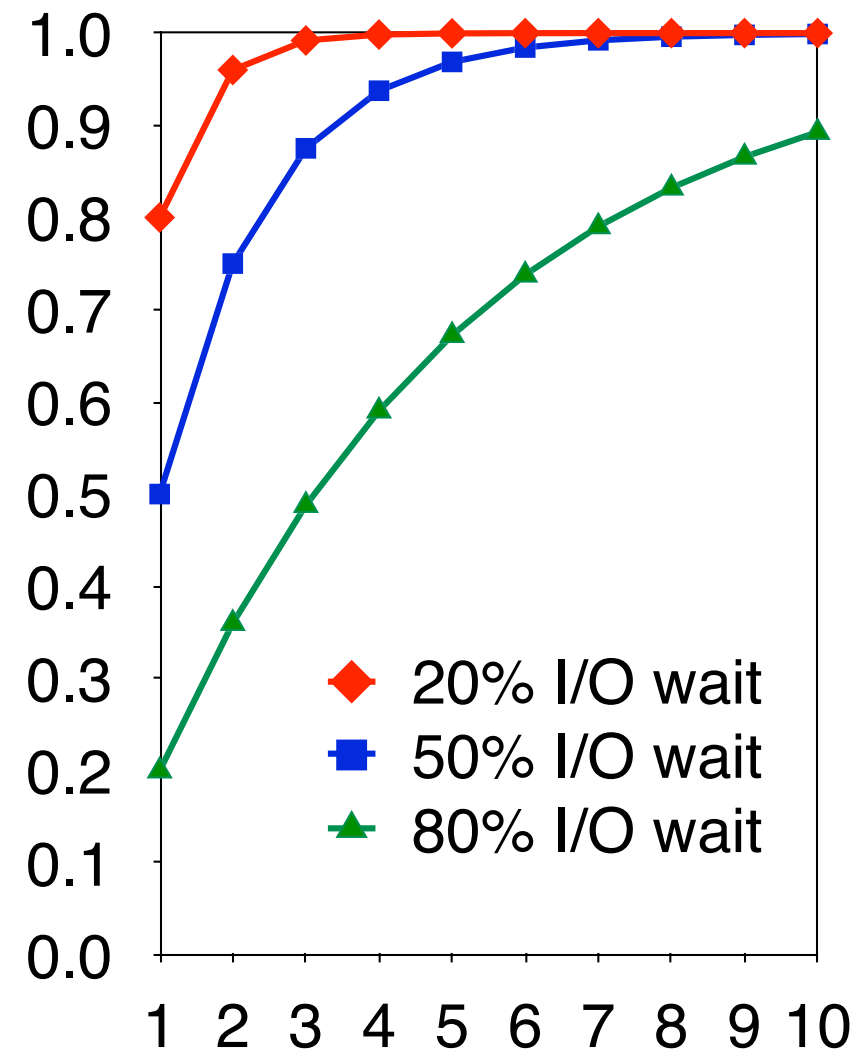


How many processes are enough?

- ✦ Several memory partitions (fixed or variable size)
- ✦ Lots of processes wanting to use the CPU
- ✦ Tradeoff
 - More processes utilize the CPU better
 - Fewer processes use less memory (cheaper!)
- ✦ How many processes do we need to keep the CPU fully utilized?
 - This will help determine how much memory we need
 - Is this still relevant with memory costing \$25/GB?

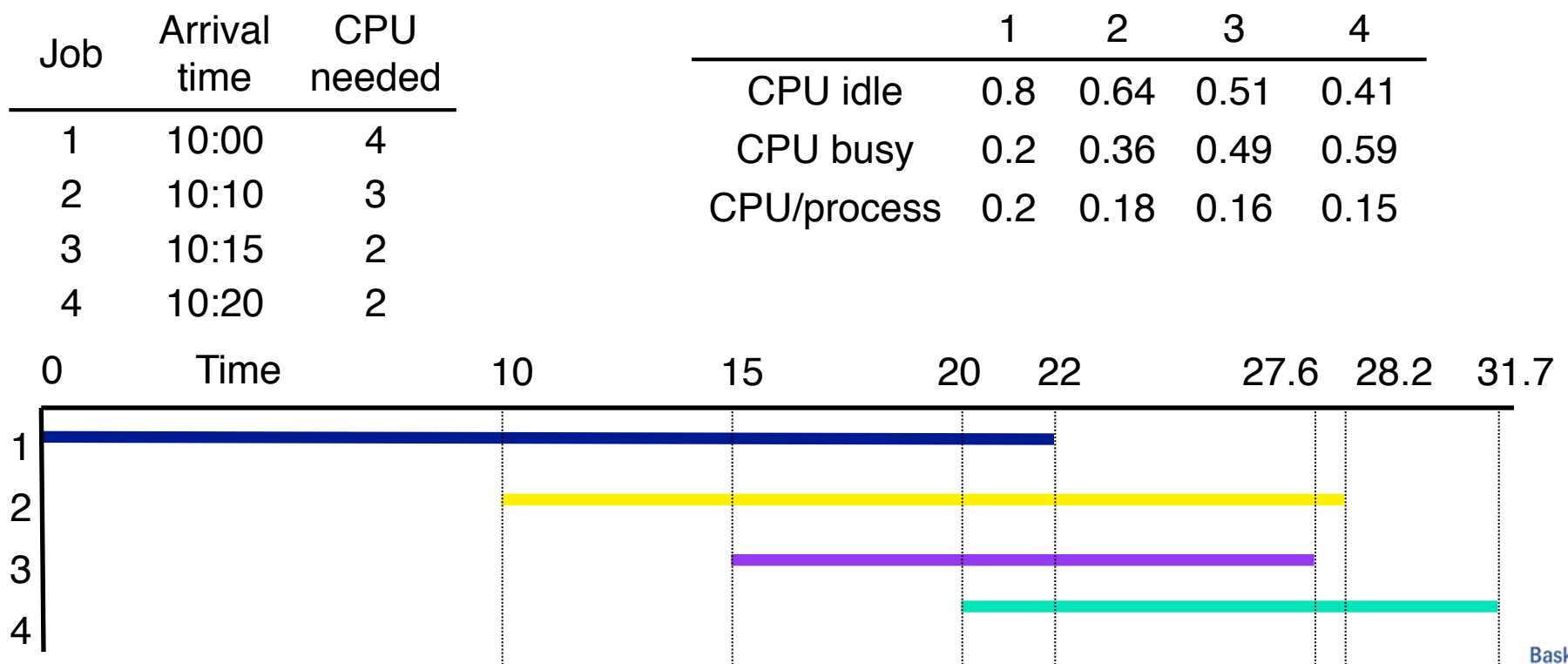
Modeling multiprogramming

- ◆ More I/O wait means less processor utilization
 - At 20% I/O wait, 3–4 processes fully utilize CPU
 - At 80% I/O wait, even 10 processes aren't enough
- ◆ This means that the OS should have more processes if they're I/O bound
- ◆ More processes → memory management & protection more important!



Multiprogrammed system performance

- ◆ Arrival and work requirements of 4 jobs
- ◆ CPU utilization for 1–4 jobs with 80% I/O wait
- ◆ Sequence of events as jobs arrive and finish
 - Numbers show amount of CPU time jobs get in each interval
 - More processes \Rightarrow better utilization, less time per process

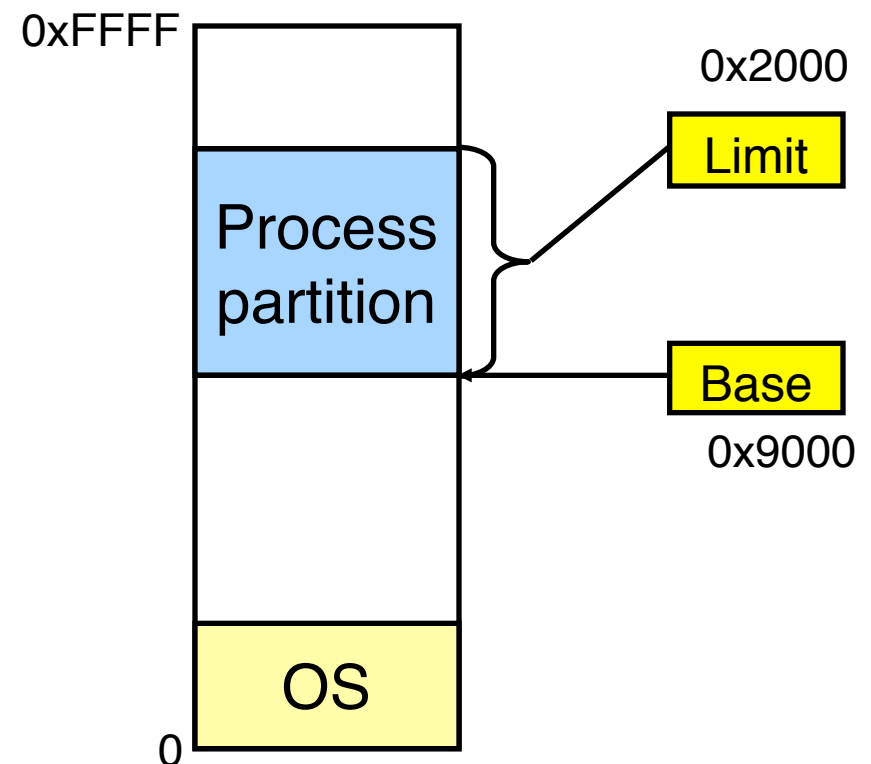


Memory and multiprocessing

- ✦ Memory needs two things for multiprocessing
 - Relocation
 - Protection
- ✦ The OS cannot be certain where a program will be loaded in memory
 - Variables and procedures can't use absolute locations in memory
 - Several ways to guarantee this
- ✦ The OS must keep processes' memory separate
 - Protect a process from other processes reading or modifying its own memory
 - Protect a process from modifying its own memory in undesirable ways (such as writing to program code)

Base and limit registers

- ◆ Special CPU registers: base & limit
 - Access to the registers limited to system mode
 - Registers contain
 - Base: start of the process's memory partition
 - Limit: length of the process's memory partition
- ◆ Address generation
 - Physical address: location in actual memory
 - Logical address: location from the process's point of view
 - Physical address = base + logical address
 - Logical address larger than limit → error

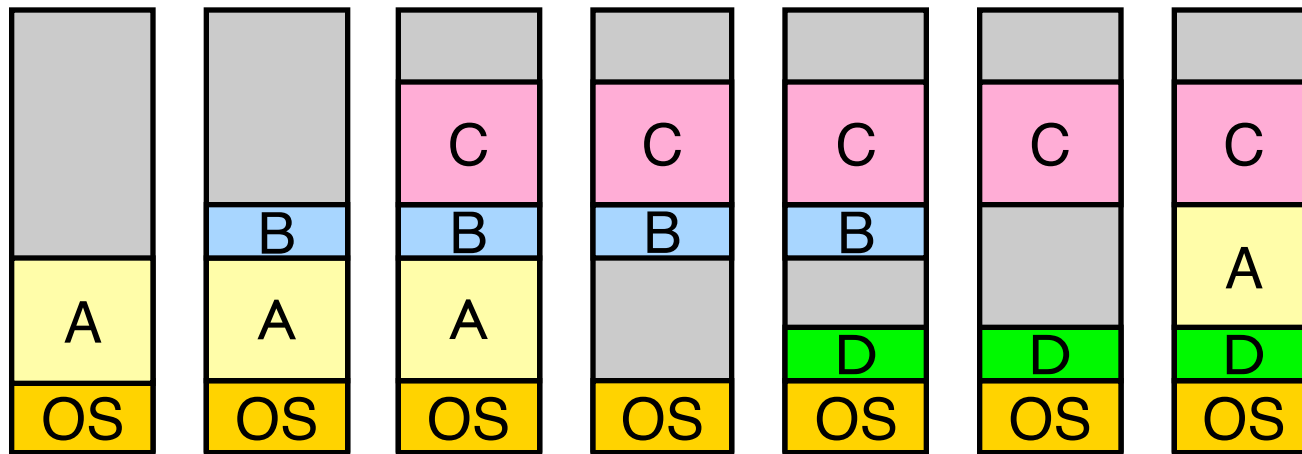


Logical address: 0x1204

Physical address:

$$0x1204 + 0x9000 = 0xa204$$

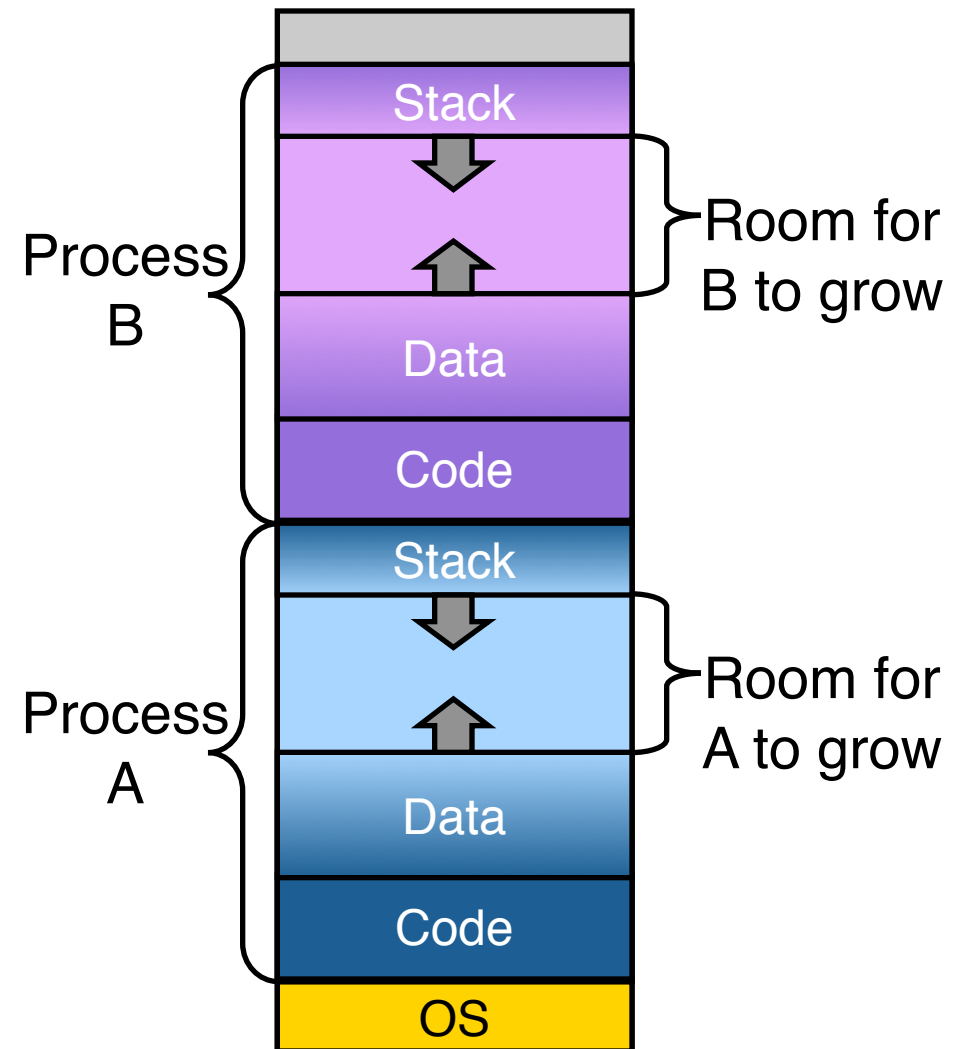
Swapping



- ✦ Memory allocation changes as
 - Processes come into memory
 - Processes leave memory
 - Swapped to disk
 - Complete execution
- ✦ Gray regions are unused memory

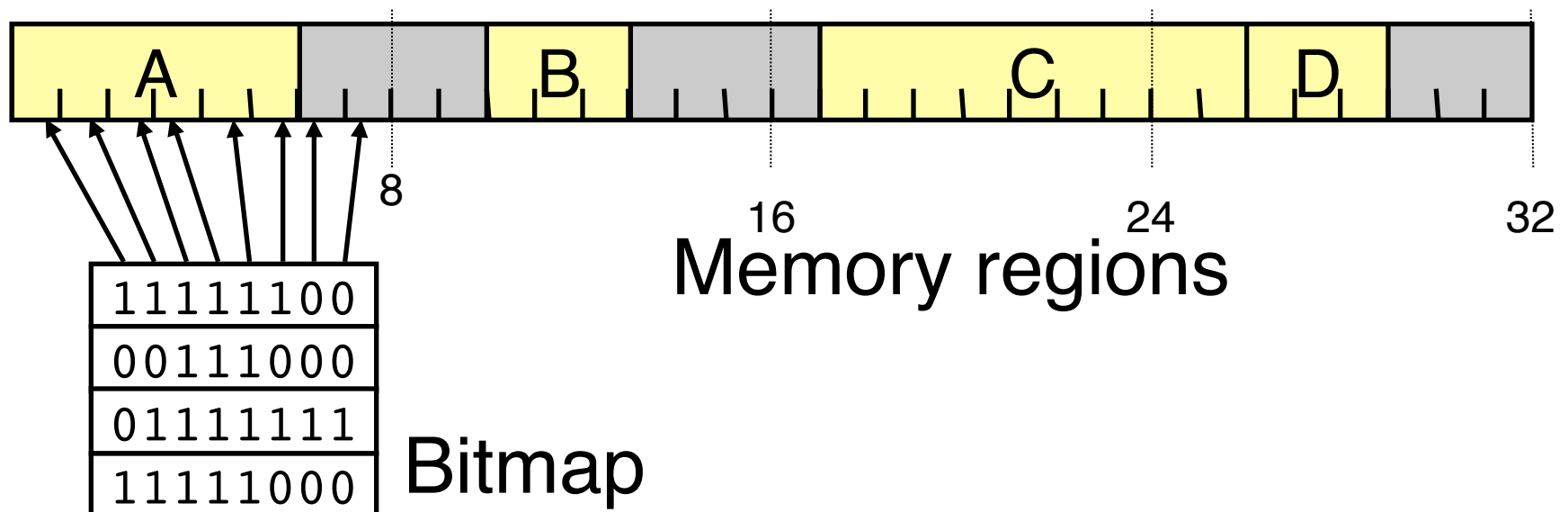
Swapping: leaving room to grow

- ◆ Need to allow for programs to grow
 - Allocate more memory for data
 - Larger stack
- ◆ Handled by allocating more space than is necessary at the start
 - Inefficient: wastes memory that's not currently in use
 - What if the process requests too much memory?



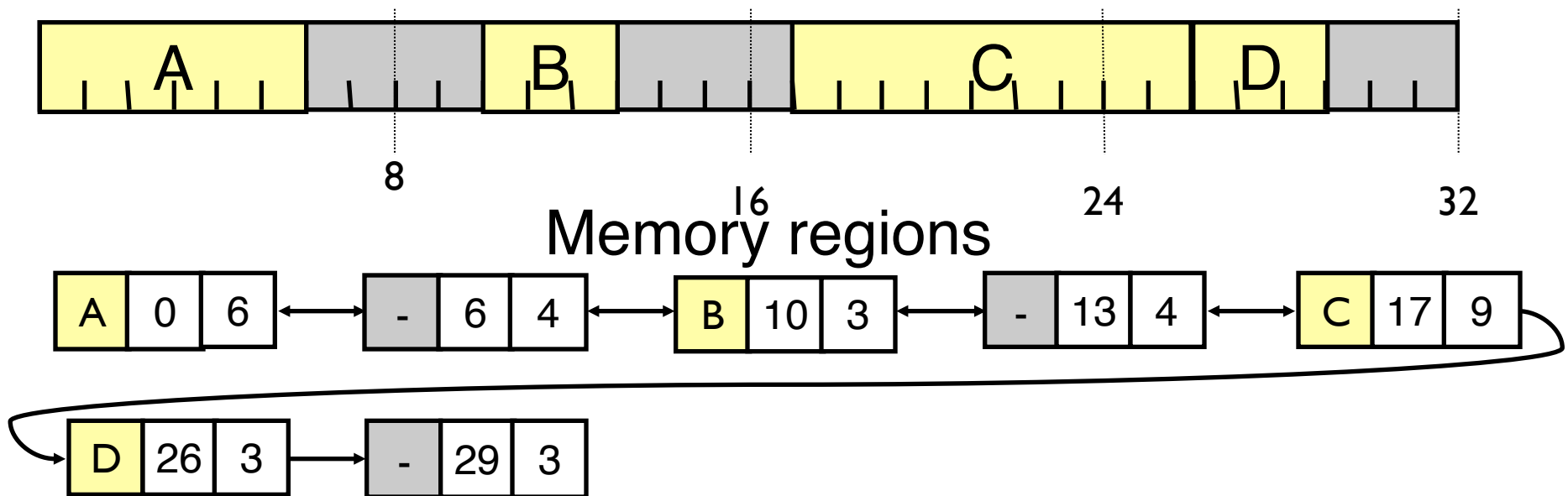
Tracking memory usage: bitmaps

- ◆ Keep track of free / allocated memory regions with a bitmap
 - One bit in map corresponds to a fixed-size region of memory
 - Bitmap is a constant size for a given amount of memory regardless of how much is allocated at a particular time
- ◆ Chunk size determines efficiency
 - At 1 bit per 4KB chunk, we need just 256 bits (32 bytes) per MB of memory
 - For smaller chunks, we need more memory for the bitmap
 - Can be difficult to find large contiguous free areas in bitmap



Tracking memory usage: linked lists

- ◆ Keep track of free / allocated memory regions with a linked list
 - Each entry in the list corresponds to a contiguous region of memory
 - Entry can indicate either allocated or free (and, optionally, owning process)
 - May have separate lists for free and allocated areas
- ◆ Efficient if chunks are large
 - Fixed-size representation for each region
 - More regions → more space needed for free lists

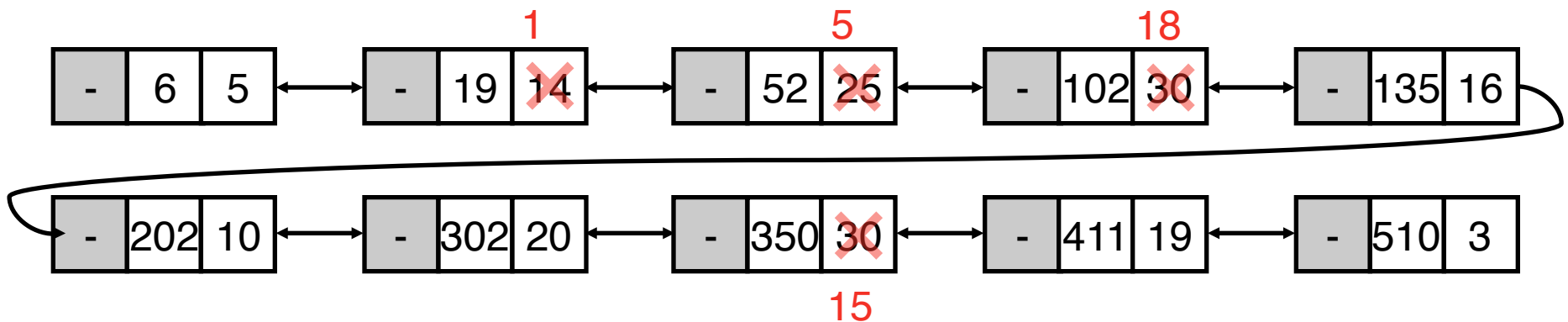


Allocating memory

- ♦ Search through region list to find a large enough space
- ♦ Suppose there are several choices: which one to use?
 - First fit: the first suitable hole on the list
 - Next fit: the first suitable after the previously allocated hole
 - Best fit: the smallest hole that is larger than the desired region (wastes least space?)
 - Worst fit: the largest available hole (leaves largest fragment)
- ♦ Option: maintain separate queues for different-size holes

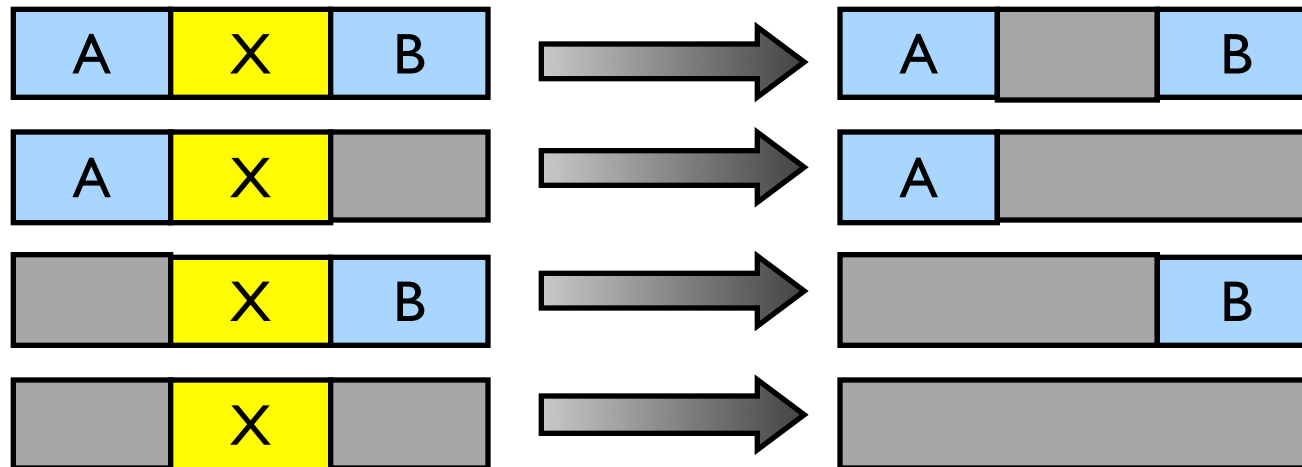
Allocate 20 blocks first fit
Allocate 12 blocks next fit

Allocate 13 blocks best fit
Allocate 15 blocks worst fit



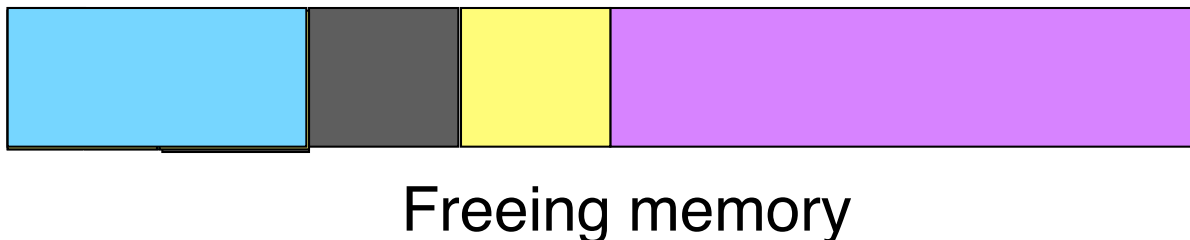
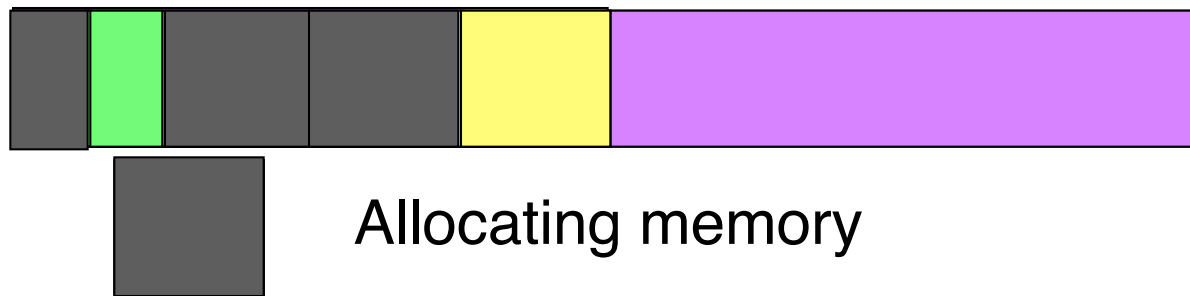
Freeing memory

- ✦ Allocation structures must be updated when memory is freed
- ✦ Easy with bitmaps: just set the appropriate bits in the bitmap
- ✦ Linked lists: modify adjacent elements as needed
 - Merge adjacent free regions into a single region
 - May involve merging two regions with the just-freed area



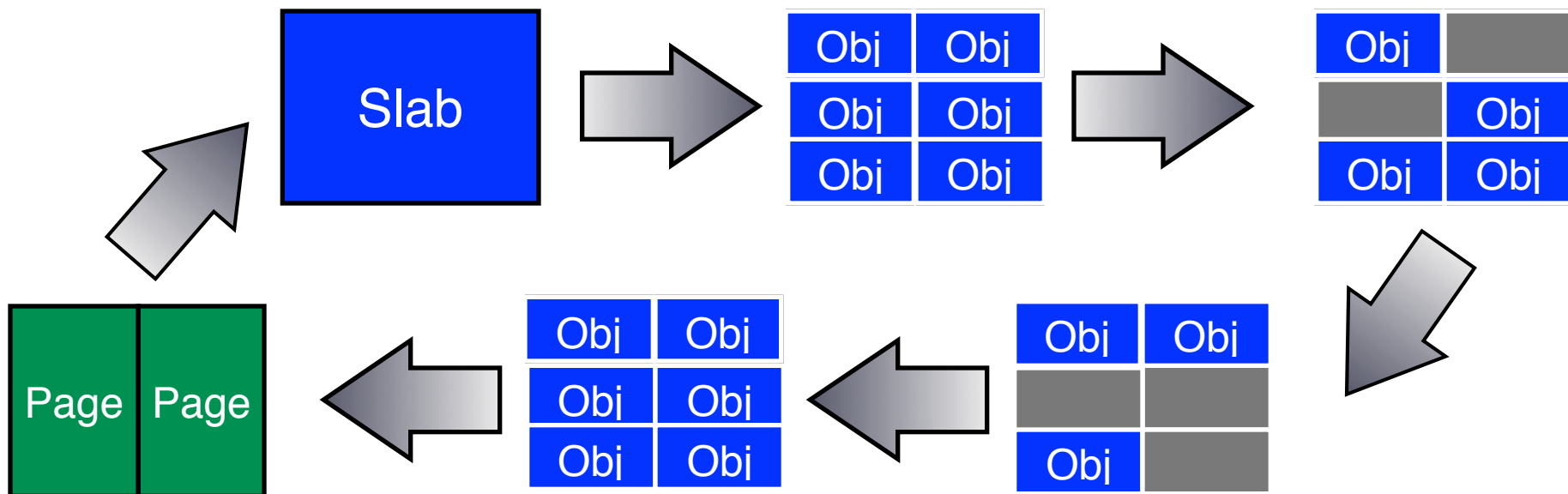
Buddy allocation

- ✦ Allocate memory in powers of two
 - Good for objects of varied sizes
- ✦ Split larger chunks to create two smaller chunks
- ✦ When chunk is freed, see if it can be combined with its buddy to rebuild a larger chunk
 - This is recursive...



Slab allocation

- ✦ Use *slabs* (each 1+ pages long) to allocate and manage a large number of equal-sized objects
 - Avoid internal fragmentation
 - Better performance: faster allocation / deallocation
- ✦ Carve slabs into objects
 - Free a slab only when all of its objects are freed
 - Keep track of free objects and hand out again when needed



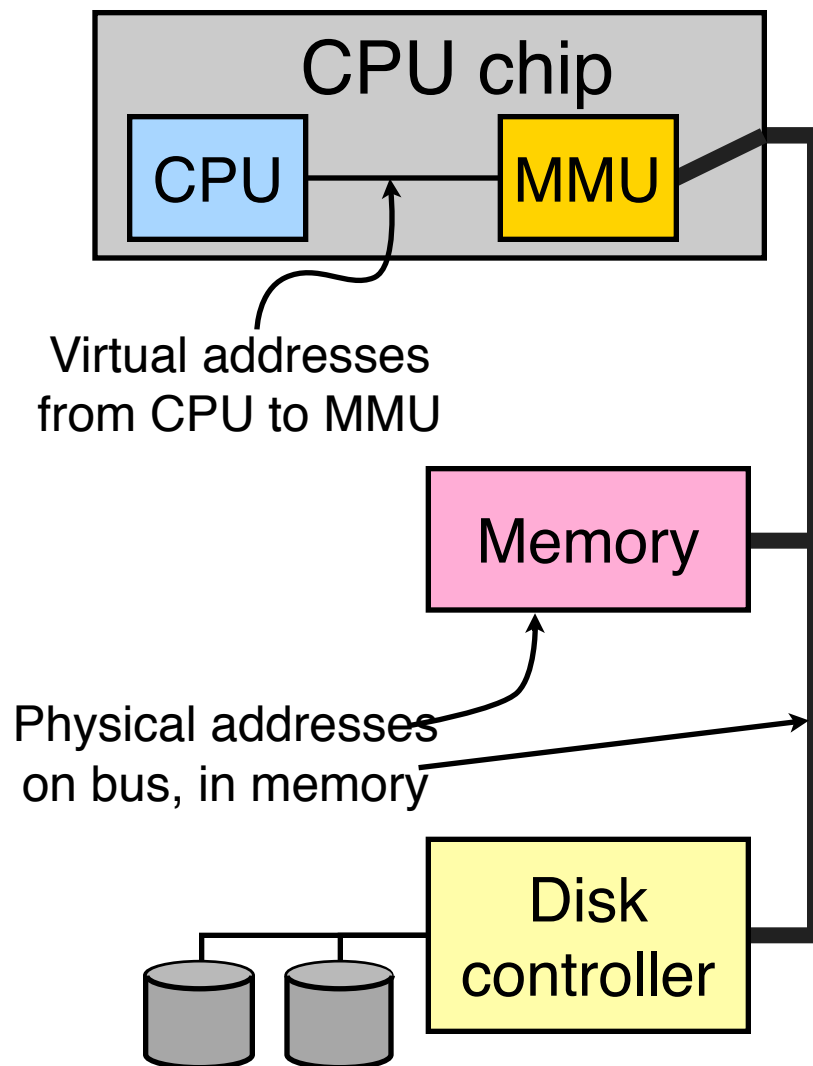
Limitations of swapping

- ◆ Problems with swapping
 - Process must fit into physical memory (impossible to run larger processes)
 - Memory becomes fragmented
 - External fragmentation: lots of small free areas
 - Compaction needed to reassemble larger free areas
 - Processes are either in memory or on disk: half and half doesn't do any good
- ◆ Overlays solved the first problem
 - Bring in pieces of the process over time (typically data)
 - Still doesn't solve the problem of fragmentation or partially resident processes

Virtual memory

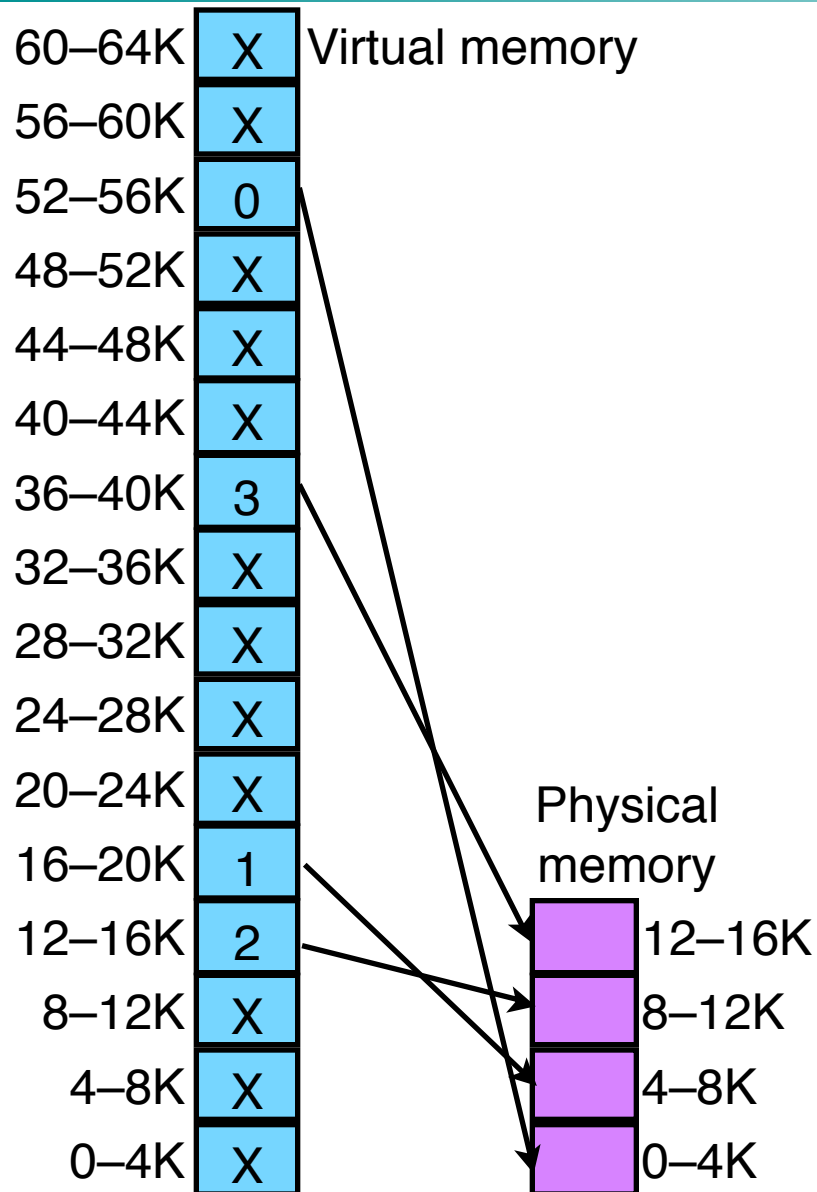
- ✦ Basic idea: allow the OS to hand out more memory than exists on the system
- ✦ Keep recently used stuff in physical memory
- ✦ Move less recently used stuff to disk
- ✦ Keep all of this hidden from processes
 - Processes still see an address space from 0–*max_address*
 - Movement of information to and from disk handled by the OS without process help
- ✦ Virtual memory (VM) especially helpful in multiprogrammed systems
 - CPU schedules process B while process A waits for its memory to be retrieved from disk

Virtual and physical addresses



- ✦ Program uses **virtual addresses**
 - Addresses local to the process
 - Hardware translates virtual address to physical address
- ✦ Translation done by the **Memory Management Unit**
 - Usually on the same chip as the CPU
 - Only physical addresses leave the CPU/MMU chip
- ✦ Physical memory indexed by physical addresses

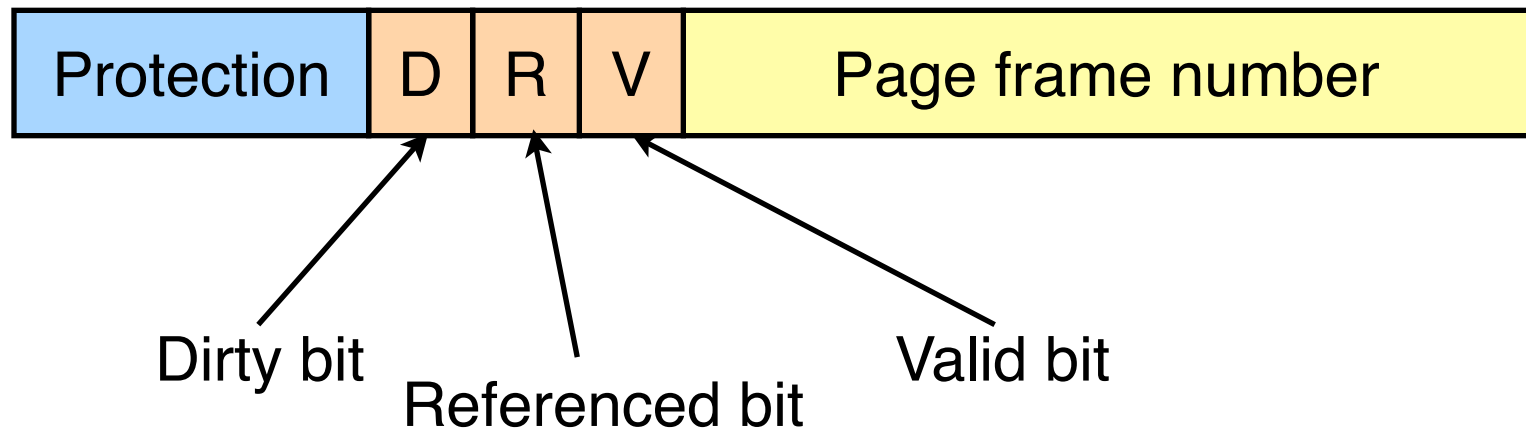
Paging and page tables



- ✦ Virtual addresses mapped to physical addresses
 - Unit of mapping is called a **page**
 - All addresses in the same virtual page are in the same physical page
 - **Page table entry (PTE)** contains translation for a single page
- ✦ Table translates virtual page number to physical page number
 - Not all virtual memory has a physical page
 - Not every physical page need be used
- ✦ Example:
 - 64 KB virtual memory
 - 16 KB physical memory

What's in a page table entry?

- ◆ Each entry in the page table contains
 - Valid bit: set if this logical page number has a corresponding physical frame in memory
 - If not valid, remainder of PTE is irrelevant
 - Page frame number: page in physical memory
 - Referenced bit: set if data on the page has been accessed
 - Dirty (modified) bit :set if data on the page has been modified
 - Protection information

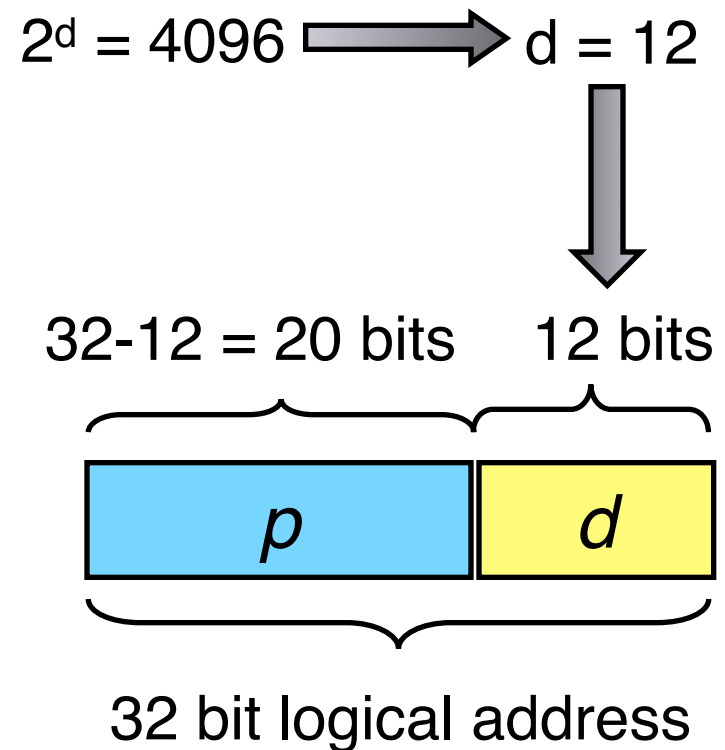


Mapping logical addresses to physical addresses

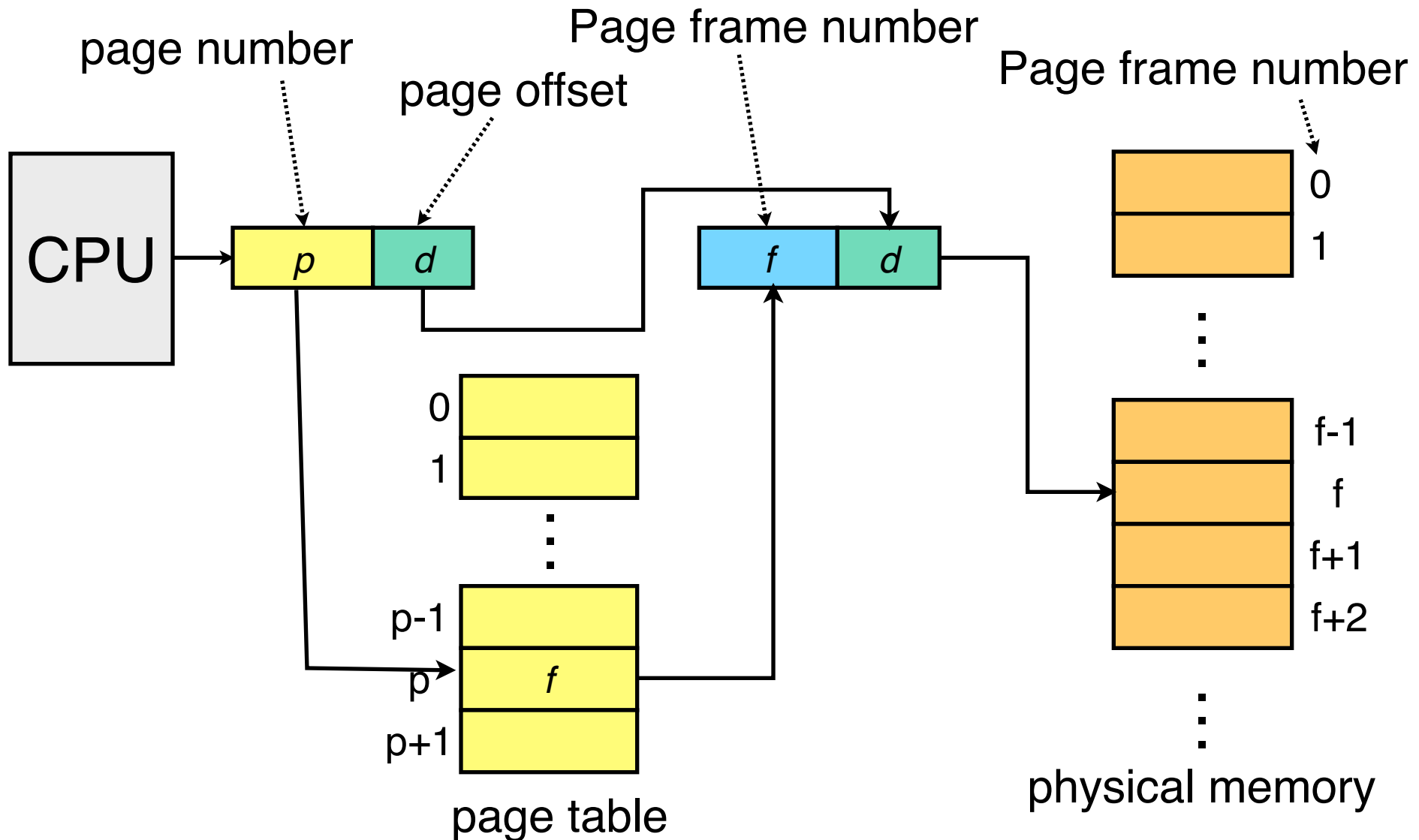
- ◆ Split address from CPU into two pieces
 - Page number (p)
 - Page offset (d)
- ◆ Page number
 - Index into page table
 - Page table contains base address of page in physical memory
- ◆ Page offset
 - Added to base address to get actual physical memory address
- ◆ Page size = 2^d bytes

Example:

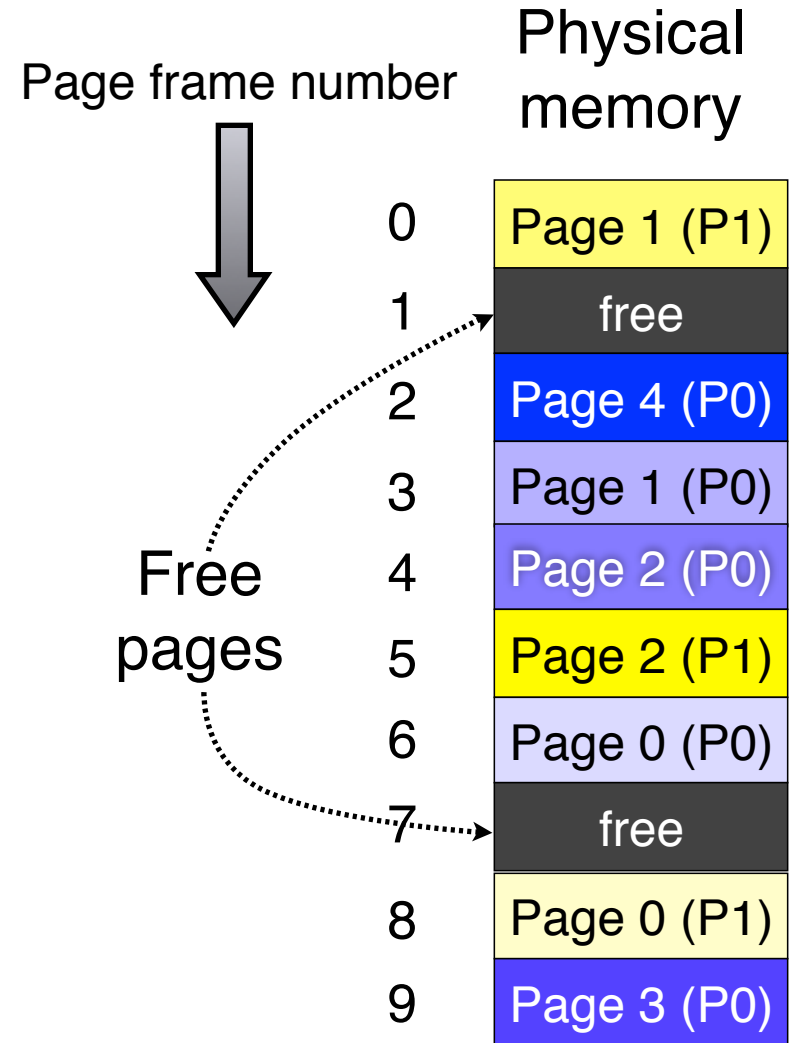
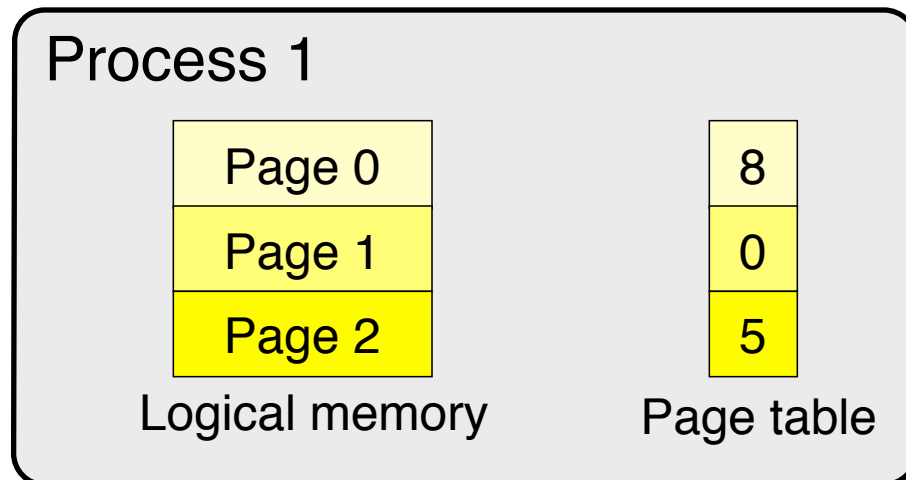
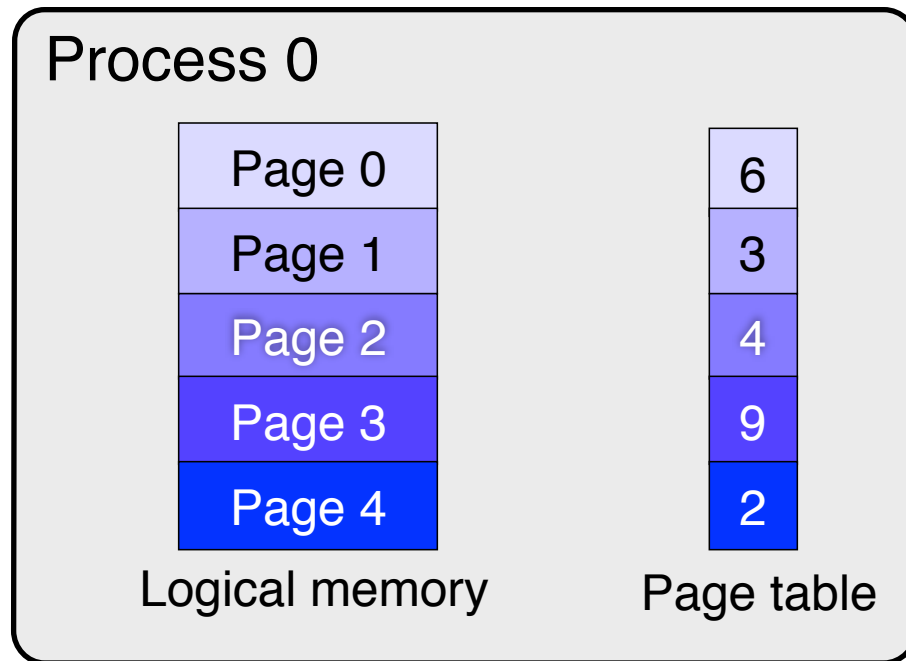
- 4 KB (=4096 byte) pages
- 32 bit logical addresses



Address translation architecture

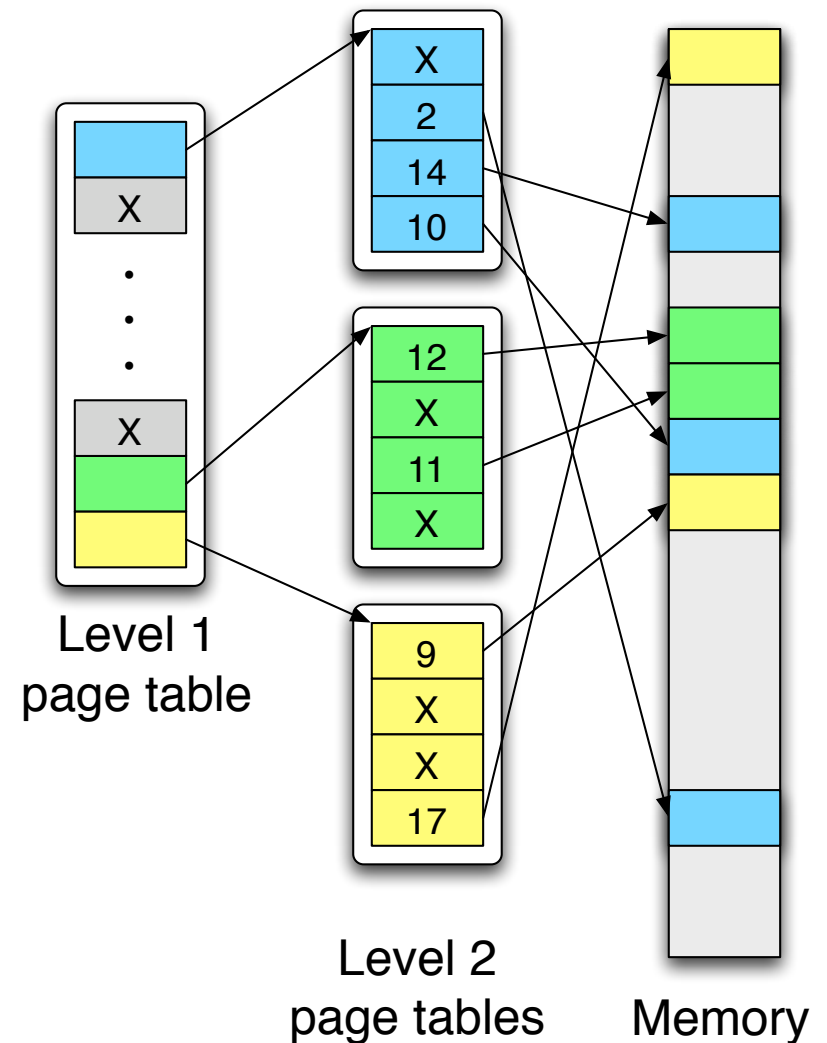


Memory & paging structures



Two-level page tables

- ◆ Problem: page tables can be too large
 - 2^{32} bytes in 4KB pages \Rightarrow 1 million PTEs
 - Worse for 64-bit addressing!
- ◆ Solution: use multi-level page tables
 - “Page size” in first page table is large (megabytes)
 - PTE marked invalid in first page table needs no 2nd level page table
- ◆ 1st level page table has pointers to 2nd level page tables
- ◆ 2nd level page table has actual physical page numbers in it

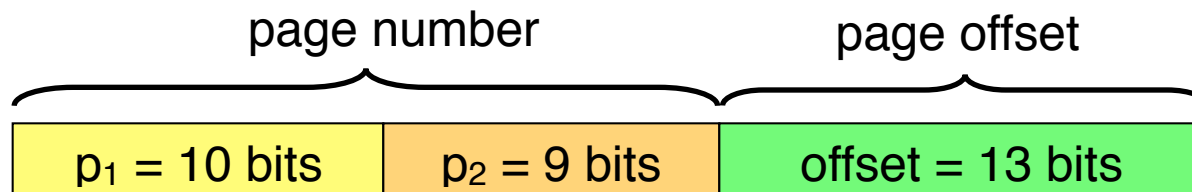


More on two-level page tables

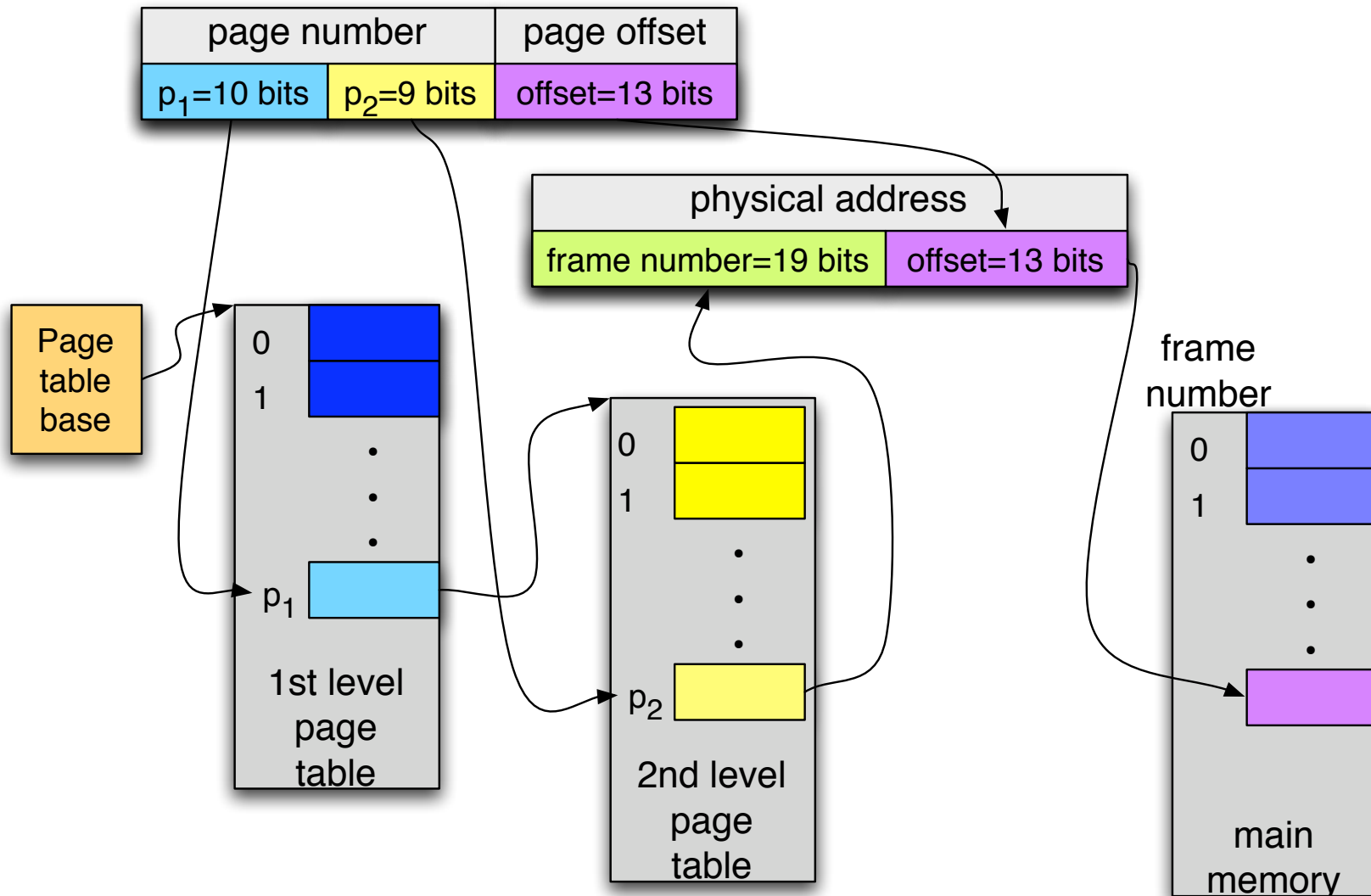
- ✦ Tradeoffs between 1st and 2nd level page table sizes
 - Total number of bits indexing 1st and 2nd level table is constant for a given page size and logical address length
 - Tradeoff between number of bits indexing 1st and number indexing 2nd level tables
 - More bits in 1st level: fine granularity at 2nd level
 - Fewer bits in 1st level: maybe less wasted space?
- ✦ All addresses in table are physical addresses
- ✦ Protection bits kept in 2nd level table

Two-level paging: example

- ◆ System characteristics
 - 8 KB pages
 - 32-bit logical address divided into 13 bit page offset, 19 bit page number
- ◆ Page number divided into:
 - 10 bit page number
 - 9 bit page offset
- ◆ Logical address looks like this:
 - p_1 is an index into the 1st level page table
 - p_2 is an index into the 2nd level page table pointed to by p_1



2-level address translation example



Implementing page tables in hardware

- ✦ Page table resides in main (physical) memory
- ✦ CPU uses special registers for paging
 - Page table base register (PTBR) points to the page table
 - Page table length register (PTLR) contains length of page table: restricts maximum legal logical address
- ✦ Translating an address requires two memory accesses
 - First access reads page table entry (PTE)
 - Second access reads the data / instruction from memory
- ✦ Reduce number of memory accesses
 - Can't avoid second access (we need the value from memory)
 - Eliminate first access by keeping a hardware cache (called a translation lookaside buffer or TLB) of recently used page table entries

Translation Lookaside Buffer (TLB)

- ✦ Search the TLB for the desired logical page number
 - Search entries in parallel
 - Use standard cache techniques
- ✦ If desired logical page number is found, get frame number from TLB
- ✦ If desired logical page number isn't found
 - Get frame number from page table in memory
 - Replace an entry in the TLB with the logical & physical page numbers from this reference

Logical page #	Physical frame #
8	3
invalid	
2	1
3	0
12	12
29	6
22	11
7	4

Example TLB



Handling TLB misses

- ✦ If PTE isn't found in TLB, OS needs to do the lookup in the page table
- ✦ Lookup can be done in hardware or software
- ✦ Hardware TLB replacement
 - CPU hardware does page table lookup
 - Can be faster than software
 - Less flexible than software, and more complex hardware
- ✦ Software TLB replacement
 - OS gets TLB exception
 - Exception handler does page table lookup & places the result into the TLB
 - Program continues after return from exception
 - Larger TLB (lower miss rate) can make this feasible

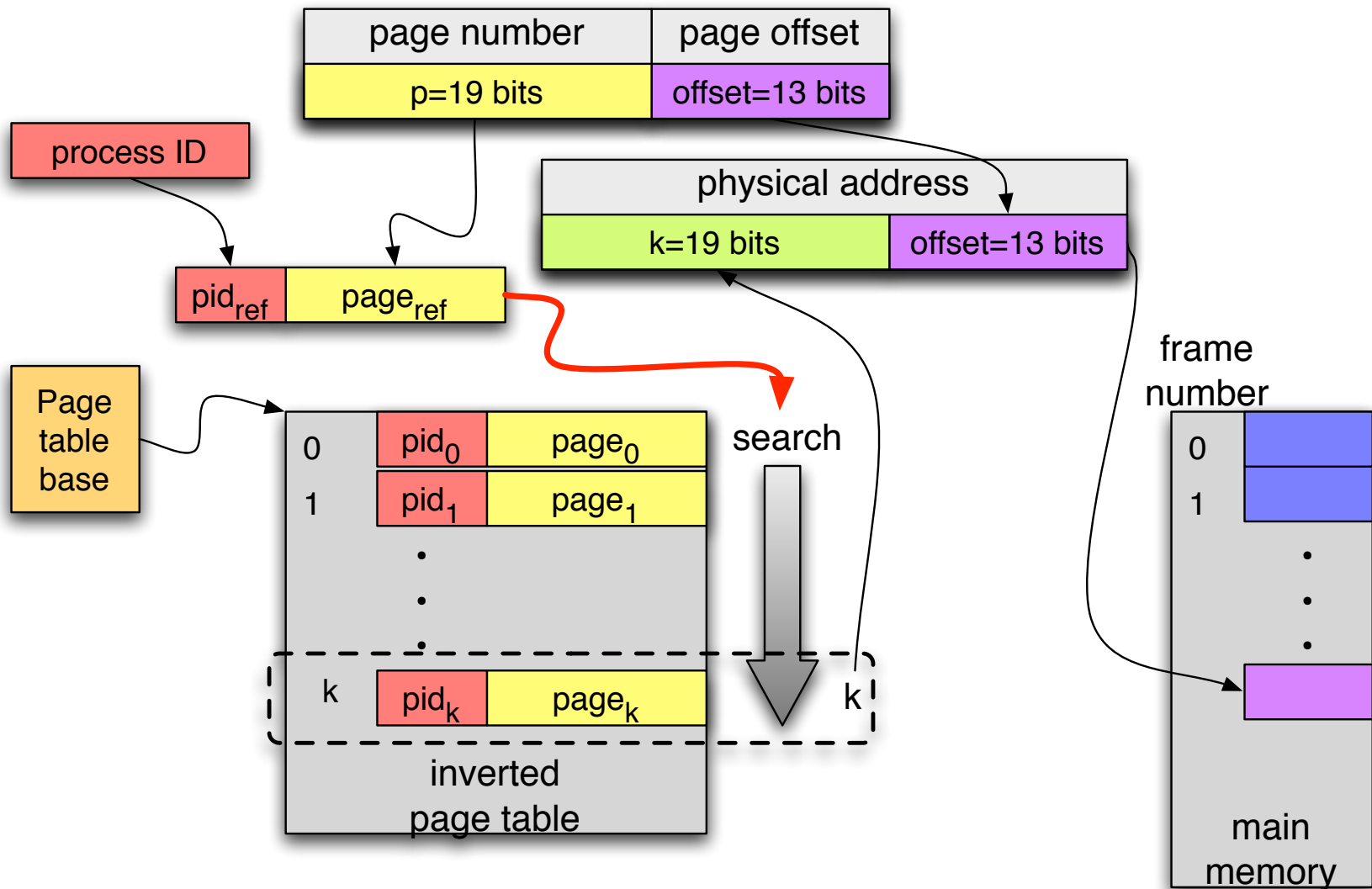
How long do memory accesses take?

- ✦ Assume the following times:
 - TLB lookup time = a (often zero—overlapped in CPU)
 - Memory access time = m
- ✦ Hit ratio (h) is percentage of time that a logical page number is found in the TLB
 - Larger TLB usually means higher h
 - TLB structure can affect h as well
- ✦ Effective access time (an average) is calculated as:
 - $EAT = (m + a)h + (m + m + a)(1-h)$
 - $EAT = a + (2-h)m$
- ✦ Interpretation
 - Reference always requires TLB lookup, 1 memory access
 - TLB misses also require an additional memory reference

Inverted page table

- ✦ Reduce page table size further: keep one entry for each frame in memory
 - Alternative: merge tables for pages in memory and on disk
- ✦ PTE contains
 - Virtual address pointing to this frame
 - Information about the process that owns this page
- ✦ Search page table by
 - Hashing the virtual page number and process ID
 - Starting at the entry corresponding to the hash result
 - Search until either the entry is found or a limit is reached
- ✦ Page frame number is index of PTE
- ✦ Improve performance by using more advanced hashing algorithms

Inverted page table architecture



One-to-one correspondence between page table entries and pages in memory

Memory Management

Part 2: Paging Algorithms and Implementation Issues

Page replacement algorithms

- ✦ Page fault forces a choice
 - No room for new page (steady state)
 - Which page must be removed to make room for an incoming page?
- ✦ How is a page removed from physical memory?
 - If the page is unmodified, simply overwrite it: a copy already exists on disk
 - If the page has been modified, it must be written back to disk: prefer unmodified pages?
- ✦ Better not to choose an often used page
 - It'll probably need to be brought back in soon

Optimal page replacement algorithm

- ✦ What's the best we can possibly do?
 - Assume perfect knowledge of the future
 - Not realizable in practice (usually)
 - Useful for comparison: if another algorithm is within 5% of optimal, not much more can be done...
- ✦ Algorithm: replace the page that will be used furthest in the future
 - Only works if we know the whole sequence!
 - Can be approximated by running the program twice
 - Once to generate the reference trace
 - Once (or more) to apply the optimal algorithm
- ✦ Nice, but not achievable in real systems!

Not-recently-used (NRU) algorithm

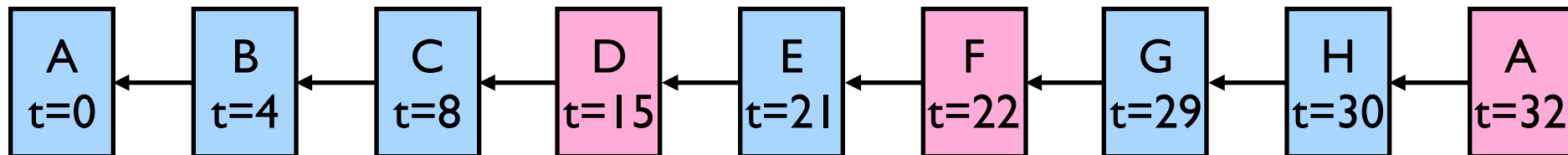
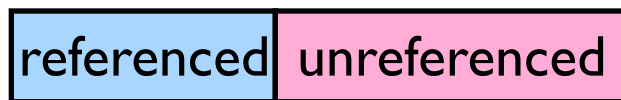
- ✦ Each page has reference bit and dirty bit
 - Bits are set when page is referenced and/or modified
- ✦ Pages are classified into four classes
 - 0: not referenced, not dirty
 - 1: not referenced, dirty
 - 2: referenced, not dirty
 - 3: referenced, dirty
- ✦ Clear reference bit for all pages periodically
 - Can't clear dirty bit: needed to indicate which pages need to be flushed to disk
 - Class 1 contains dirty pages where reference bit has been cleared
- ✦ Algorithm: remove a page from the lowest non-empty class
 - Select a page at random from that class
- ✦ Easy to understand and implement

First-In, First-Out (FIFO) algorithm

- ✦ Maintain a linked list of all pages
 - Maintain the order in which they entered memory
- ✦ Page at front of list replaced
- ✦ Advantage: (really) easy to implement
- ✦ Disadvantage: page in memory the longest may be often used
 - This algorithm forces pages out regardless of usage
 - Usage may be helpful in determining which pages to keep

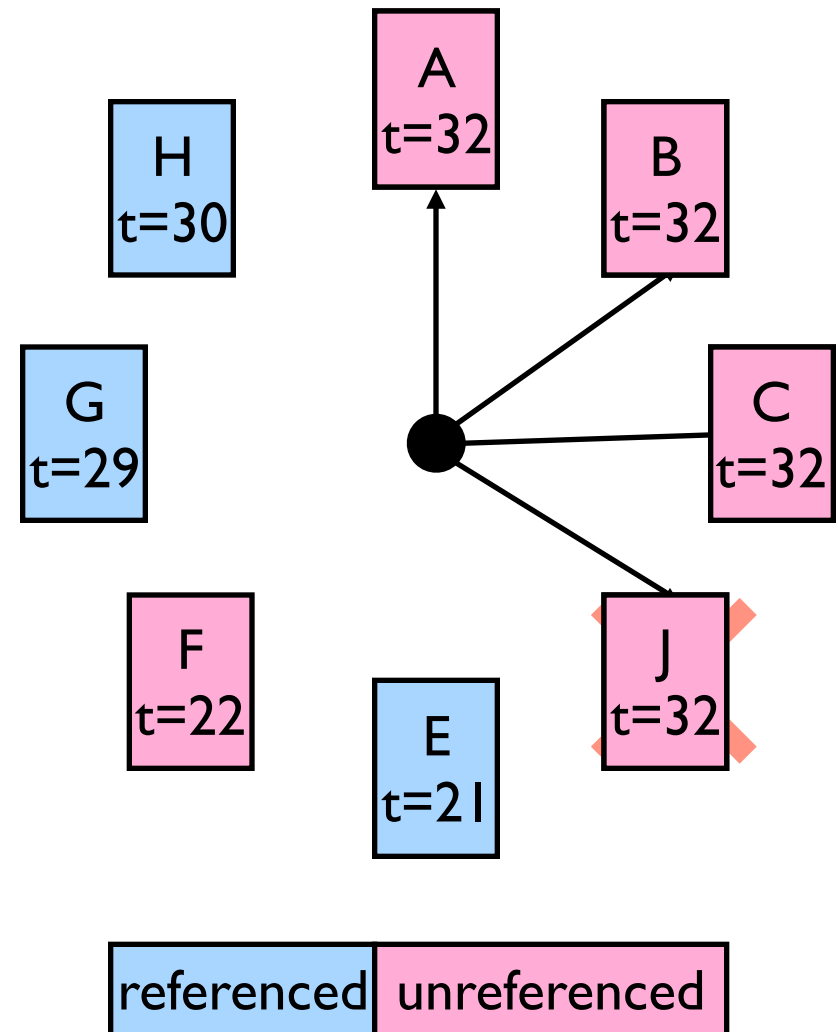
Second chance page replacement

- ✦ Modify FIFO to avoid throwing out heavily used pages
 - If reference bit is 0, throw the page out
 - If reference bit is 1
 - Reset the reference bit to 0
 - Move page to the tail of the list



Clock algorithm

- ✦ Same functionality as second chance
- ✦ Simpler implementation
 - “Clock” hand points to next page to replace
 - If $R=0$, replace page
 - If $R=1$, set $R=0$ and advance the clock hand
- ✦ Continue until page with $R=0$ is found
 - This may involve going all the way around the clock...



Least Recently Used (LRU)

- ✦ Assume pages used recently will be used again soon
 - Throw out page that has been unused for longest time
- ✦ Must keep a linked list of pages
 - Most recently used at front, least at rear
 - Update this list every memory reference!
 - This can be somewhat slow: hardware has to update a linked list on every reference!
- ✦ Alternatively, keep counter in each page table entry
 - Global counter increments with each CPU cycle

Simulating LRU in software

- ✦ Few computers have the necessary hardware to implement full LRU
 - Linked-list method impractical in hardware
 - Counter-based method could be done, but it's slow to find the desired page
- ✦ Approximate LRU with Not Frequently Used (NFU) algorithm
 - At each clock interrupt, scan through page table
 - If $R=1$ for a page, add one to its counter value
 - On replacement, pick the page with the lowest counter value
- ✦ Problem: no notion of age—pages with high

Aging replacement algorithm

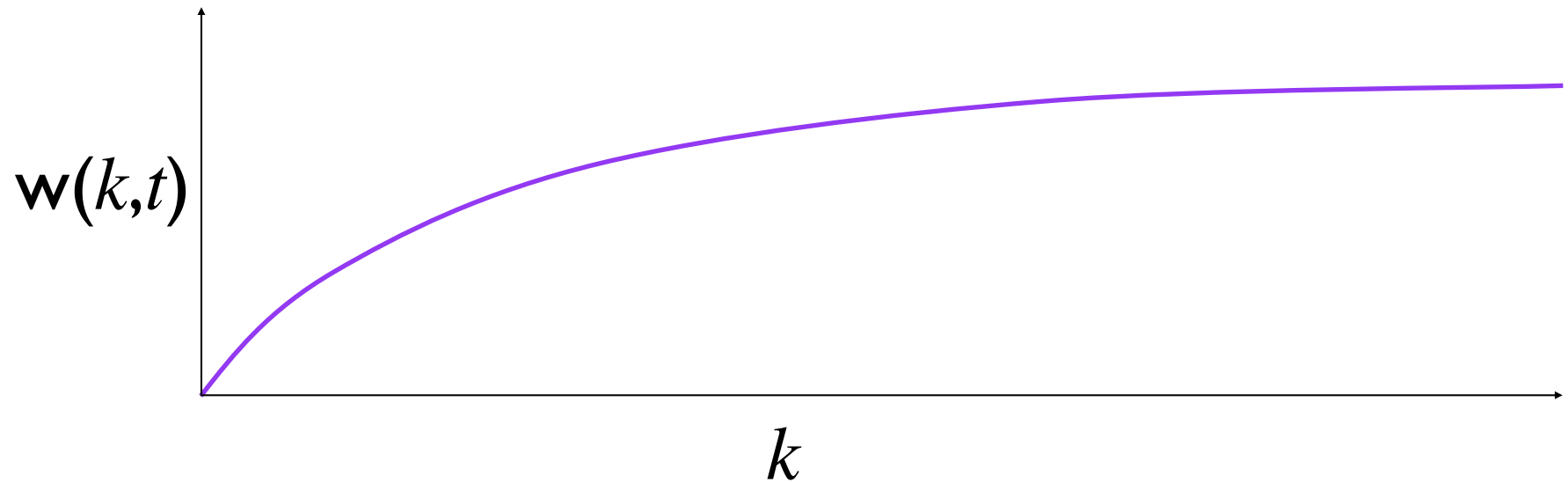
- ✦ Reduce counter values over time
 - Divide by two every clock cycle (use right shift)
 - More weight given to more recent references!
- ✦ Select page to be evicted by finding the lowest counter value
- ✦ Algorithm is:
 - Every clock tick, shift all counters right by 1 bit
 - On reference, set leftmost bit of a counter (can be done by copying the reference bit to the counter at the clock tick)

Referenced this tick	Tick 0	Tick 1	Tick 2	Tick 3	Tick 4
Page 0	10000000	11000000	11100000	01110000	10111000
Page 1	00000000	10000000	01000000	00100000	00010000
Page 2	10000000	01000000	00100000	10010000	01001000
Page 3	00000000	00000000	00000000	10000000	01000000
Page 4	10000000	01000000	10100000	11010000	01101000
Page 5	10000000	11000000	01100000	10110000	11011000

Working set

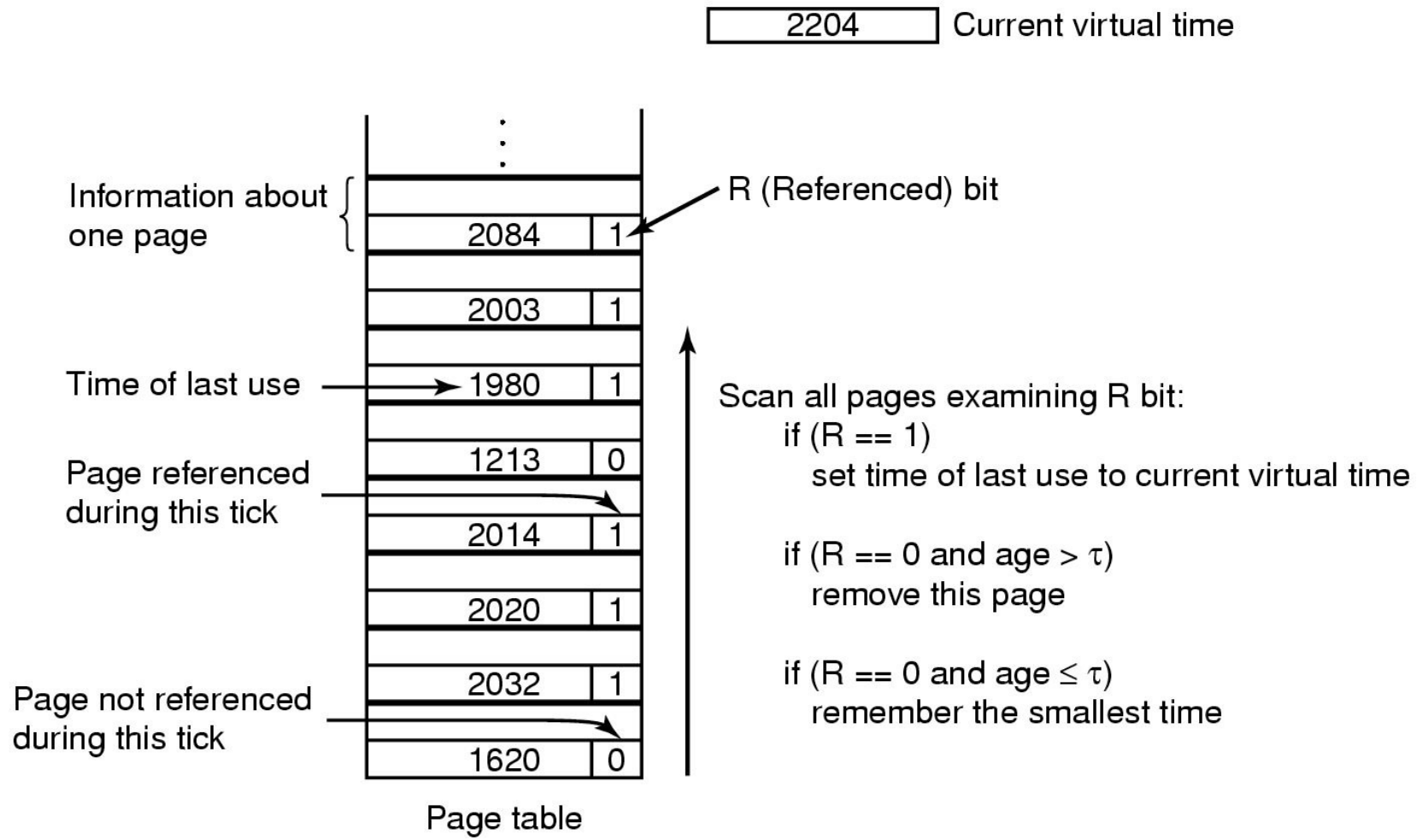
- ✦ Demand paging: bring a page into memory when it's requested by the process
- ✦ How many pages are needed?
 - Could be all of them, but not likely
 - Instead, processes reference a small set of pages at any given time—locality of reference
 - Set of pages can be different for different processes or even different times in the running of a single process
- ✦ Set of pages used by a process in a given interval of time is called the **working set**
 - If entire working set is in memory, no page faults!
 - If insufficient space for working set, thrashing may occur
 - Goal: keep most of working set in memory to minimize the number of page faults suffered by a process

How big is the working set?



- ✦ Working set is the set of pages used by the k most recent memory references
- ✦ $w(k, t)$ is the size of the working set at time t
- ✦ Working set may change over time
 - Size of working set can change over time as well...

Working set page replacement algorithm



Question

Which of these paging algorithms is impractical to use in a modern operating system?

Page replacement algorithms: summary

Algorithm	Comment
OPT (Optimal)	Not implementable, but useful as a benchmark
NRU (Not Recently Used)	Crude
FIFO (First-In, First Out)	Might throw out useful pages
Second chance	Big improvement over FIFO
Clock	Better implementation of second chance
LRU (Least Recently Used)	Excellent, but hard to implement exactly
NFU (Not Frequently Used)	Poor approximation to LRU
Aging	Good approximation to LRU, inefficient to implement
Working Set	Somewhat expensive to implement
WSClock	Implementable version of Working Set

Modeling page replacement algorithms

- ✦ Goal: provide quantitative analysis (or simulation) showing which algorithms do better
 - Workload (page reference string) is important: different strings may favor different algorithms
 - Show tradeoffs between algorithms
- ✦ Compare algorithms to one another
- ✦ Model parameters within an algorithm
 - Number of available physical pages
 - Number of bits for aging

How is modeling done?

- ♦ Generate a list of references
 - Artificial (made up)
 - Trace a real workload (set of processes)
- ♦ Use an array (or other structure) to track the pages in physical memory at any given time
 - May keep other information per page to help simulate the algorithm (modification time, time when paged in, etc.)
- ♦ Run through references, applying the replacement algorithm
- ♦ Example: FIFO replacement on reference string 0 1 2 3 0 1 4 0 1 2 3 4
 - Page replacements highlighted in yellow

Page referenced	0	1	2	3	0	1	4	0	1	2	3	4
Youngest page	0	1	2	3	0	1	4	4	4	2	3	3
		0	1	2	3	0	1	1	1	4	2	2
Oldest page			0	1	2	3	0	0	0	1	4	4

Belady's anomaly

- ✦ Reduce the number of page faults by supplying more memory
 - Use previous reference string and FIFO algorithm
 - Add another page to physical memory (total 4 pages)
- ✦ More page faults (10 vs. 9), not fewer!
 - This is called Belady's anomaly
 - Adding more pages shouldn't result in worse performance!
- ✦ Motivated the study of paging algorithms

Page referenced	0	1	2	3	0	1	4	0	1	2	3	4
Youngest page	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 page				0	0	0	1	2	3	4	0	1

Modeling more replacement algorithms

- ✦ Paging system characterized by:
 - Reference string of executing process
 - Page replacement algorithm
 - Number of page frames available in physical memory (m)
- ✦ Model this by keeping track of all n pages referenced in array M
 - Top part of M has m pages in memory
 - Bottom part of M has $n-m$ pages stored on disk
- ✦ Page replacement occurs when page moves from top to bottom
 - Top and bottom parts may be rearranged without

Example: LRU

- ✦ Model LRU replacement with
 - 8 unique references in the reference string
 - 4 pages of physical memory
- ✦ Array state over time shown below
- ✦ LRU treats list of pages like a stack

Page referenced	0	2	1	3	5	4	6	3	7	4	7	3	3	5	5	3	1	1	1	7	1	3	4	1
	0	2	1	3	5	4	6	3	7	4	7	3	3	5	5	3	1	1	1	7	1	3	4	1
Pages in RAM		0	2	1	3	5	4	6	3	7	4	7	7	3	3	5	3	3	3	1	7	1	3	4
			0	2	1	3	5	4	6	3	3	4	4	7	7	7	5	5	5	3	3	7	1	3
				0	2	1	3	5	4	6	6	6	6	4	4	4	7	7	7	5	5	5	7	7
					0	2	1	1	5	5	5	5	5	6	6	6	4	4	4	4	4	4	5	5
Pages on disk						0	2	2	1	1	1	1	1	1	1	1	6	6	6	6	6	6	6	6
							0	0	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
									0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Stack algorithms

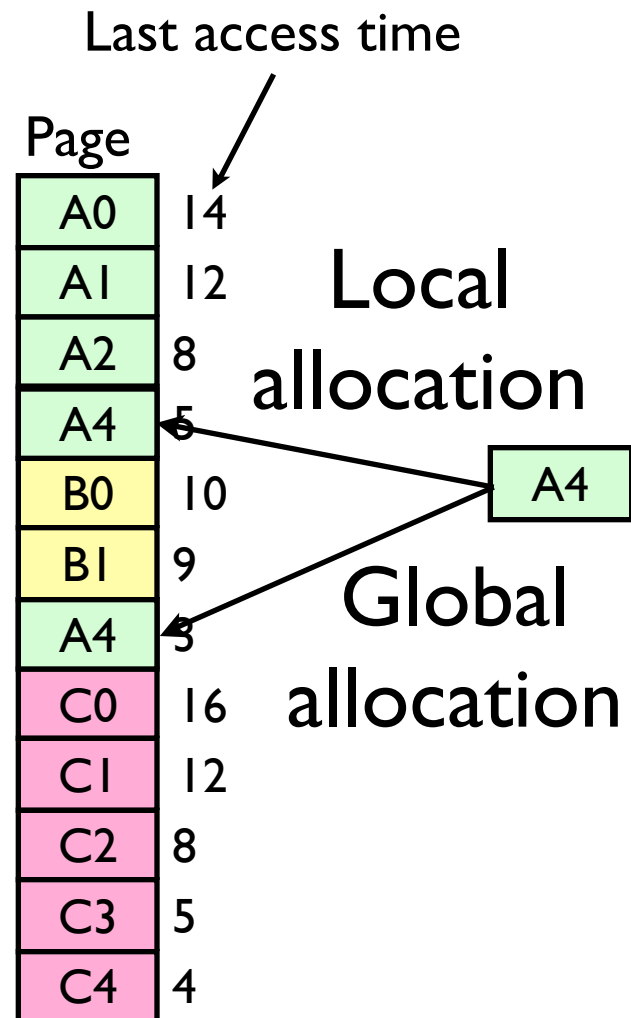
- ✦ LRU is an example of a stack algorithm
- ✦ For stack algorithms
 - Any page in memory with m physical pages is also in memory with $m+1$ physical pages
 - Increasing memory size is guaranteed to reduce (or at least not increase) the number of page faults
- ✦ Stack algorithms do not suffer from Belady's anomaly
- ✦ Distance of a reference == position of the page in the stack before the reference was made
 - Distance is ∞ if no reference had been made before
 - Distance depends on reference string and paging

Predicting page fault rates using distance

- ✦ Distance can be used to predict page fault rates
- ✦ Make a single pass over the reference string to generate the distance string on-the-fly
- ✦ Keep an array of counts
 - Entry j counts the number of times distance j occurs in the distance string
- ✦ The number of page faults for a memory of size m is the sum of the counts for $j > m$
 - This can be done in a single pass!
 - Makes for fast simulations of page replacement algorithms

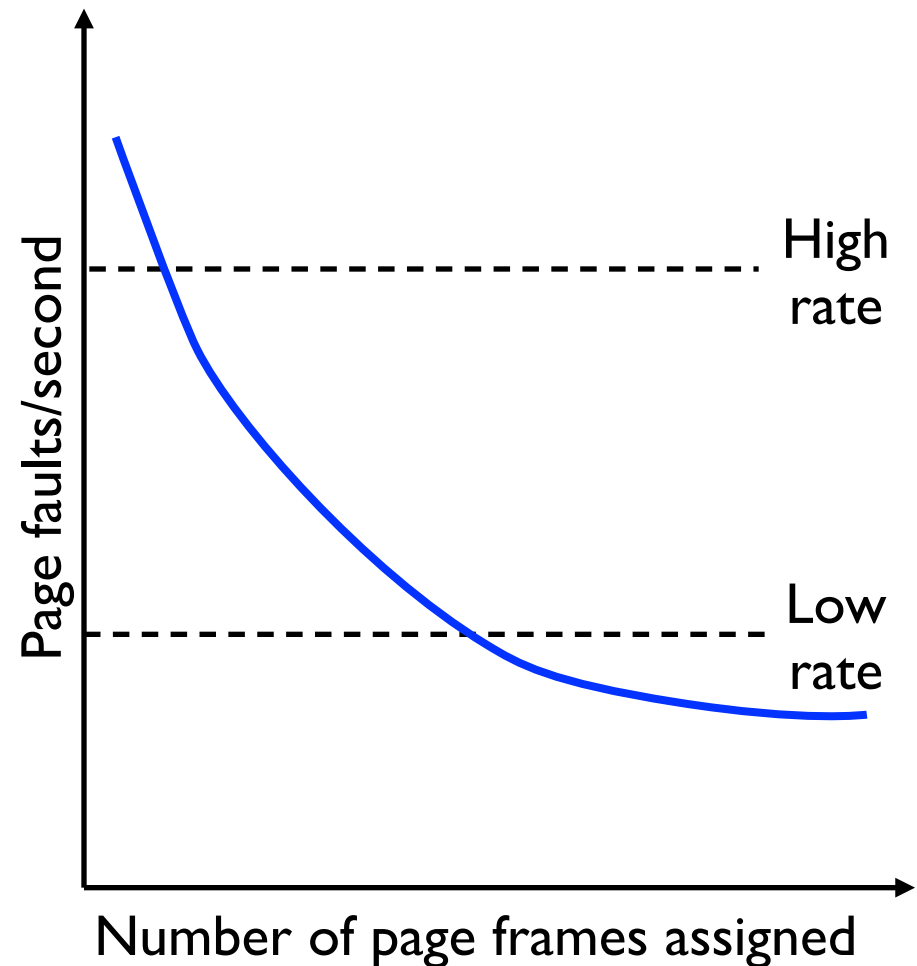
Local vs. global allocation policies

- What is the pool of pages eligible to be replaced?
 - Pages belonging to the process needing a new page
 - All pages in the system
- Local allocation: replace a page from this process
 - May be more “fair”: penalize processes that replace many pages
 - Can lead to poor performance: some processes need more pages than others
- Global allocation: replace a page from any process



Page fault rate vs. allocated frames

- ◆ Local allocation may be more “fair”
 - Don’t penalize other processes for high page fault rate
- ◆ Global allocation is better for overall system performance
 - Take page frames from processes that don’t need them as much
 - Reduce the overall page fault rate (even though rate



Control overall page fault rate

- ✦ Despite good designs, system may still thrash
- ✦ Most (or all) processes have high page fault rate
 - Some processes need more memory, ...
 - but no processes need less memory (and could give some up)
- ✦ Problem: no way to reduce page fault rate
- ✦ Solution :
Reduce number of processes competing for memory
 - Swap one or more to disk, divide up pages they held
 - Reconsider degree of multiprogramming

Question

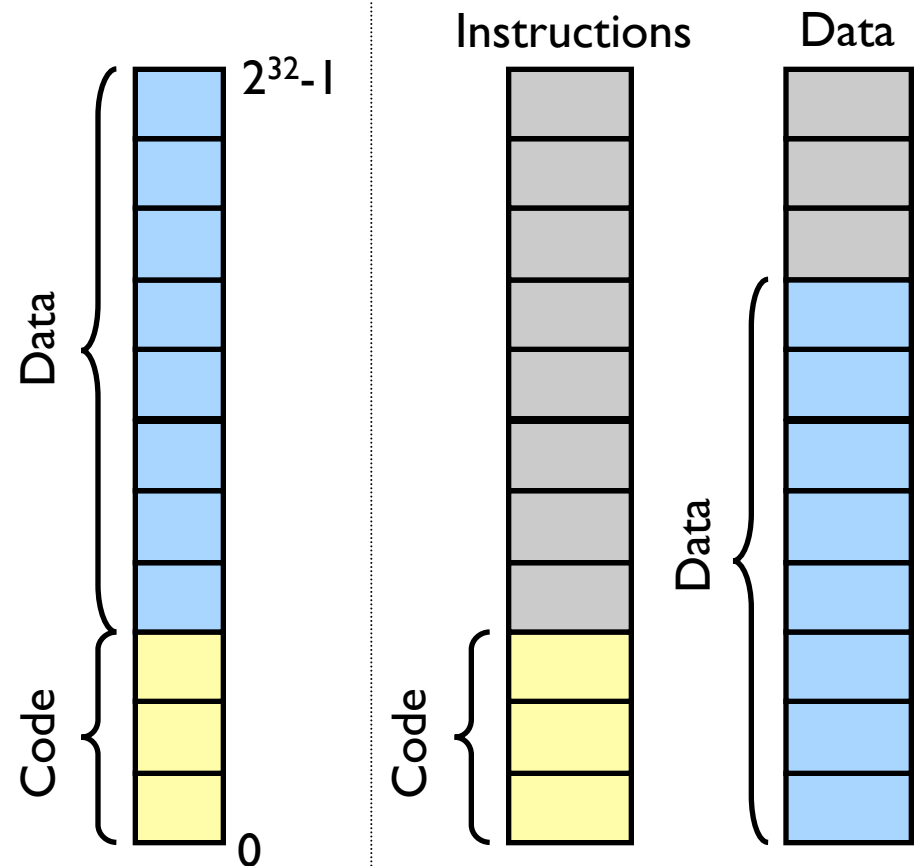
- ◆ How big should a page be, and why?
- A. Relatively large (8KB+): page tables are smaller
- B. Relatively large (8KB+): fewer page faults
- C. Relatively large (8KB+): paging algorithms run faster
- D. Relatively small (1-4KB): less internal fragmentation
- E. Relatively small (1-4KB): more accurate protection

How big should a page be?

- ✦ Smaller pages have advantages
 - Less internal fragmentation
 - Better fit for various data structures, code sections
 - Less unused physical memory (some pages have 20 useful bytes and the rest isn't needed currently)
- ✦ Larger pages are better because
 - Less overhead to keep track of them
 - Smaller page tables
 - TLB can point to more memory (same number of pages, but more memory per page)
 - Faster paging algorithms (fewer table entries to look through)
 - More efficient to transfer larger pages to and from

Separate I & D address spaces

- ◆ One user address space for both data & code
 - Simpler
 - Code/data separation harder to enforce
 - More address space?
- ◆ One address space for data, another for code
 - Code & data separated
 - More complex in hardware
 - Less flexible
 - CPU must handle instructions & data



Sharing pages

- ✦ Processes can share pages
 - Entries in page tables point to the same physical page frame
 - Easier to do with code: no problems with modification
- ✦ Virtual addresses in different processes can be...
 - The same: easier to exchange pointers, keep data structures consistent
 - Different: may be easier to actually implement
 - Not a problem if there are only a few shared regions
 - Can be very difficult if many processes share regions with each other

When are dirty pages written to disk?

- ✦ On demand (when they're replaced)
 - Fewest writes to disk
 - Slower: replacement takes twice as long (must wait for disk write and disk read)
- ✦ Periodically (in the background)
 - Background process scans through page tables, writes out dirty pages that are pretty old
- ✦ Background process also keeps a list of pages ready for replacement
 - Page faults handled faster: no need to find space on demand
 - Cleaner may use the same structures discussed

Implementation issues

- ✦ Four times when OS involved with paging
- ✦ Process creation
 - Determine program size
 - Create page table
- ✦ During process execution
 - Reset the MMU for new process
 - Flush the TLB (or reload it from saved state)
- ✦ Page fault time
 - Determine virtual address causing fault
 - Swap target page out, needed page in
- ✦ Process termination time
 - Release page table

How is a page fault handled?

- ◆ Hardware causes a page fault
- ◆ General registers saved (as on every exception)
- ◆ OS determines which virtual page needed
 - Actual fault address in a special register
 - Address of faulting instruction in register
 - Page fault was in fetching instruction, or
 - Page fault was in fetching operands for instruction
 - OS must figure out which...
- ◆ OS checks validity of address
- ◆ Process killed if address was illegal
- ◆ OS finds a place to put new page frame
- ◆ If frame selected for replacement is dirty, write it out to disk
- ◆ OS requests the new page from disk
- ◆ Page tables updated
- ◆ Faulting instruction backed up so it can be restarted
- ◆ Faulting process scheduled
- ◆ Registers restored
- ◆ Program continues

Backing up an instruction

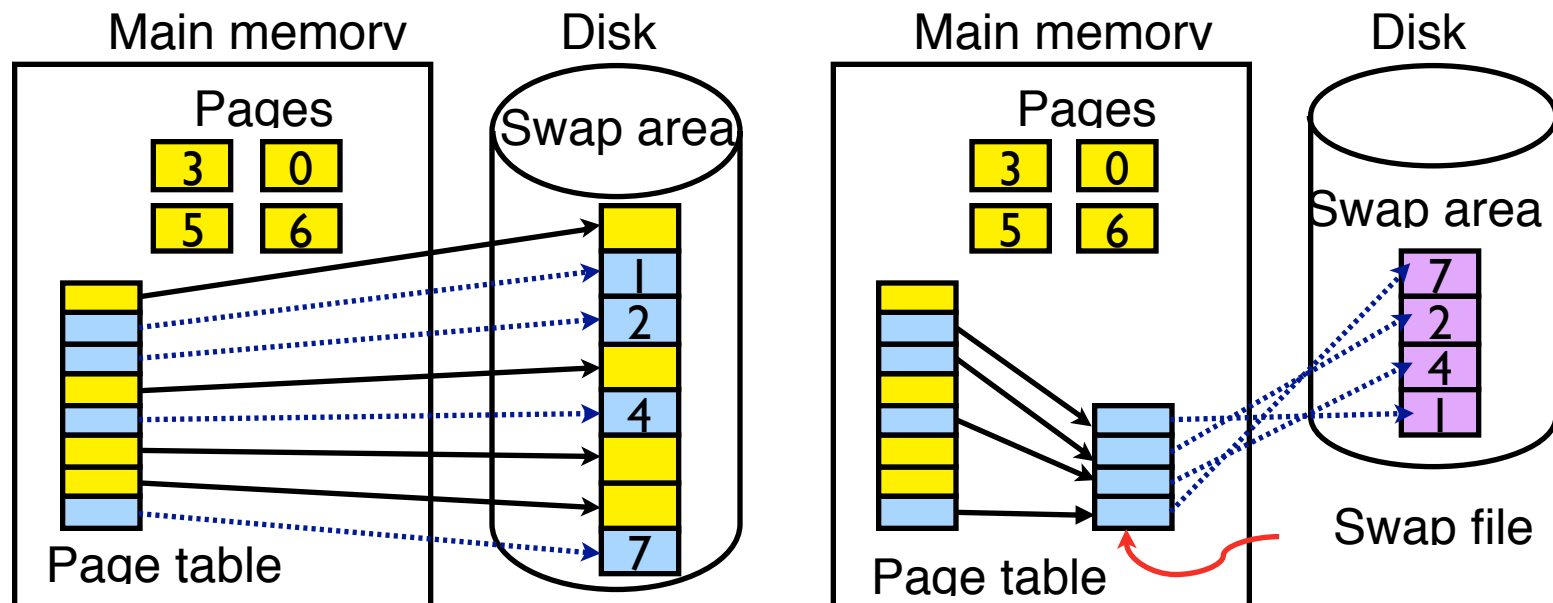
- ✦ Problem: page fault happens in the middle of instruction execution
 - Some changes may have already happened
 - Others may be waiting for VM to be fixed
- ✦ Solution: undo all of the changes made by the instruction
 - Restart instruction from the beginning
 - This is easier on some architectures than others
- ✦ Example: LW R1, 12(R2)
 - Page fault in fetching instruction: nothing to undo
 - Page fault in getting value at 12(R2): restart instruction
- ✦ Example: ADD (Rd)+,(Rs1)+,(Rs2)+
 - Page fault in writing to (Rd): may have to undo an awful

Locking pages in memory

- ✦ Virtual memory and I/O occasionally interact
- ✦ P1 issues call for read from device into buffer
 - While it's waiting for I/O, P2 runs
 - P2 has a page fault
 - P1's I/O buffer might be chosen to be paged out
 - This can create a problem because an I/O device is going to write to the buffer on P1's behalf
- ✦ Solution: allow some pages to be locked into memory
 - Locked pages are immune from being replaced
 - Pages only stay locked for (relatively) short periods

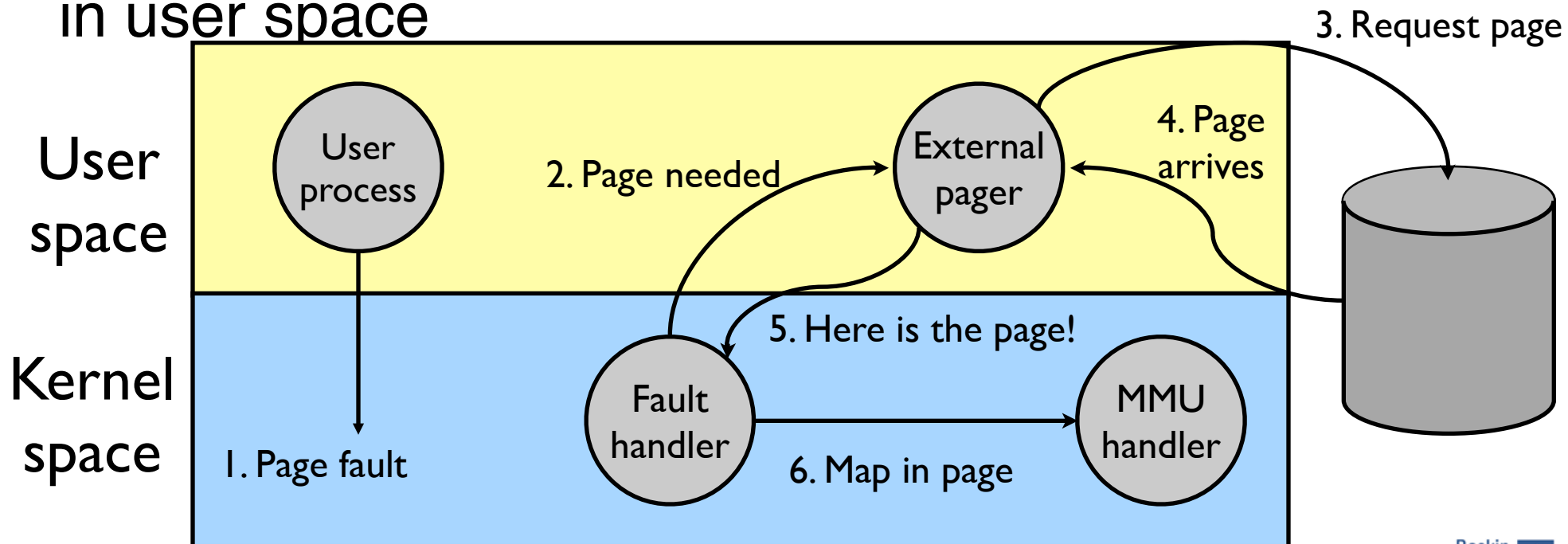
Storing pages on disk

- ◆ Pages removed from memory are stored on disk
- ◆ Where are they placed?
 - Static swap area: easier to code, less flexible
 - Dynamically allocated space: more flexible, harder to locate a page
 - Dynamic placement often uses a special file (managed by the file system) to hold pages
- ◆ Need to keep track of which pages are where within the on-disk storage



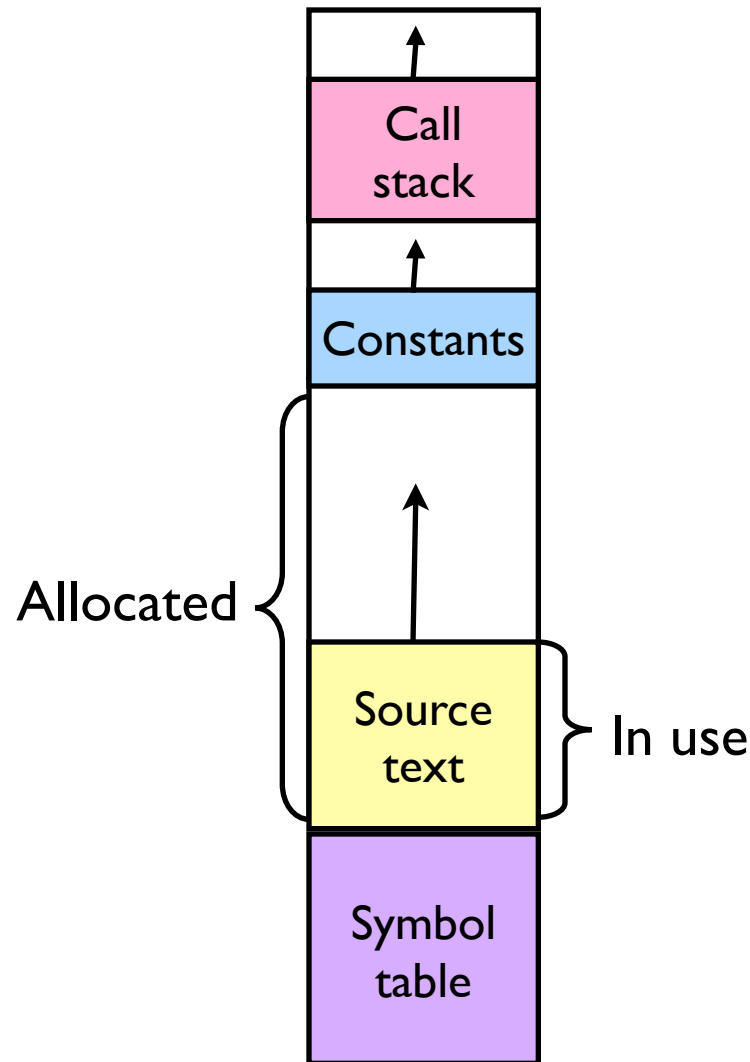
Separating policy and mechanism

- ✦ Mechanism for page replacement has to be in kernel
 - Modifying page tables
 - Reading and writing page table entries
- ✦ Policy for deciding which pages to replace could be in user space



Why use segmentation?

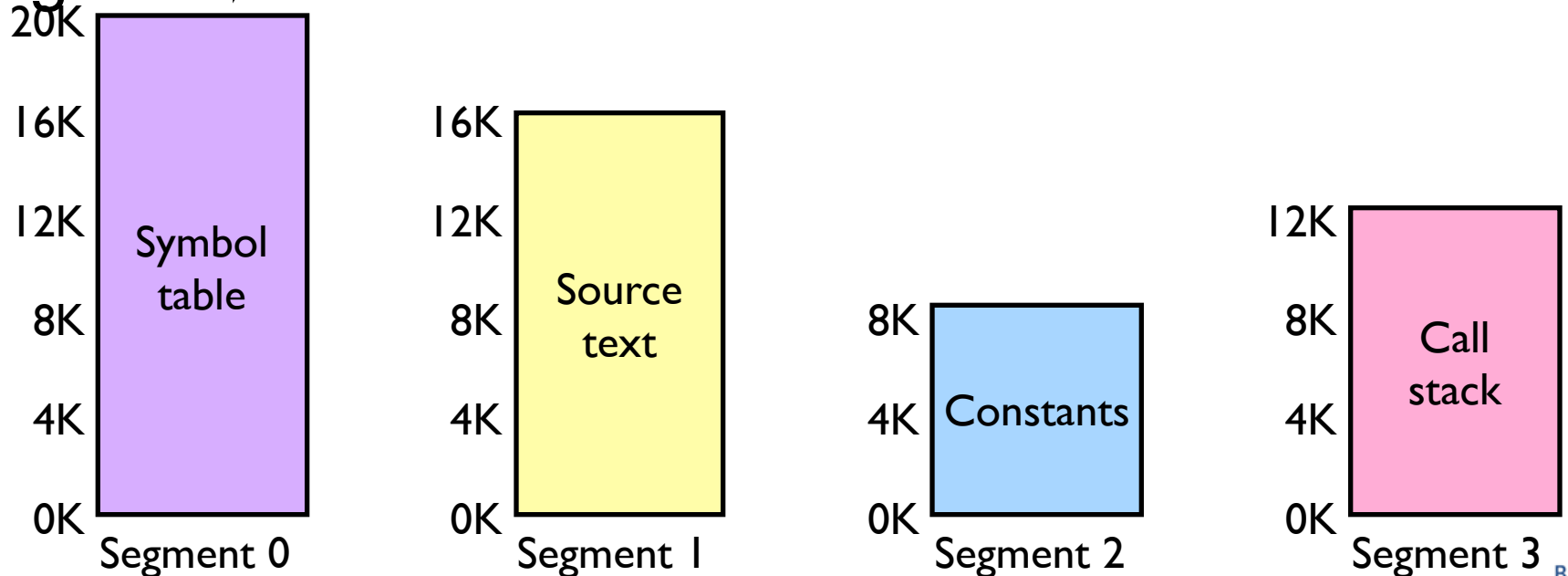
Virtual address space



- ✦ Different “units” in a single virtual address space
 - Each unit can grow
 - How can they be kept apart?
 - Example: symbol table is out of space
- ✦ Solution: segmentation
 - Give each unit its own address space

Using segments

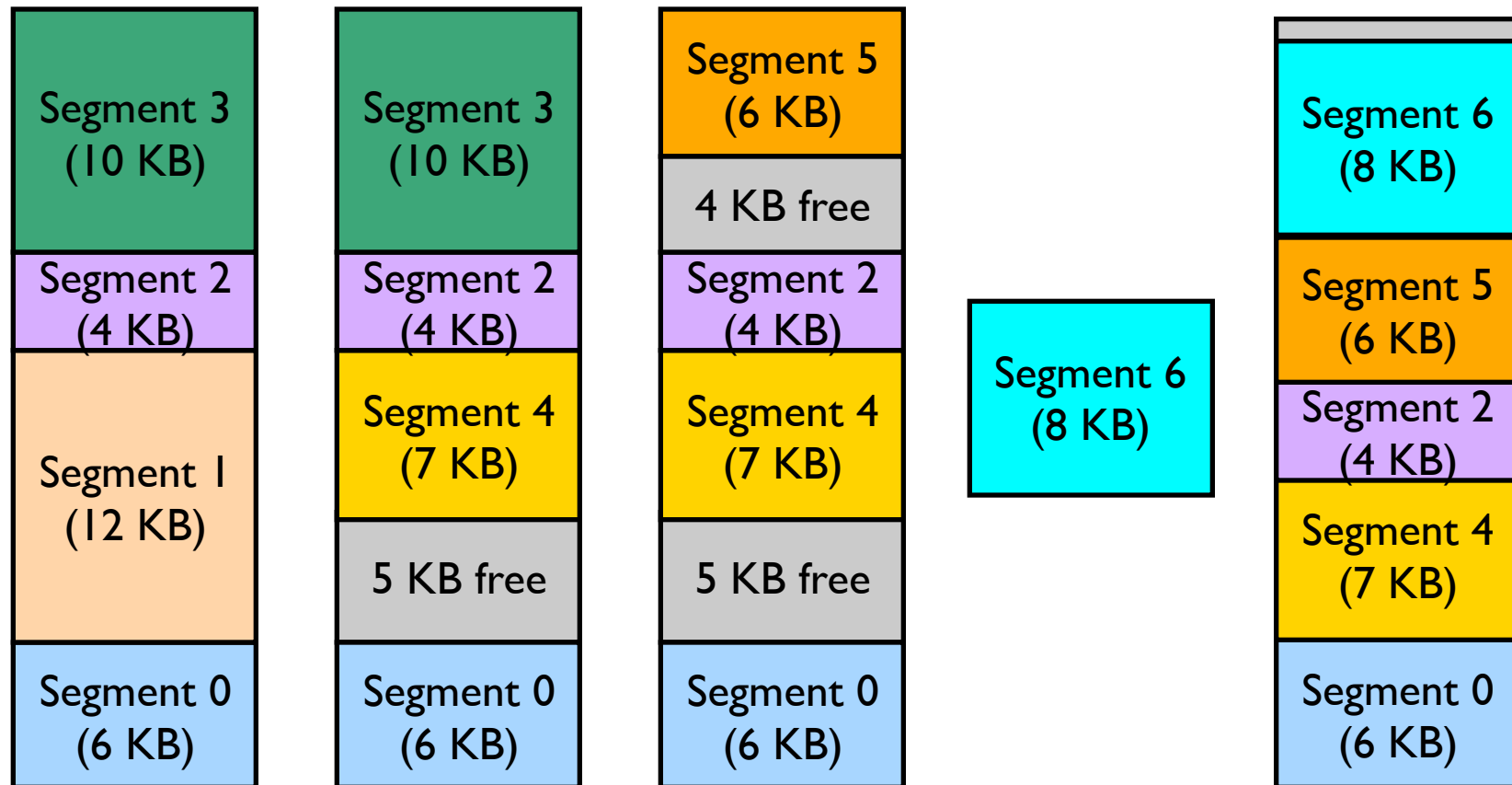
- ◆ Each region of the process has its own segment
- ◆ Each segment can start at 0
 - Addresses within the segment relative to the segment start
- ◆ Virtual addresses are $\langle \text{segment \#}, \text{offset within segment} \rangle$



Paging vs. segmentation

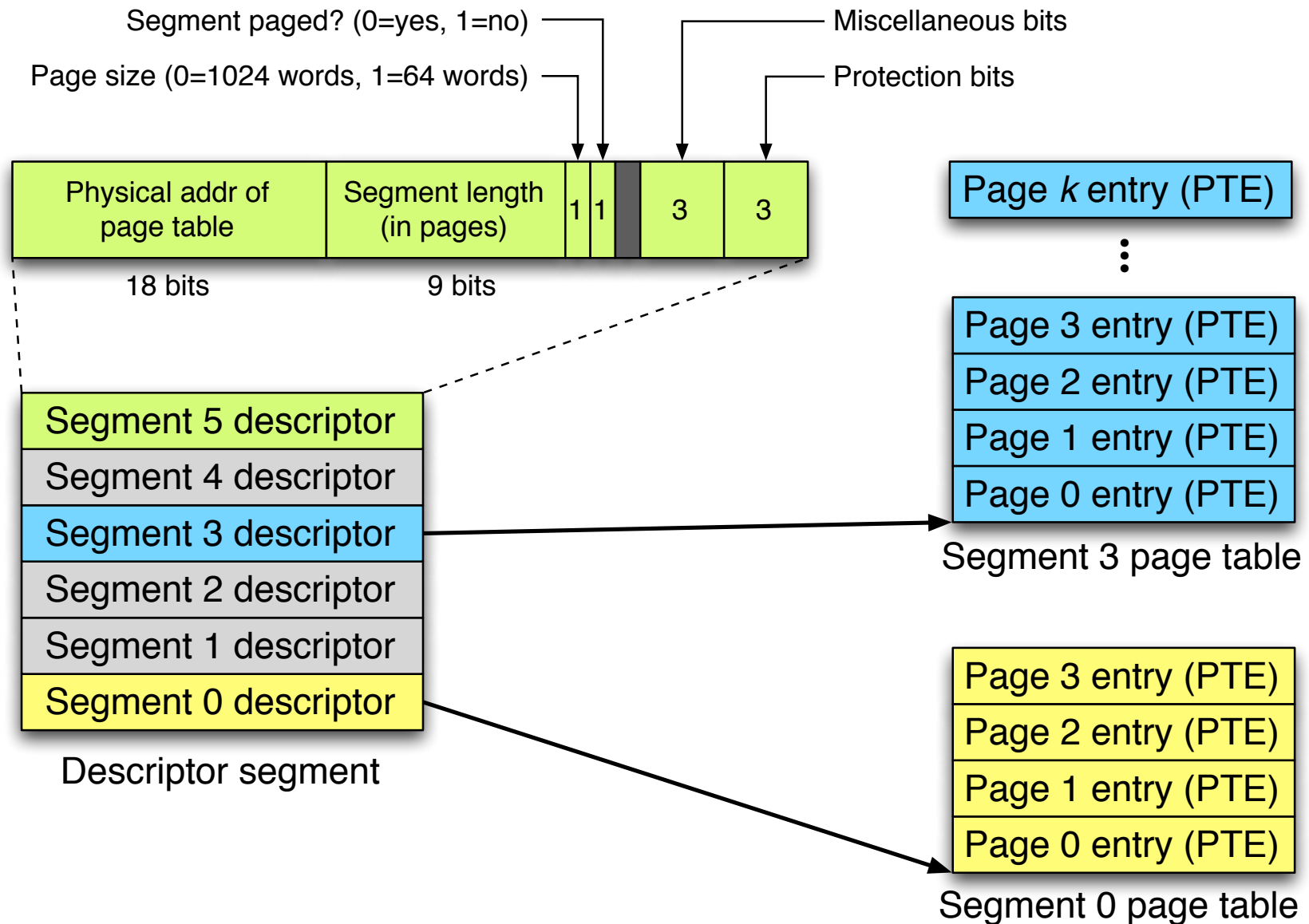
What?	Paging	Segmentation
Need the programmer know about it?	No	Yes
How many linear address spaces?	One	Many
More addresses than physical memory?	Yes	Yes
Separate protection for different objects?	Not really	Yes
Variable-sized objects handled with ease?	No	Yes
Is sharing easy?	No	Yes
Why use it?	More address space without buying more memory	Break programs into logical pieces that are handled separately.

Implementing segmentation

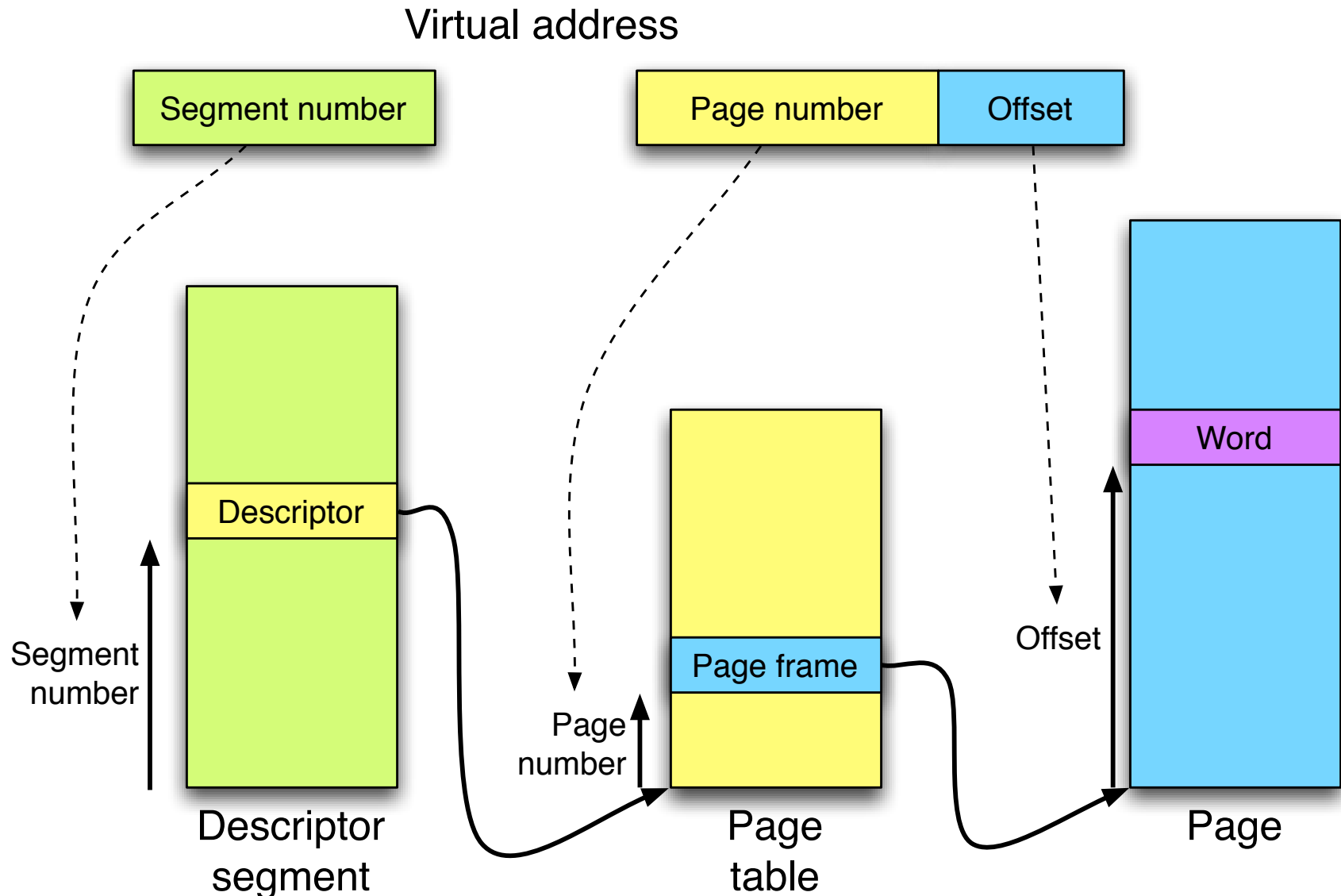


→ Need to do memory compaction!

A better way: segmentation *with* paging

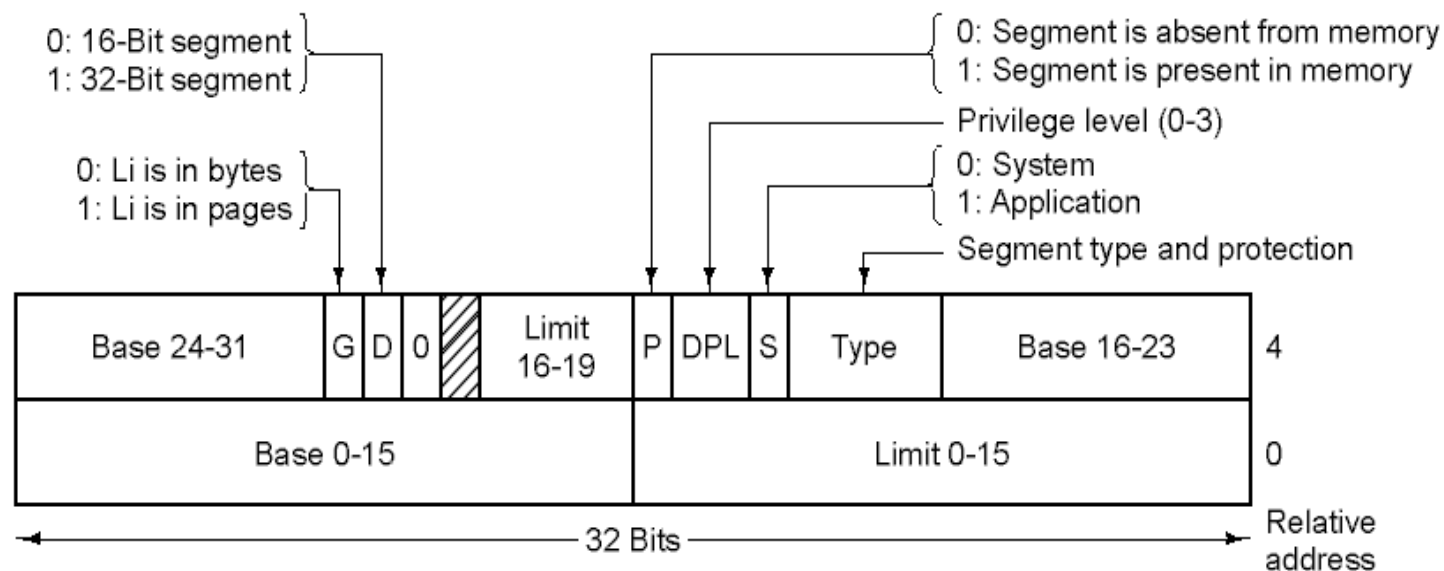


Translating an address in MULTICS



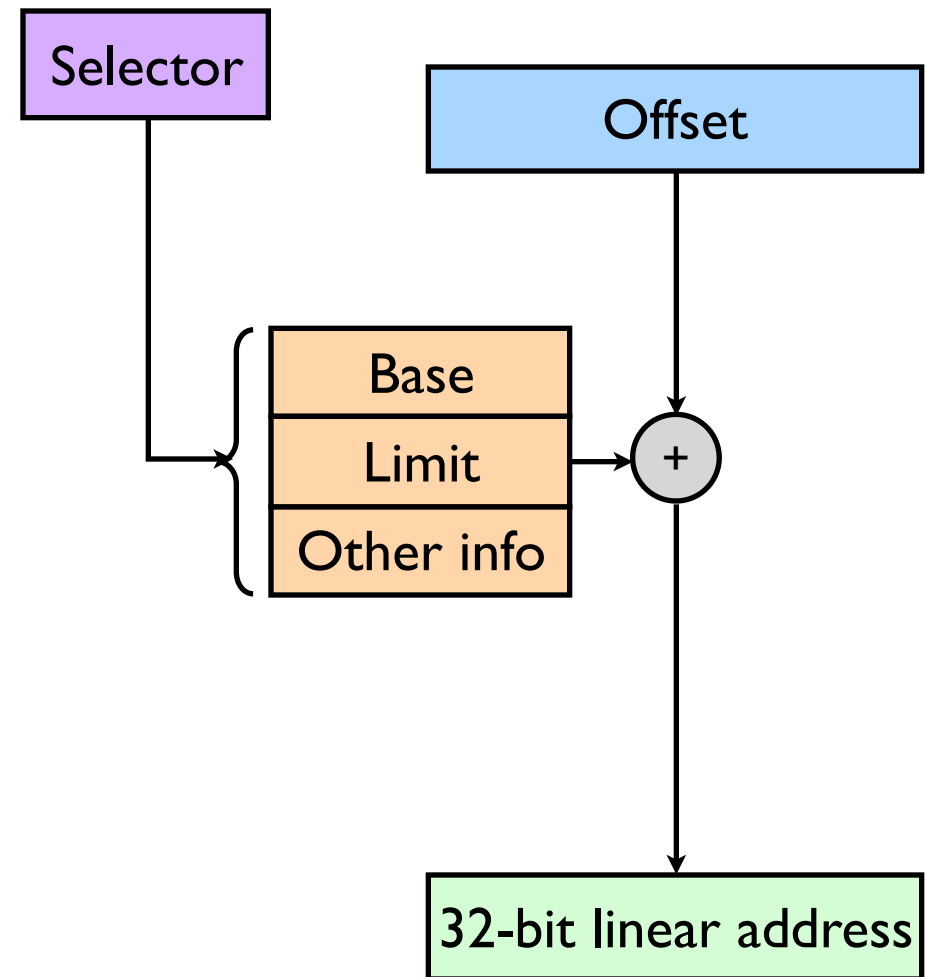
Memory management in the Pentium

- ◆ Memory composed of segments
 - Segment pointed to by segment descriptor
 - Segment selector used to identify descriptor
- ◆ Segment descriptor describes segment
 - Base virtual address
 - Size
 - Protection
 - Code / data



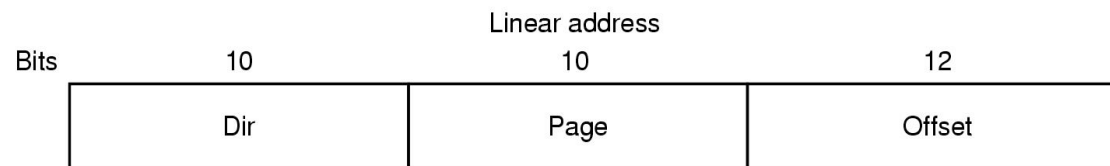
Converting segment to linear address

- ✦ Selector identifies segment descriptor
 - Limited number of selectors available in the CPU
- ✦ Offset added to segment's base address
- ✦ Result is a virtual address that will be translated by paging

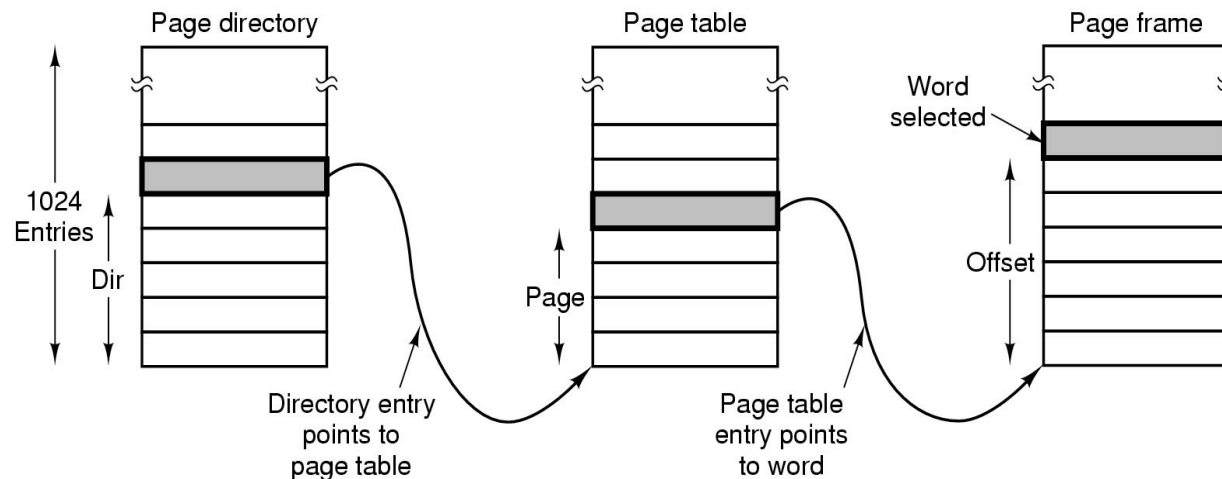


Translating virtual to physical addresses

- ◆ Pentium uses two-level page tables
 - Top level is called a “page directory” (1024 entries)
 - Second level is called a “page table” (1024 entries each)
 - 4 KB pages



(a)



(b)