Copy-on-write & Fork

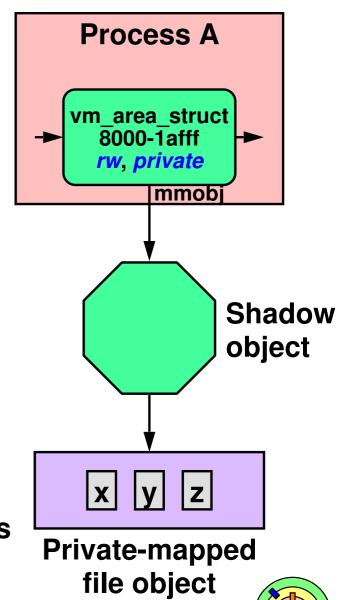


Shadow Objects

- indirection
- keep track of pages that were originally copy-on-write but have been modified
- A page in a memory map, into which an object was mapped *private* (e.g., data region), has an associated *shadow object*
 - if a page is "managed by a shadow object" (or "referenced in a shadow object"), it has been modified
 - otherwise, the page is managed by the *original* object (file or a "zero/anonymous" object)
 - x, y, z on the right are pages / page frames



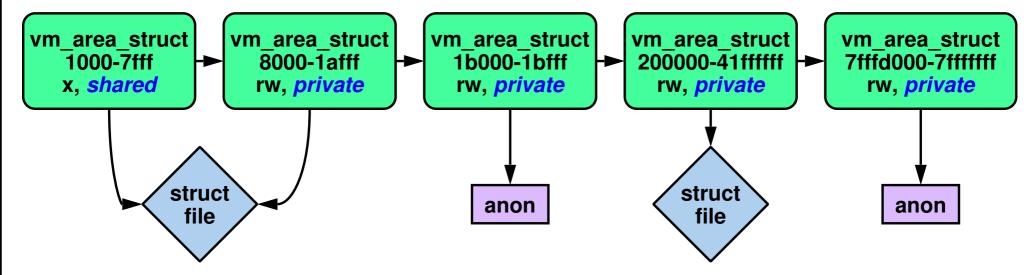
Shadow object tells you where to copy from when you need to perform copy-on-write



Address-Space Representation



Remember this?

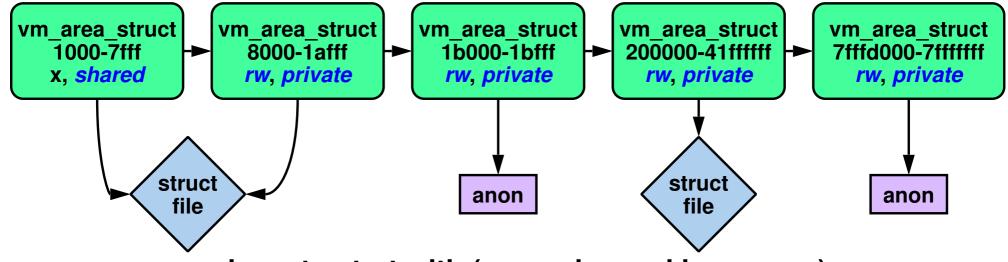




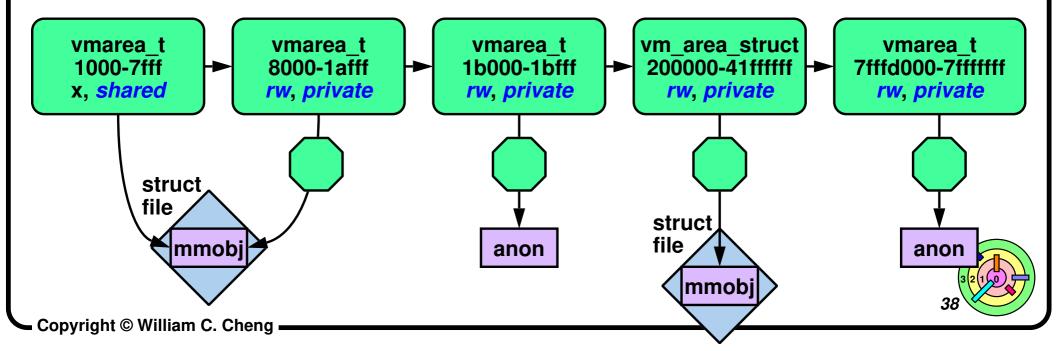
Address-Space Representation



Remember this?



now we have to start with (mmobj is used in weenix):



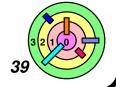
Share Mapping (1)

Process A shared **Process A has** shared-mapped a file into its address space |y| |z|

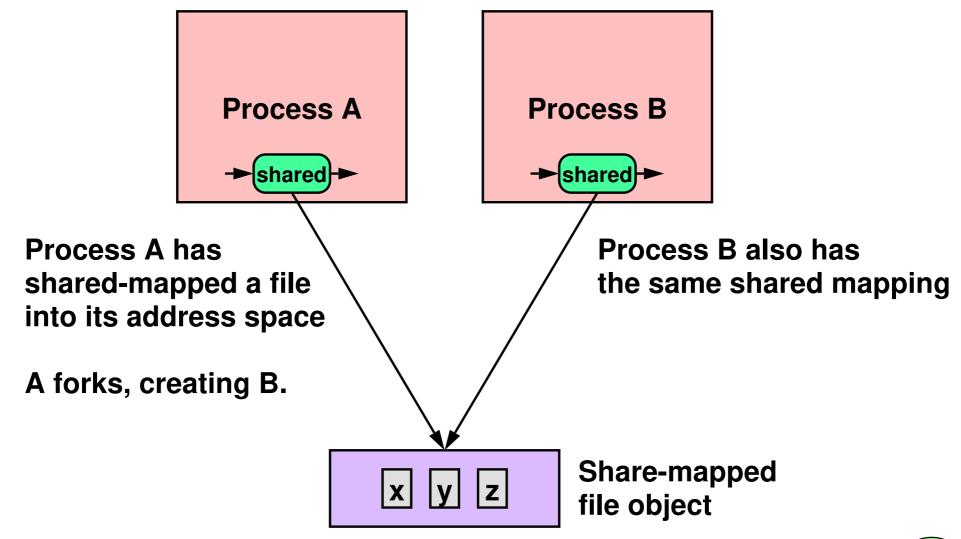
in weenix

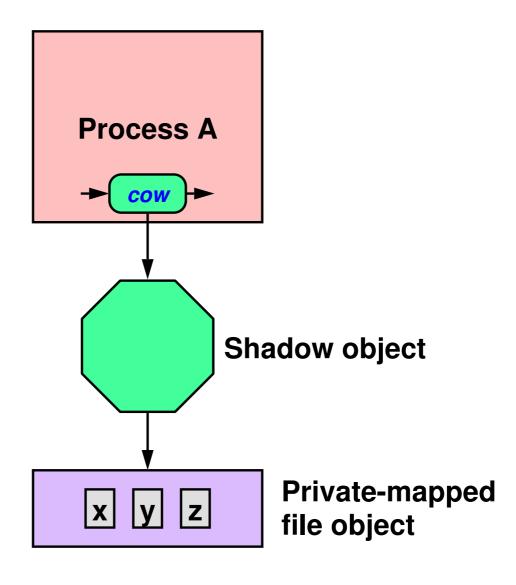
- instead of pointing to a File object, it's pointing to an mmobj inside a vnode inside a File object
- mmobj is used to manage page frames

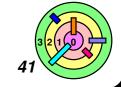
File object

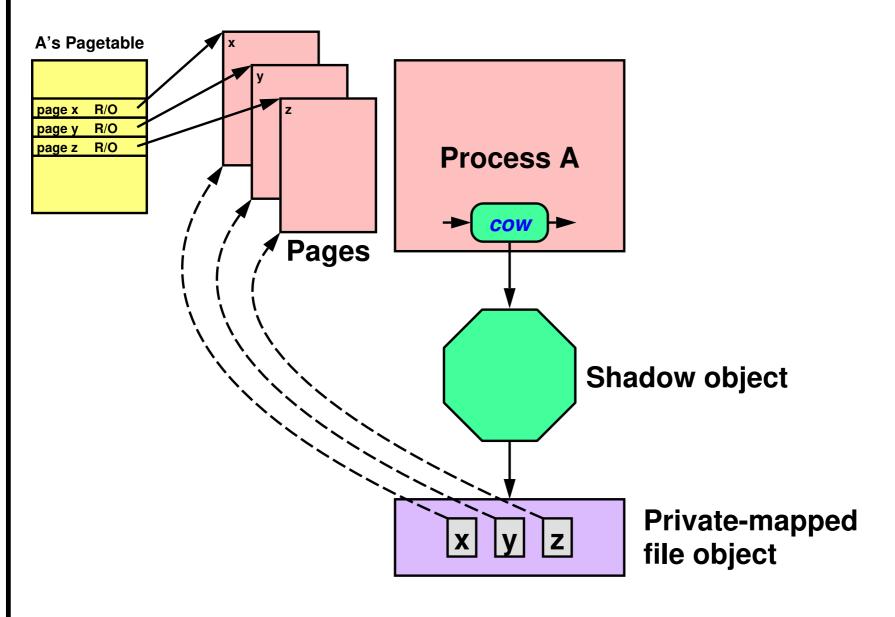


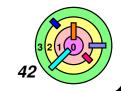
Share Mapping (2)

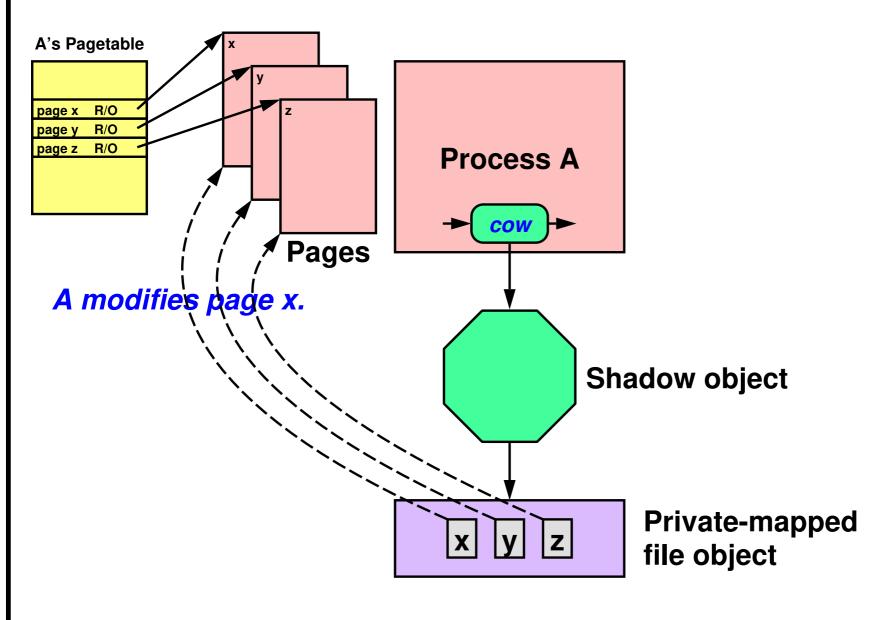


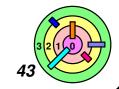


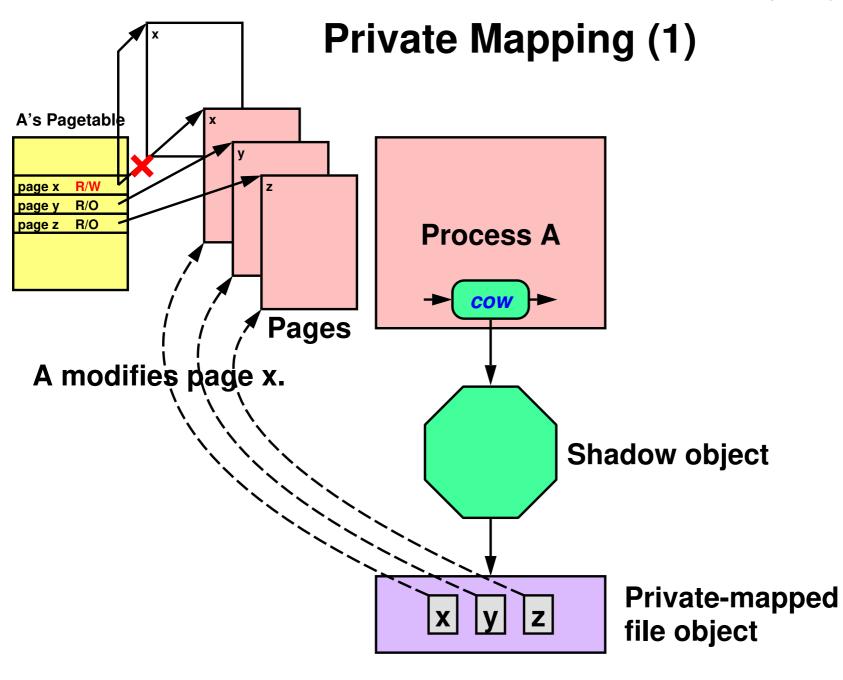




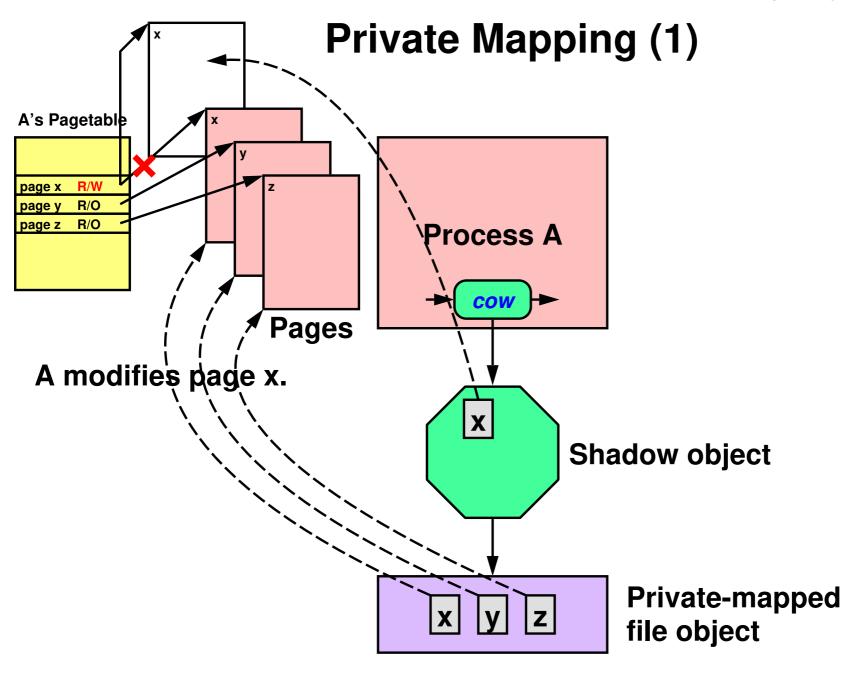


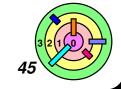


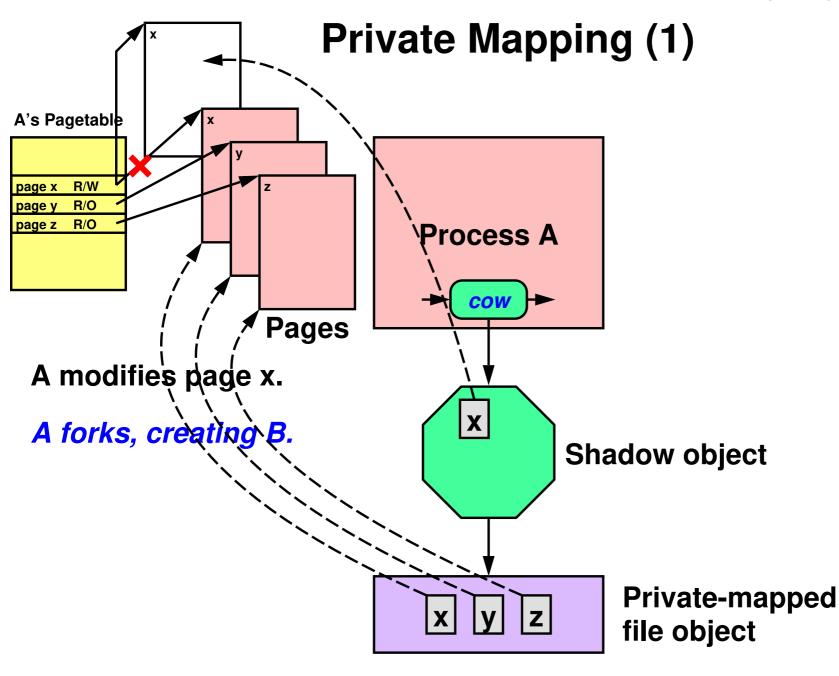


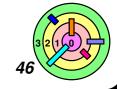


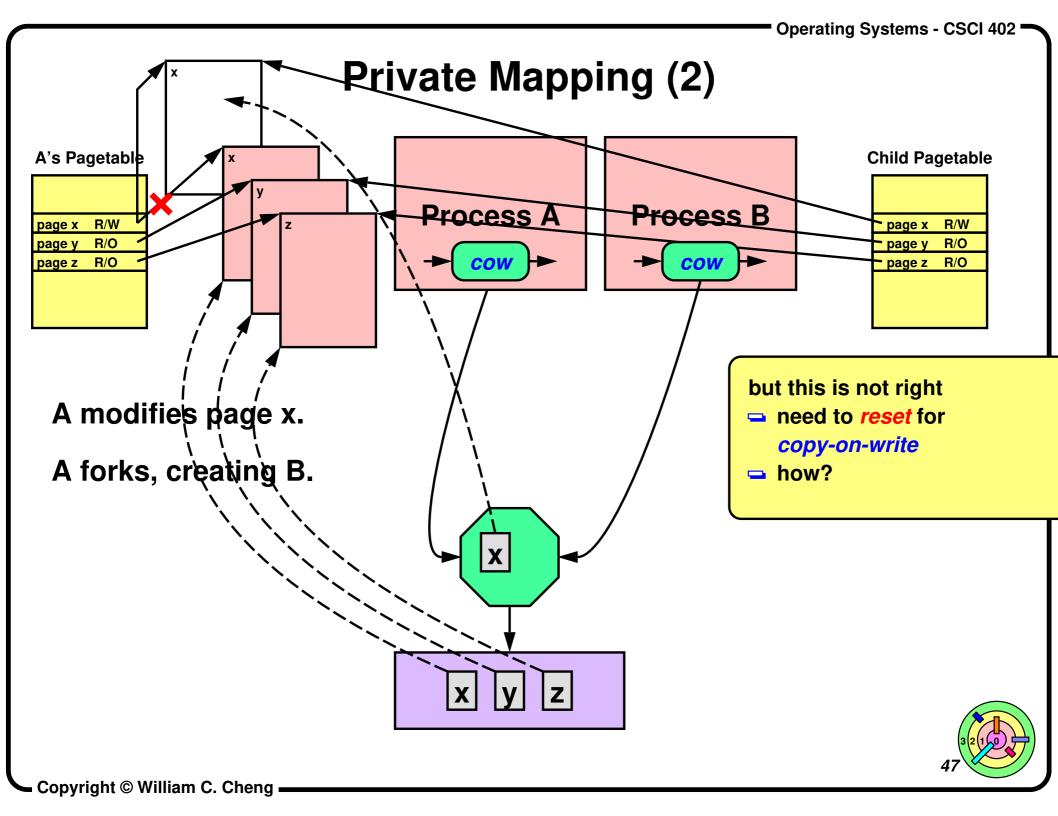


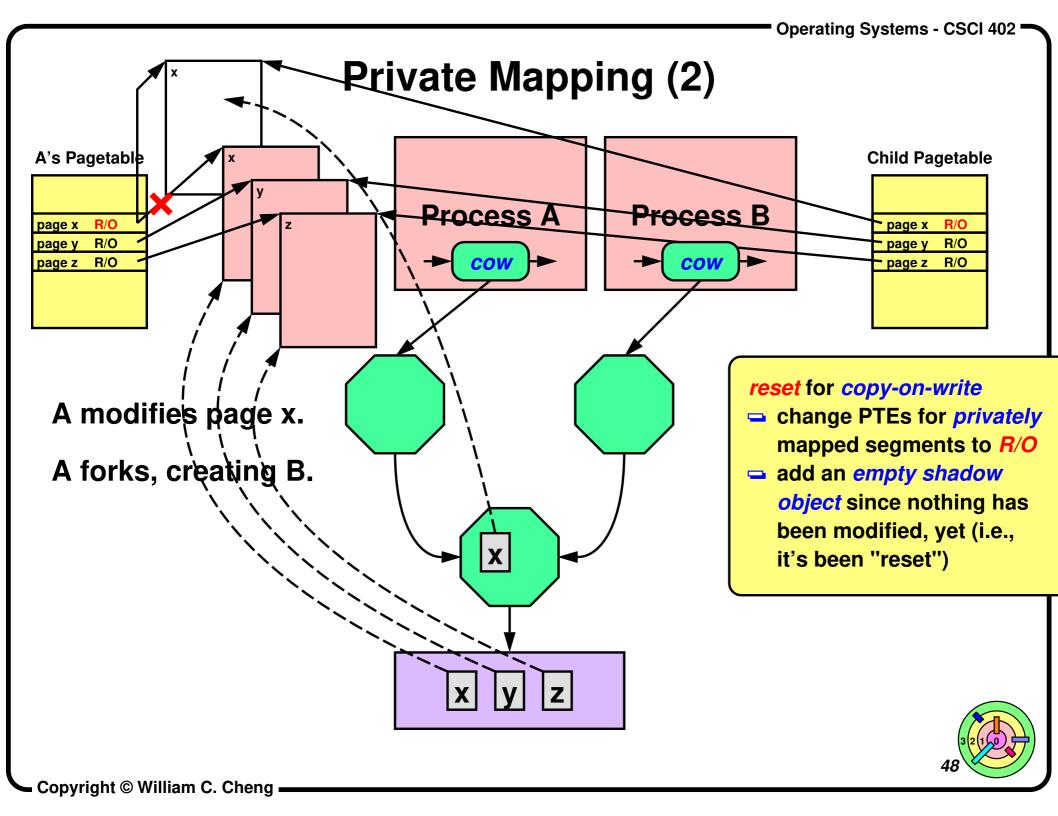


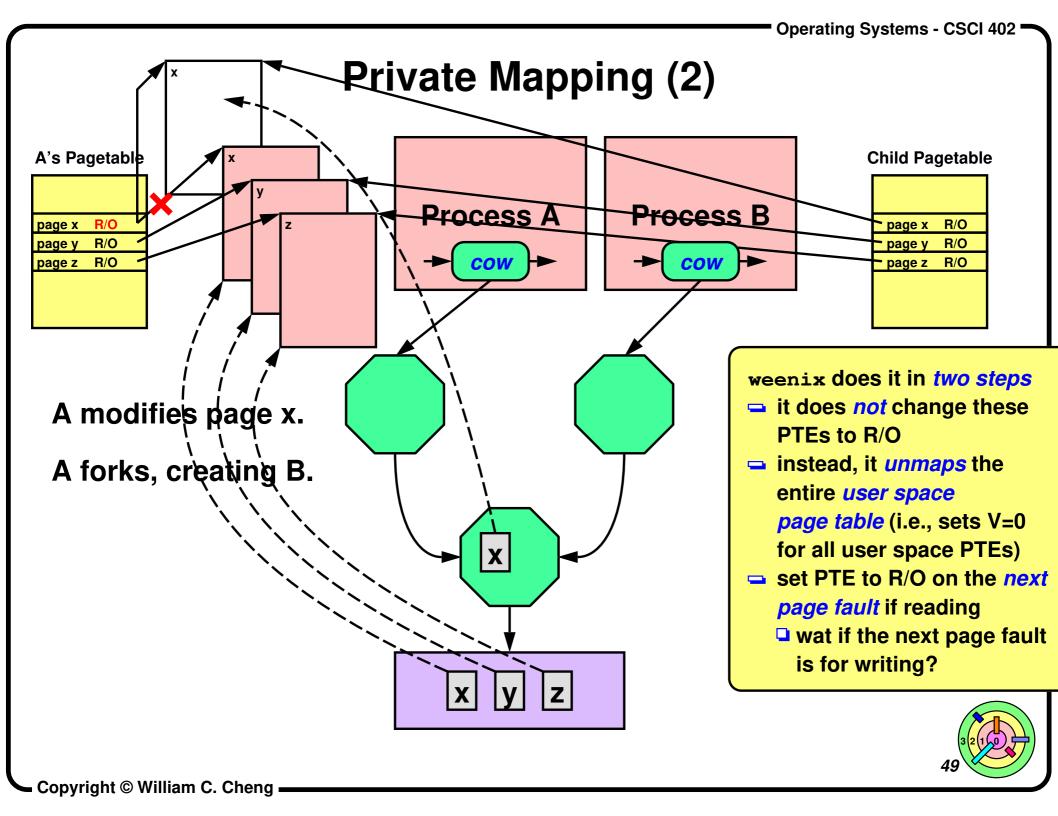


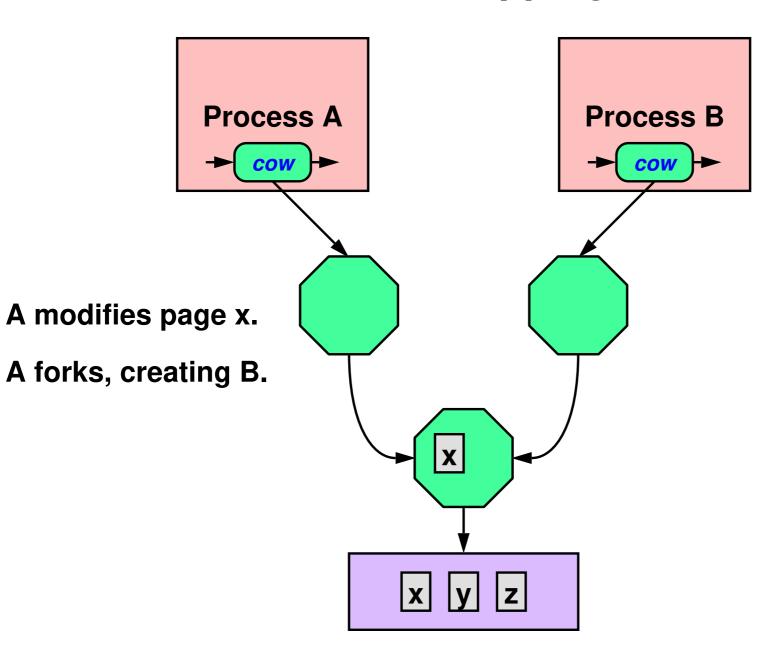


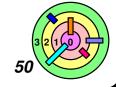


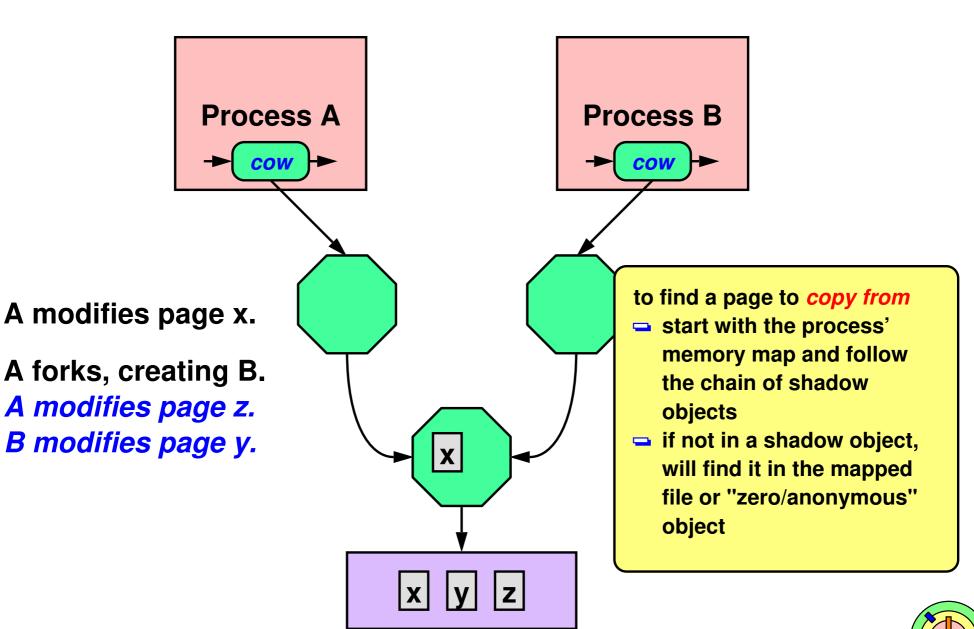


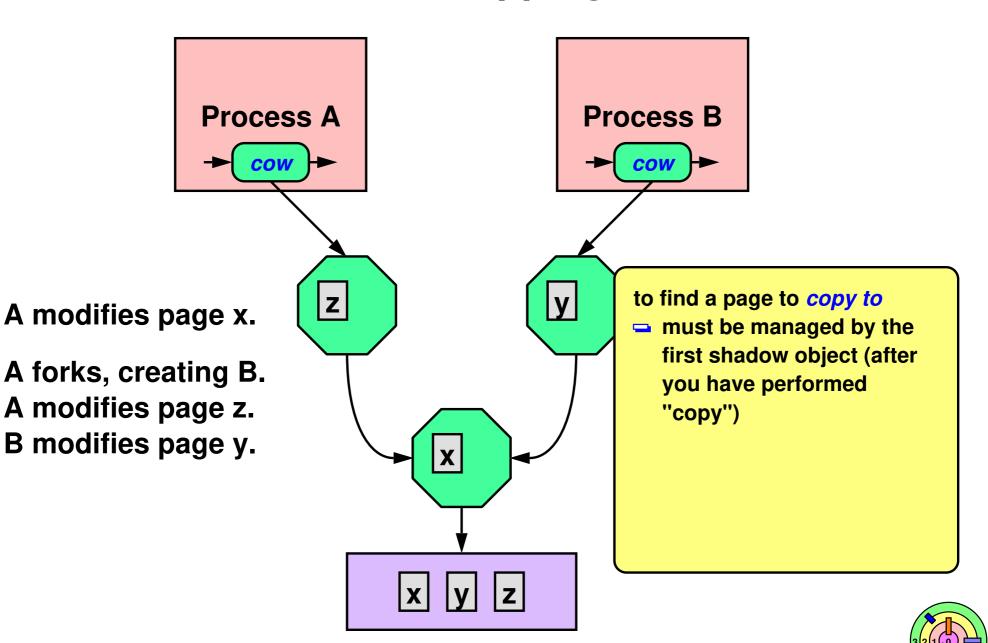


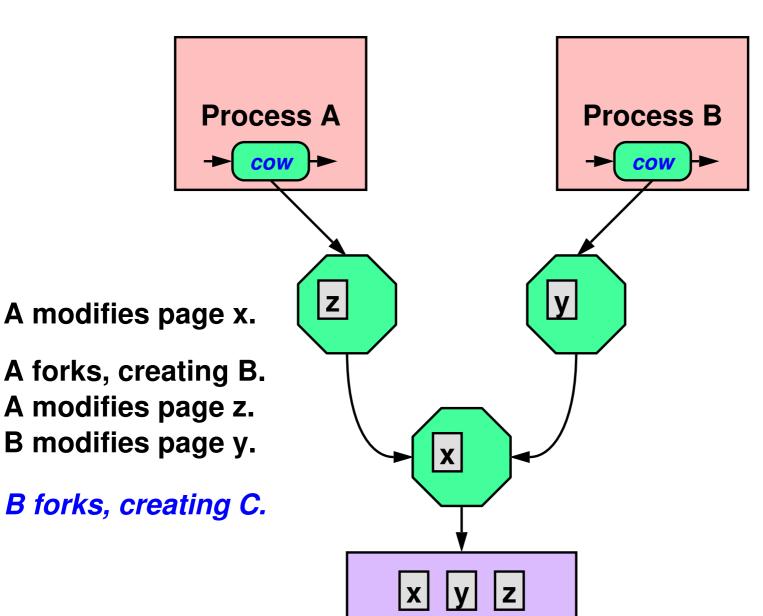


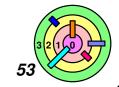


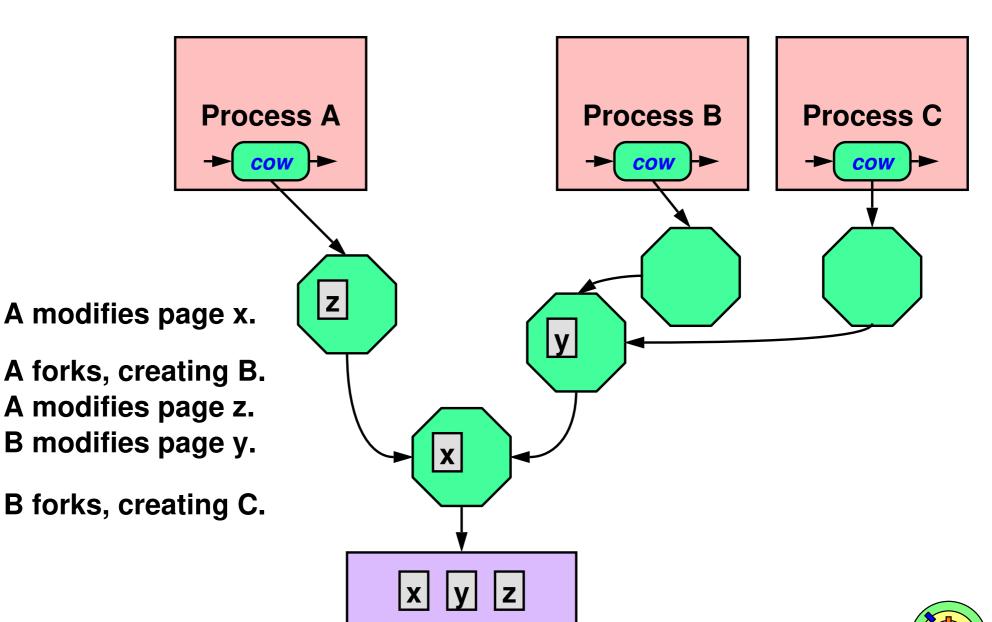


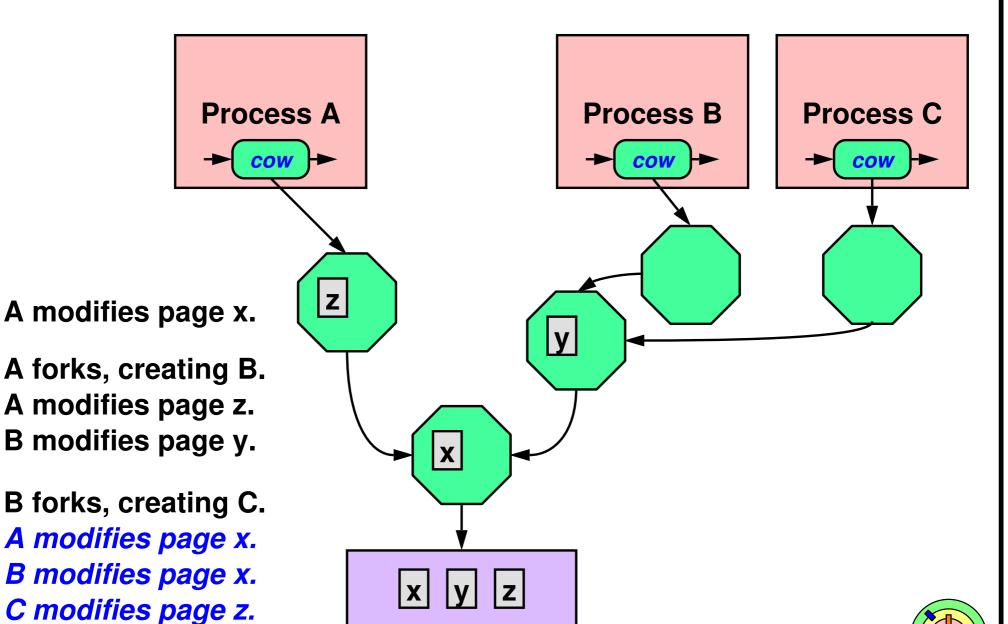




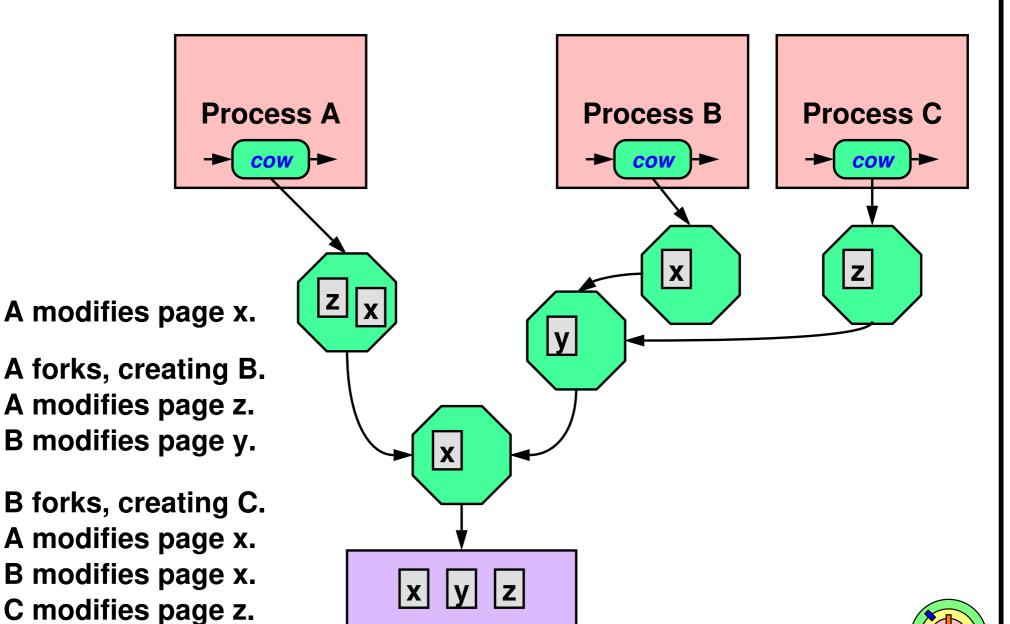




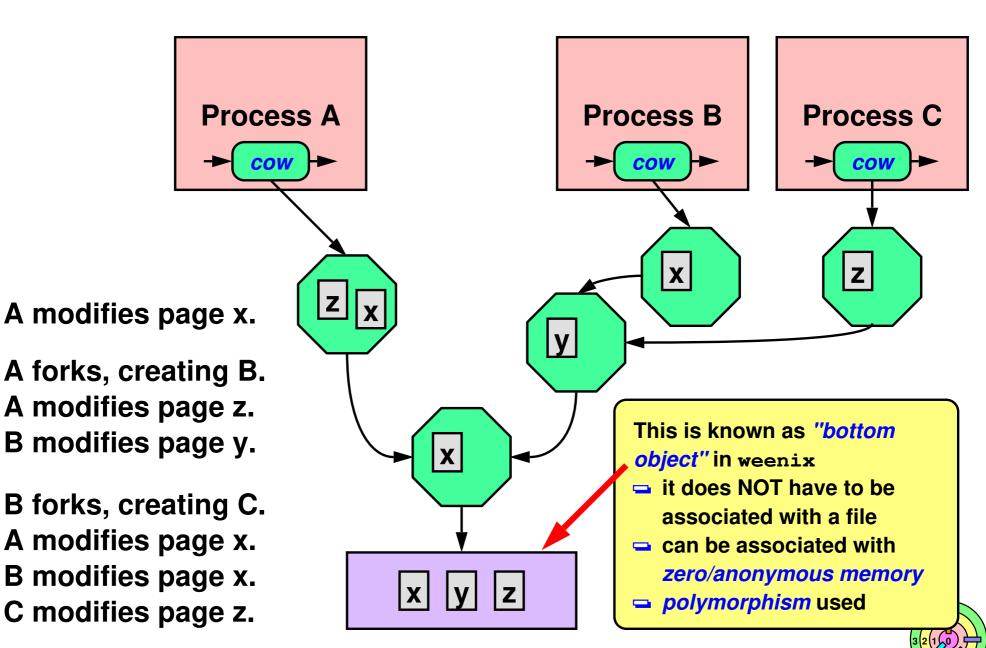


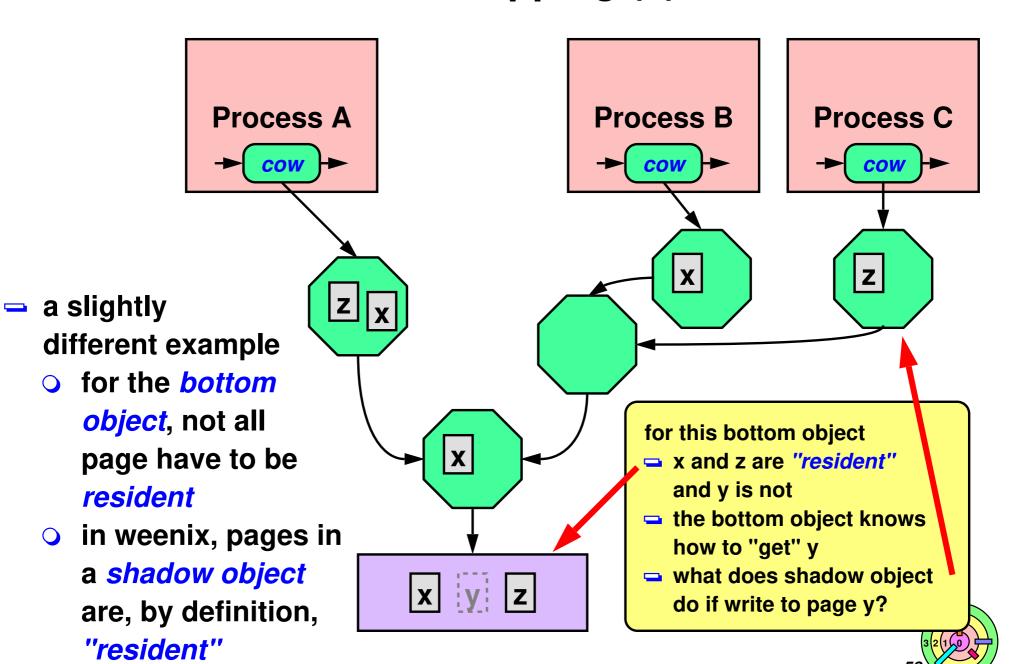


Copyright © William C. Cheng



Copyright © William C. Cheng





Virtual Copy

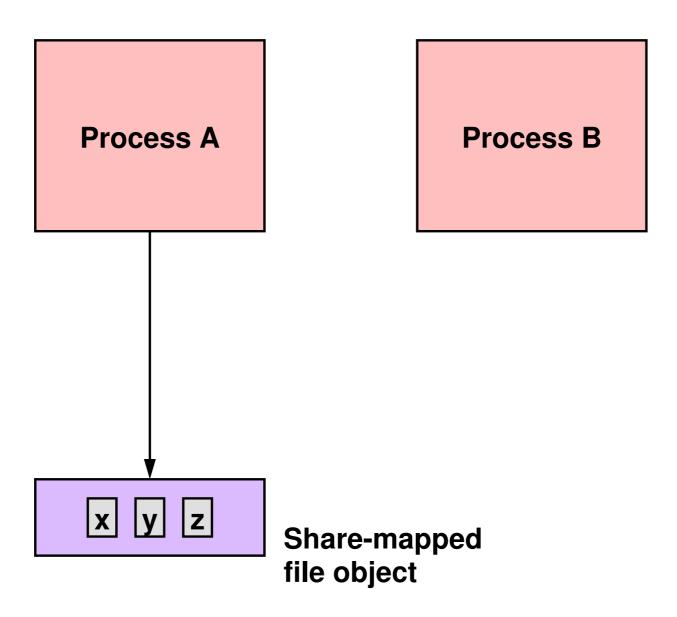


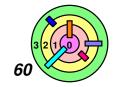
Local RPC

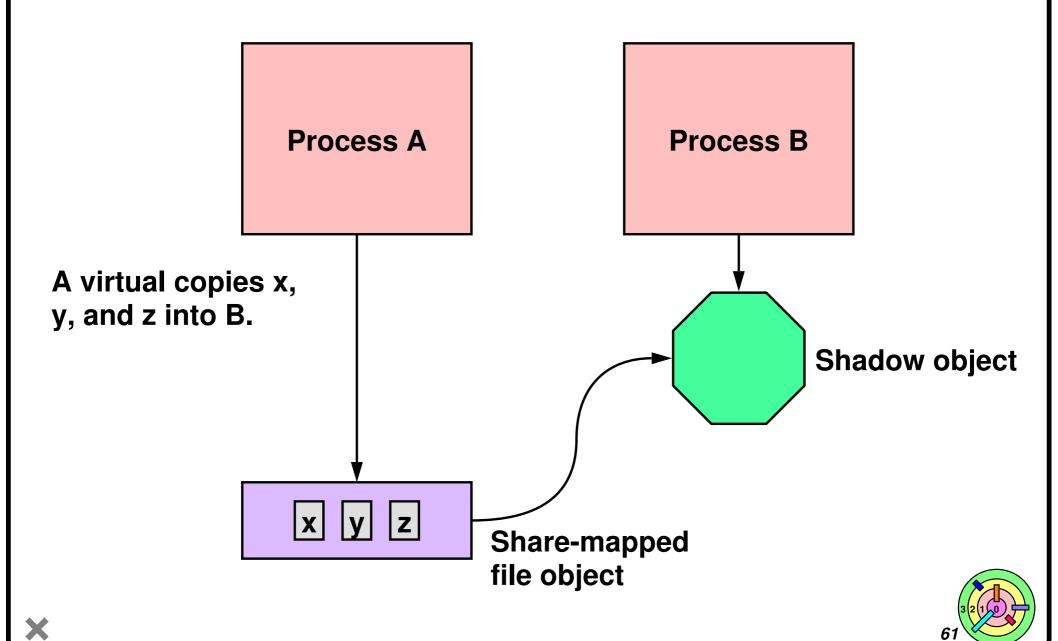
- "copy" arguments from one process to another
- assume arguments are page-aligned and page-sized
- map pages into both caller and callee, copy-on-write
 - works in most cases, except when the page corresponds to a shared memory-mapped file
 - in this case, the sender does not have a shadow object!

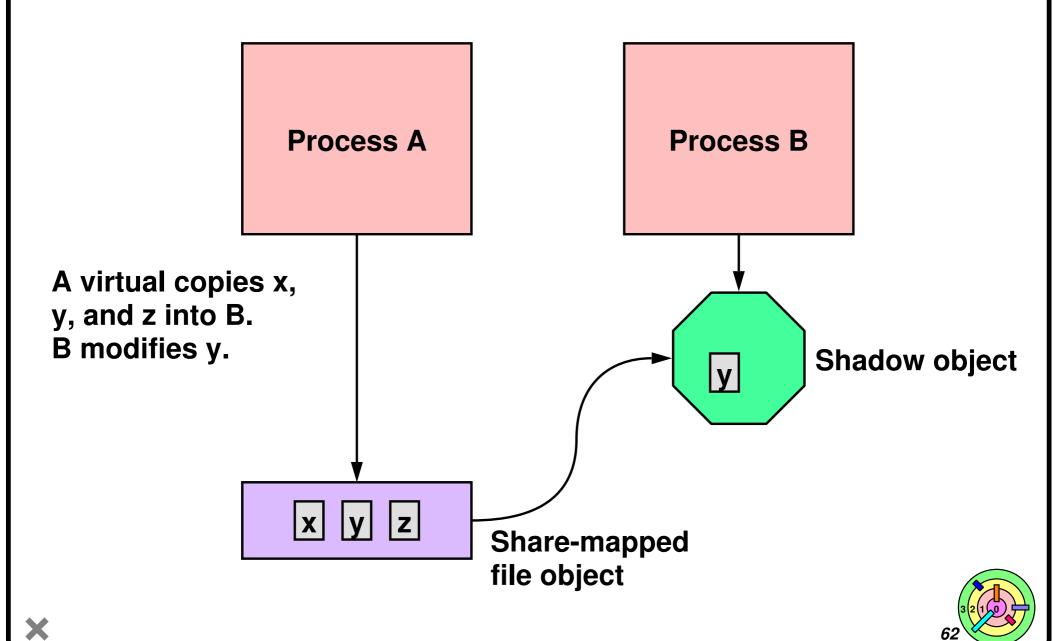


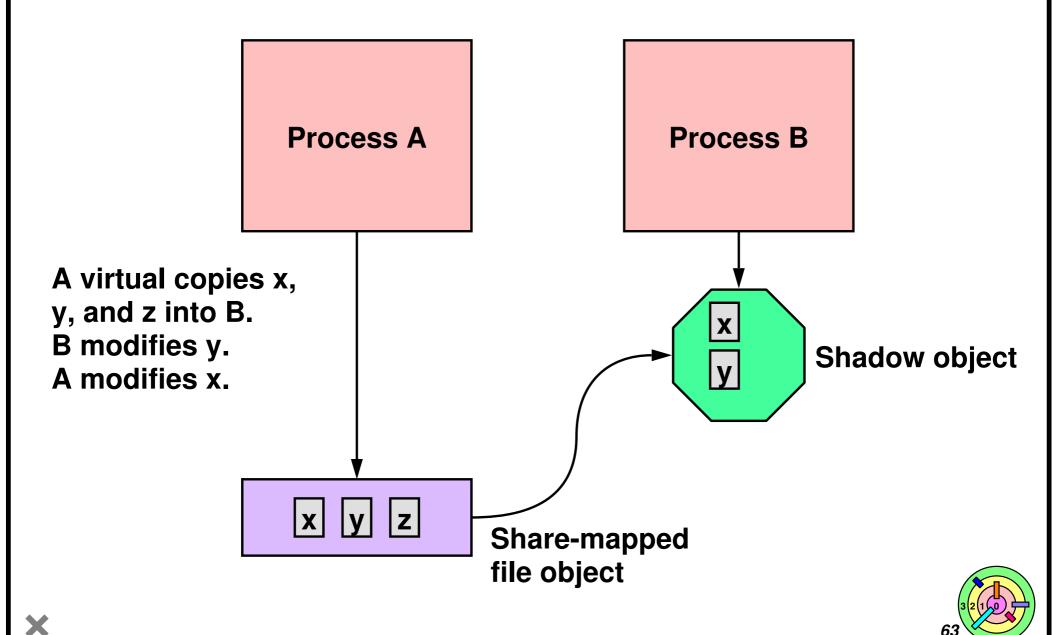




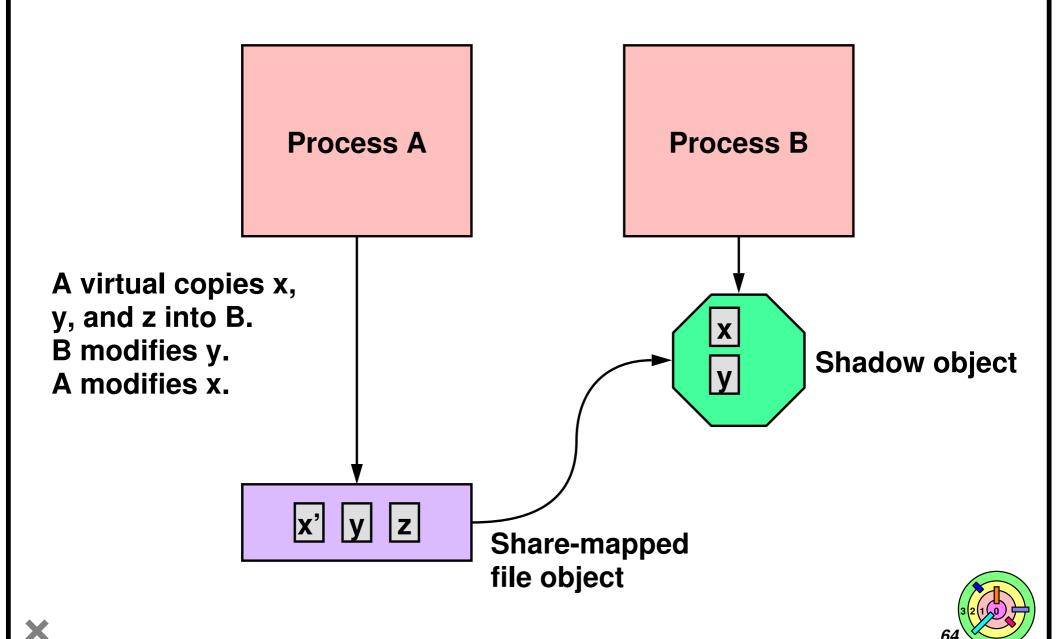








Copyright © William C. Cheng



Copyright © William C. Cheng

Shadow Objects Summary



Why go through all this trouble?

- because we want to implement copy-on-write together with fork ()
 - a variable (such as Data a few slides back) can exist in many different physical pages simultaneously
 - each contains a different version of this variable

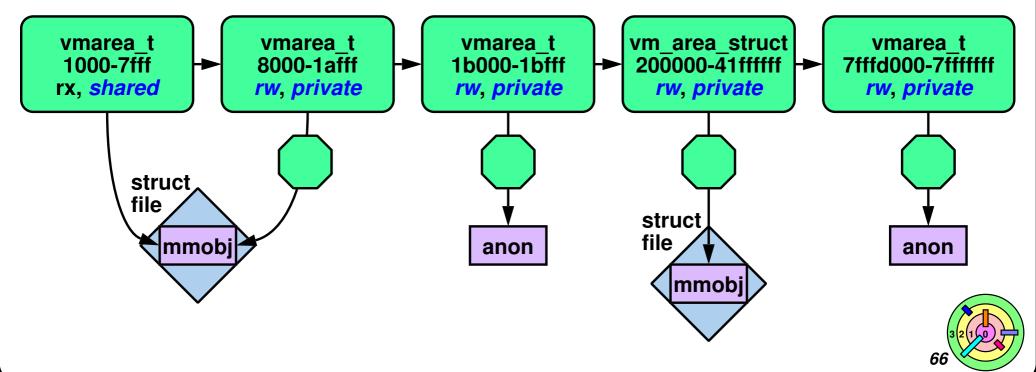


To manage this mess, weenix uses the idea of Shadow Objects

- what is the "idea" of Shadow Objects?
 - organize a tree of shadow objects using an inverted tree data structure
 - where the root is the bottom object
 - the rule for finding page frame / physical page that contains the global variable in question for a particular process
 - traversing shadow object pointers on the inverted tree
 - when and how to perform copy-on-write
- you have to implement what's described on these slides in kernel 3

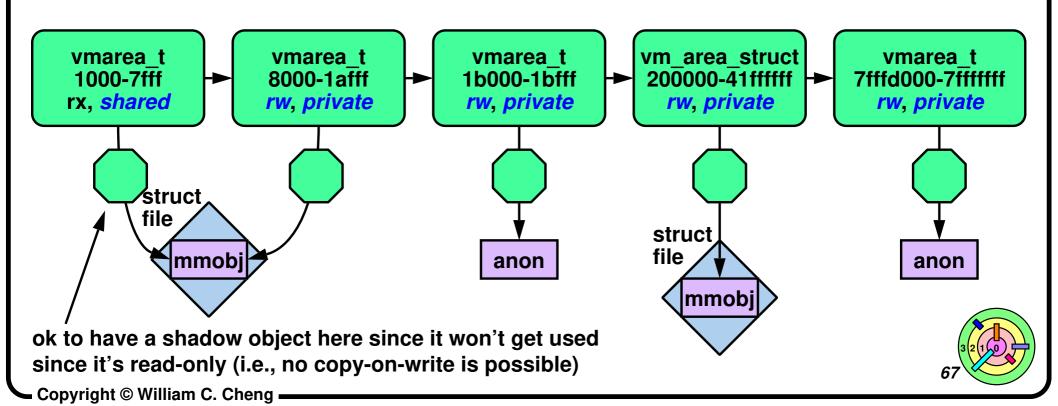


- types of mmobj in kernel assignments are:
 - there's one that lives inside a vnode (vn->vn_mmobj)
 - a shadow object is an mmobj
 - an anonymous object (meaning not associated with a file and not a shadow object) is an mmobj
- a vmarea is supported by one of these 3 mmobjs

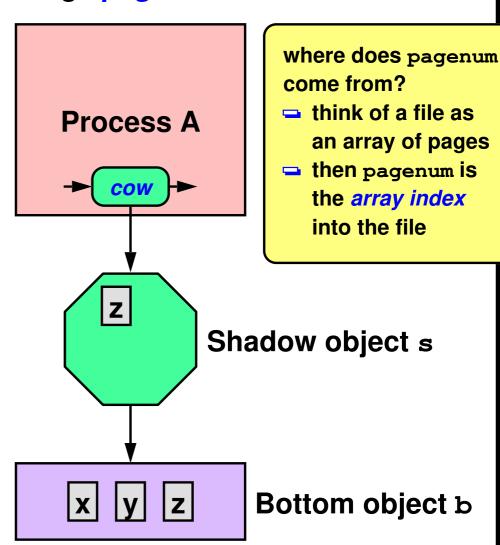


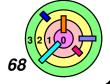


- types of mmobj in kernel assignments are:
 - there's one that lives inside a vnode (vn->vn_mmobj)
 - a shadow object is an mmobj
 - an anonymous object (meaning not associated with a file and not a shadow object) is an mmobj
- a vmarea is supported by one of these 3 mmobjs



- a page frame is uniquely identified by an mmobj and a pagenum
 - o notation:
 (mmobj, pagenum)

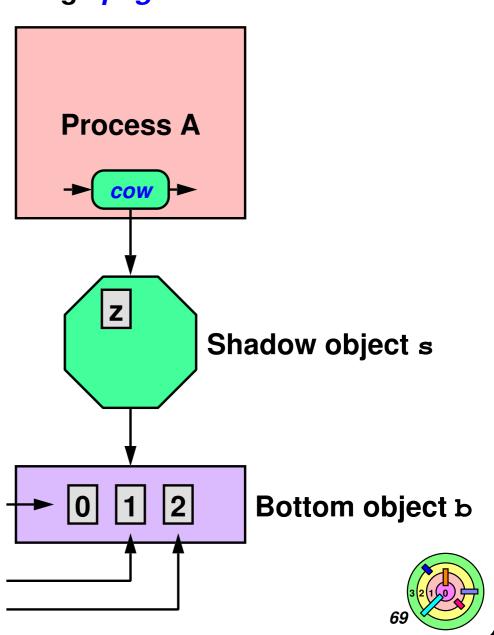




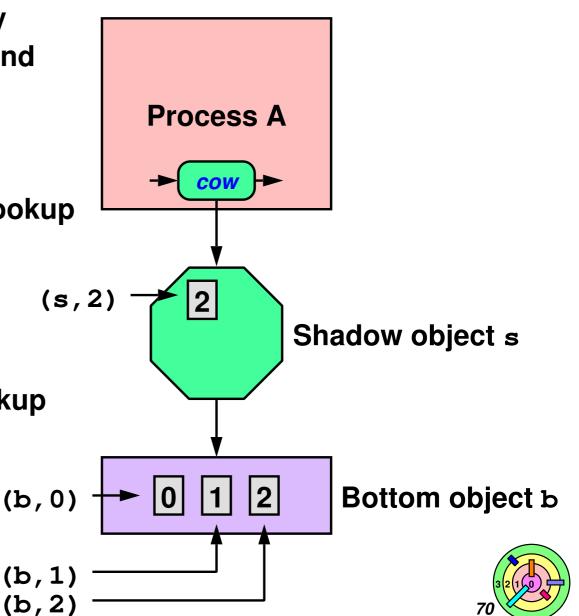
(b, 0)

(b, 1)

- a page frame is uniquely identified by an mmobj and a pagenum
 - o notation:
 (mmobj, pagenum)
 - if you map part of a file (say a page) into your address space, you need to remember which page
 - pagenum is then a page index in that file



- a page frame is uniquely identified by an mmobj and a pagenum
 - o notation:
 (mmobj, pagenum)
 - hash table used for lookup
 - read kernel 3 FAQ
- sometimes, you know the exact name of a page frame
 - use hash table to lookup
- sometimes, you only know pagenum (e.g., "where is page z?")
 - need to search

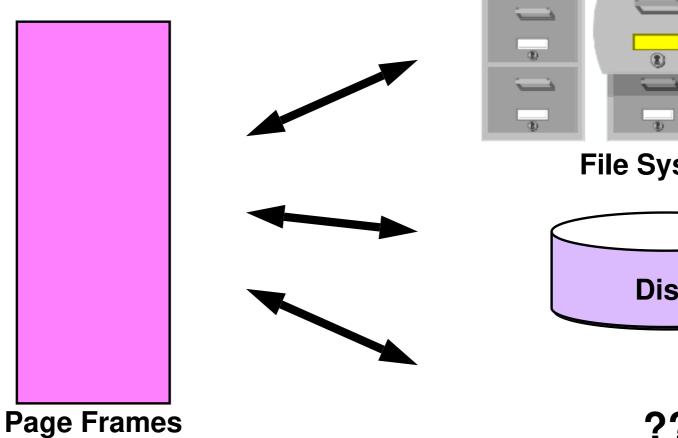


7.3 Operating System Issues

- General Concerns
- Representative Systems
- Copy on Write and Fork
- Backing Store Issues

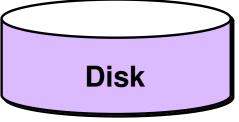


The Backing Store





File System



??



Backing Up Pages (1)



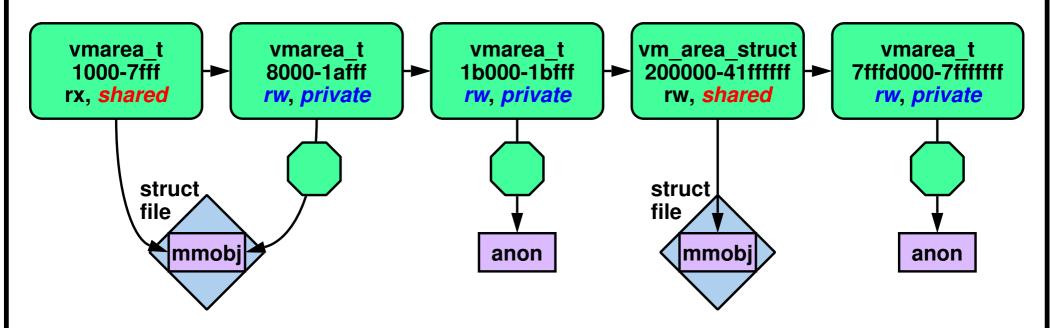
Read-only mapping of a file (e.g. text)

 pages come from the file, but, since they are never modified, they never need to be written back



Read-write shared mapping of a file (e.g. via mmap () system call)

 pages come from the file, modified pages are written back to the file





weenix supports this type of "backing store"



Backing Up Pages (2)



- Read-write *private* mapping of a file (e.g. the data section as well as memory mapped private by the mmap () system call)
- pages come from the file, but modified pages, associated with shadow objects, must be backed up in swap space



- Anonymous memory (e.g. bss, stack, and shared memory), also *privately* mapped
- pages are created as zero fill on demand; they must be backed up in swap space
 - modified pages of these, associated with shadow objects, must be backed up in swap space



- weenix does not support this type of backing store
- need to prevent the pageout daemon to free up these pages accidentically
 - simply move them out of the pageout daemon's way using pframe_pin()

Swap Space



Swap space management possibilities

- radically-conservative approach: Eager Evaluation (or pre-allocation)
 - backing-store space is allocated when virtual memory is allocated
 - page outs always succeed
 - disadvantage: might need to have much more backing store than needed
- radically-liberal approach: Lazy Evaluation
 - backing-store space is allocated only when needed
 - advantage: can get by with minimal backing-store space
 - disadvantage: page outs could fail because of no space and process gets killed at a seemingly random time



Swap Space



Space management possibilities

- mixed approach: e.g., reserve stack space for a thread in Windows
 - the address space for the thread stack is first "reserved"
 - no backing store actually created, but space is reserved so no other thread can use the reserved space
 - when part of this address space is used, it's "committed" (backing store is actually allocated)



- by default, done with eager evaluation in Windows and most Unix/Linux systems
- both systems provide means for lazy evaluation as well



Space Allocation in Linux



Total memory = primary + swap space



System-wide parameter: overcommit_memory

- three possibilities
 - maybe (default)
 - always
 - never



mmap has MAP_NORESERVE flag

don't worry about over-committing





Space Allocation in Windows



allocation of virtual memory



- reservation of physical resources
 - paging space + physical memory



no over-commitment

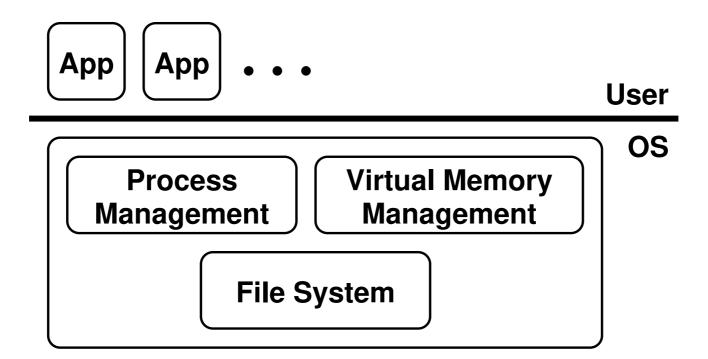


creator specifies both reservation and commitment for stack pages





Summary



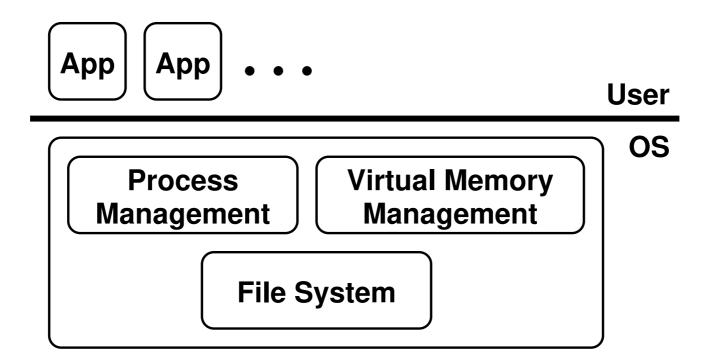


The subsystems are inter-related

- file systems uses threads managed by the process subsystem
- file systems uses buffer cache (managed by the memory subsystem)
- memory subsystem uses threads to do background work
- process subsystem keeps track of data structures related to files and virtual memory on behalf of processes



Summary





- think of everything that happens in these subsystems when you type "1s" into a console
- Kernel 3 is where everything comes together
 - although we are already using page tables in earlier assignments (see pt_init())

