

4.2 Rethinking Operating-System Structure

- Virtual Machines
- Microkernel



Copyright © William C. Cheng



- Virtual Machines
- A nicely designed and implemented monolithic OS is great
 - but that's not the reality
- Major problem with a monolithic OS *implementation*
 - bugs in one component can adversely affect another component
 - worse if large number of programmers contribute code
 - some coders are not as good as others
 - good coders have bad days
- Modern OSs *isolate* applications from one another
 - code executing in the privileged mode can do things the user mode code cannot
 - e.g., invoking privileged instructions
 - if you invoke a privileged instruction in user mode, you will cause a violation and trap into the kernel
- Can the same kind of *isolation* be provided for *OS components*?
 - if yes, at what cost? (there is no free lunch)

Copyright © William C. Cheng

Monolithic Kernel

- Major advantage of monolithic kernel
 - performance
- Major down side of monolithic kernel
 - reliability (i.e., buggy kernel)
- Proposal to fix the reliability problem
 - shrink the code in "privileged mode"
- Two major approaches
 - virtual machines*
 - microkernel*



Copyright © William C. Cheng

It's 1964 ...

- IBM wants a *multituser time-sharing system*

- TSS (Time-Sharing System) project
 - it's a very difficult system to build
 - large, monolithic system
 - lots of people working on it
 - for years
 - total, complete flop

- CMS
 - single-user time-sharing system* for IBM 360
- CP67
 - virtual machine monitor (VMM)*
 - supports multiple virtual IBM 360s

- Put the two together ...
 - a (working) *multituser time-sharing system*



Copyright © William C. Cheng

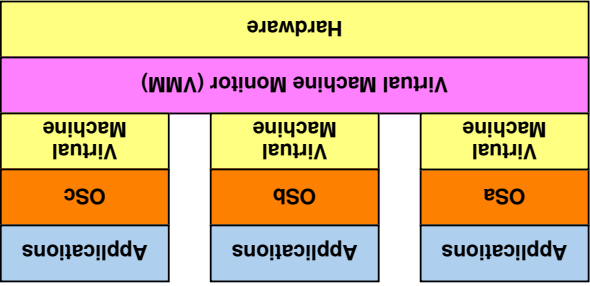
Virtual Machines

Virtual Machines Part I: > 45 Years Ago

- Had a different motivation



Copyright © William C. Cheng



- In OS, a "*monitor*" is a *synchronization construct* that allows executing entities to have both *mutual exclusion* and the ability to *wait* (block) for a certain condition to become true
- What abstraction does a *virtual machine* provide?



Copyright © William C. Cheng

Copyright © William C. Cheng

11

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Processor in the virtual machine is the *real* processor

instructions are *executed* (and not interpreted or emulated)

traps are generated just as they are on real machines

in a real machine, trap handler is indexed by the trap number into a hardware-mandated jump table

the VMM needs to find the address of the virtual machine's trap handler in the table and transfer control to it

interrupts pretty much work the same way

"Virtual Machine" in the picture contains:

- virtual CPU, virtual disk, virtual keyboard, etc.
- data structures* that *represent HW components*

VMM

Operating Systems - CSCI 402

Copyright © William C. Cheng

12

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Execute a *non-privileged* instruction

e.g., "add", "mul"

VMM Operations

Operating Systems - CSCI 402

Copyright © William C. Cheng

9

Applications

Guest OS

Virtual Machine

VMM

Hardware

Applications

Guest OS

Virtual Machine

VMM

Hardware

Run the *entire virtual machine* in *user mode* of the real machine

VMM runs in the *privileged* mode of the *real* machine

VMM keeps track of whether each virtual machine is in the *virtual privileged mode* or in the *virtual user mode*

OS runs in the *(virtual) privileged* mode of the *virtual* machine

Applications runs in the *(virtual) user* mode of the *virtual* machine

How?

Operating Systems - CSCI 402

Copyright © William C. Cheng

10

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Processor in the virtual machine is the *real* processor

instructions are *executed* (and not interpreted or emulated)

traps are generated just as they are on real machines

in a real machine, trap handler is indexed by the trap number into a hardware-mandated jump table

VMM

Operating Systems - CSCI 402

Copyright © William C. Cheng

7

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Applications

Guest Osa

Virtual Machine

VMM

Hardware

In OS, a "*monitor*" is a *synchronization construct* that allows executing entities to have both *mutual exclusion* and the ability to *wait* (block) for a certain condition to become true

What abstraction does a *virtual machine* provide?

= *hardware*

Virtual Machines

Operating Systems - CSCI 402

Copyright © William C. Cheng

8

Applications

Guest Osa

Virtual Machine

VMM

Hardware

Applications

Guest Osa

Virtual Machine

VMM

Hardware

A single user time-sharing system could be developed independently of the VMM

and it can be tested on a real machine (which behaves identical to the VM)

no ambiguity about the interface th VMM must provide to its applications - *identical to the real machine!*

Virtual Machines

Operating Systems - CSCI 402

Copyright © William C. Cheng

17

Execute a *privileged* instruction (e.g., "trap" (due to "fork"))

the VMM is invoked
the VMM figures out which VM is currently executing
the VMM then asks the corresponding VM to deliver the trap to its OS

VMM Operations

Operating Systems - CSCI 402

Copyright © William C. Cheng

15

Execute a *non-privileged* instruction (e.g., "add", "mul")

executes directly on hardware or on hardware perspective, no difference running in VM

VMM Operations

Operating Systems - CSCI 402

Copyright © William C. Cheng

13

Execute a *non-privileged* instruction (e.g., "add", "mul")

executes directly on hardware or on hardware perspective, no difference running in VM

VMM Operations

Operating Systems - CSCI 402

Copyright © William C. Cheng

16

Execute a *privileged* instruction (e.g., "trap" (due to "fork"))

the VMM is invoked
the VMM figures out which VM is currently executing
the VMM then asks the corresponding VM to deliver the trap to its OS

VMM Operations

Operating Systems - CSCI 402

Copyright © William C. Cheng

16

Execute a *privileged* instruction (e.g., "trap" (due to "fork"))

the VMM is invoked
the VMM figures out which VM is currently executing
the VMM then asks the corresponding VM to deliver the trap to its OS

VMM Operations

Operating Systems - CSCI 402

Copyright © William C. Cheng

14

Execute a *non-privileged* instruction (e.g., "add", "mul")

executes directly on hardware or on hardware perspective, no difference running in VM

VMM Operations

Operating Systems - CSCI 402

Applications

OS

VM

VMM

Hardware

Disk

What about I/O?

- e.g., read ()
- read () eventually reaches the OS
- the OS asks for a block on the virtual disk

VMM Operations

Copyright © William C. Cheng

Applications

OS

VM

VMM

Hardware

Disk

What about I/O?

- e.g., read ()

real disk is divvy up among the virtual machines

- each VM has a virtual disk

VMM Operations

Copyright © William C. Cheng

Applications

OS

VM

VMM

Hardware

Disk

What about I/O?

- e.g., read ()

VMM Operations

Copyright © William C. Cheng

Applications

Guest OSa

Virtual Machine

VMM

Hardware

Disk

What about I/O?

- e.g., read ()

real disk is divvy up among the virtual machines

- each VM has a virtual disk

VMM Operations

Copyright © William C. Cheng

Applications

Guest OSa

Virtual Machine

VMM

Hardware

Disk

What about I/O?

- e.g., read ()

real disk is divvy up among the virtual machines

- each VM has a virtual disk

VMM Operations

Copyright © William C. Cheng

Applications

Guest OSa

Virtual Machine

VMM

Hardware

Disk

What about I/O?

- e.g., read ()

real disk is divvy up among the virtual machines

- each VM has a virtual disk

VMM Operations

Copyright © William C. Cheng

Copyright © William C. Cheng 29

Processor Virtualization Requirements

- Processor in the virtual machine is the real processor
- instructions are executed (and not interpreted or emulated)
- traps are generated just as they are on real machines
- in a real machine, trap handler is indexed by the trap number into a hardware-mandated jump table
- the VMM needs to find the address of the virtual machine's trap handler in the table and transfer control to it
- interrupts pretty much work the same way

➡ Pretty much everything can be worked out except for one problem

- if a virtual machine is executing in the privileged mode, what's to prevent it from changing how the memory-mapping resources that have been set up?

➡ Need to distinguish between *sensitive instructions* and *privileged instructions*

Operating Systems - CSCI 402

Copyright © William C. Cheng 27

Requirements

- A virtual machine is an efficient, isolated duplicate of real machine
- requires faithful virtualization of pretty much all components
- processor
- memory
- interval timers
- I/O devices
- etc.
- this is "*pure*" virtualization
- costly

➡ *Paravirtualization*:

- virtualized entity is a bit different from the real entity
- so as to enhance scalability, performance, and simplicity
- it is probably aware that it's not running on a real machine

Operating Systems - CSCI 402

Copyright © William C. Cheng 25

VMM Operations

What about I/O?

- e.g., read ()
- read () eventually reaches the OS
- the OS asks for a block on the virtual disk

traps into VMM

- the VMM *emulates* the instruction (i.e., translates it into a request for the real disk)
- there's really *no* disk in the VM

Operating Systems - CSCI 402

Copyright © William C. Cheng 30

Processor Virtualization Requirements

- cause privileged-instruction trap when executed in user mode
- execute fully when the processor is in privileged mode

➡ *Privileged Instructions*:

- Control-sensitive instructions*:
- instructions that affect allocation of (real) system resources
- such as instructions that change the mapping of virtual to real memory

➡ *Behavior-sensitive instructions*:

- instructions whose effect depends on the allocation of (1) such as instructions that returns the real address of a location in virtual memory
- instructions whose effect depends on the current processor mode
- such as x86's popf instructions that sets a set of processor flags when run in privileged mode, but set a different set of flags otherwise

Operating Systems - CSCI 402

Copyright © William C. Cheng 28

Processor Virtualization Requirements

- Virtualizing the processor requires:
- 1) multiplexing the real processor among the virtual machines
- relatively straightforward
- 2) making each virtual machine behaves just like a real machine
- all instructions must work identically
- generation of and response to traps and interrupts must be identical as well

Operating Systems - CSCI 402

Copyright © William C. Cheng 26


Why Virtual Machine?

- it's a good structuring technique for a multi-user system
- many advantages
- OS debugging and testing
- run a production OS in a VM, accessible to users
- test a new OS in a separate VM, accessible to developers
- Adapt to hardware changes in software
- Multiple OSes on one machine
- one type of applications run really well in one OS
- another type of applications run really well in a different OS
- one physical machine can support both, no user need to suffer today, it's common that a machine in the cloud would run multiple Linux OS instances and multiple Windows OS instances
- Server consolidation and service isolation
- web hosting, security concerns
- cloud computing

Operating Systems - CSCI 402

Virtual Machines

Part 2: Now



Copyright © William C. Cheng

36

Operating Systems - CSC1402

Virtual Devices?

- Terminals
 - connecting (real) people
- Networks
 - didn't exist in the 60s
 - (how did virtual machines communicate?)
- Disk drives
 - CP67 supported "mini disks"
 - extended at Brown into "segment system"
- Interval timer
 - virtual or real?

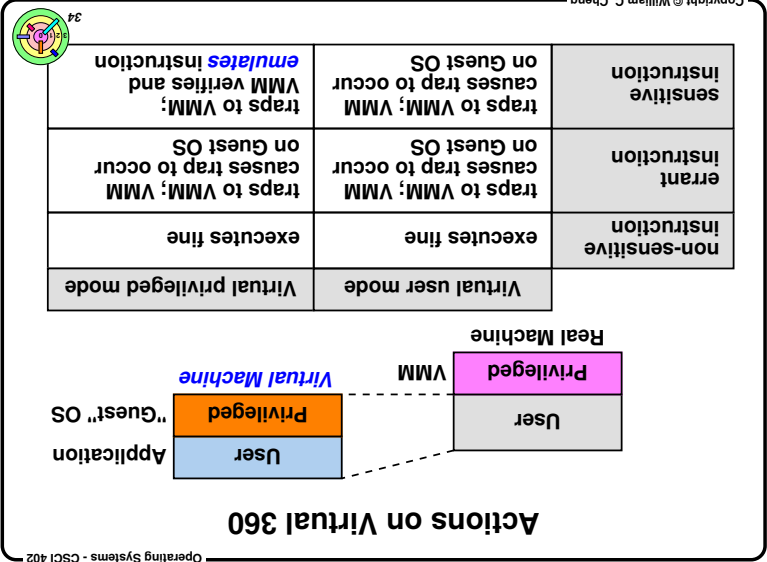
Copyright © William C. Cheng

- ## Virtual Devices?
- Terminals
 - connecting (real) people
 - Networks
 - didn't exist in the 60s
 - (how did virtual machines communicate?)
 - Disk drives
 - CP67 supported "mini disks"
 - extended at Brown into "segment system"
 - Interval timer
 - virtual or real?
- Copyright © William C. Cheng

[illegible]

	User	Privileged
Virtual Machine	executes fine	traps to kernel
Real Machine	traps to kernel	traps to kernel

The diagram illustrates how actions performed by user or privileged processes are handled differently depending on whether they are running in a virtual machine or directly on the real hardware. In a virtual machine, the hypervisor allows user-level operations to proceed normally while intercepting and handling privileged operations. On real hardware, no such abstraction exists; all privileged operations must be managed by the operating system's kernel.



The diagram illustrates the mapping between Virtual Machine (VM) components and Real Machine components during actions on Virtual 360.

Virtual Machine Components (Left):

- User** (Blue box)
- Privileged** (Orange box)
- "Guest" OS** (Label for the Privileged box)
- Application** (Label for the User box)

Real Machine Components (Right):

- User** (Grey box)
- Privileged** (Pink box)
- VMM** (Label for the Privileged box)
- Real Machine** (Label for the entire right side)

Actions on Virtual 360:

- Virtual Machine Privileged** (Label for the Privileged box)
- Virtual Machine** (Label for the entire left side)

Mapping:

- The **User** box in the Virtual Machine maps to the **User** box in the Real Machine.
- The **Privileged** box in the Virtual Machine maps to the **Privileged** box in the Real Machine.
- The **VMM** label in the Real Machine points to the **Privileged** box.

The diagram illustrates the mapping between Virtual Machine (VM) components and Real Machine components during actions on Virtual 360.

Virtual Machine Components (Left):

- User** (Blue box)
- Privileged** (Orange box)
- "Guest" OS** (Label for the Privileged box)
- Application** (Label for the User box)

Real Machine Components (Right):

- User** (Grey box)
- Privileged** (Pink box)
- VMM** (Label for the Privileged box)
- Real Machine** (Label for the entire right side)

Actions on Virtual 360:

- Virtual Machine Privileged** (Label for the Privileged box)
- Virtual Machine** (Label for the entire left side)

Mapping:

- The **User** box in the Virtual Machine maps to the **User** box in the Real Machine.
- The **Privileged** box in the Virtual Machine maps to the **Privileged** box in the Real Machine.
- The **VMM** label in the Real Machine points to the **Privileged** box.



Processor Virtualization Requirements

[Popek and Goldberg, 1974] *proved* that the *sufficient condition* to be able to construct a virtual machine is simply the following:



- a computer's set of *sensitive instructions* is a *subset* of its *privileged instructions*
- i.e., if you execute a sensitive instruction in user mode, you will trap into the kernel
- more importantly, if you execute a sensitive instruction in *virtual user or virtual privileged mode*, you will trap into *VMM*

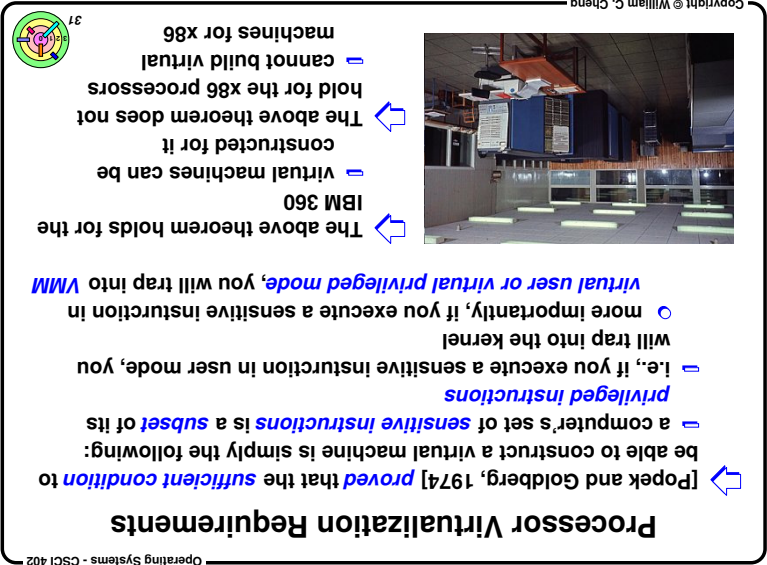
The above theorem holds for the IBM 360

- virtual machines can be constructed for it
- The above theorem does not hold for the x86 processors
- cannot build virtual machines for x86



Copyright © William C. Chen

- # Processor Virtualization Requirements
- [Popek and Goldberg, 1974] *proved* that the *sufficient condition* to be able to construct a virtual machine is simply the following:
- a computer's set of *sensitive instructions* is a *subset* of its *privileged instructions*
 - i.e., if you execute a sensitive instruction in user mode, you will trap into the kernel
 - more importantly, if you execute a sensitive instruction in *virtual user or virtual privileged mode*, you will trap into *VMM*
- The above theorem holds for the IBM 360
- virtual machines can be constructed for it
 - The above theorem does not hold for the x86 processors
 - cannot build virtual machines for x86
- 
- 
- Copyright © William C. Chen



The (Real) 360 Architecture

- Two execution modes
 - supervisor and problem (user)
 - all sensitive instructions are privileged instructions
- Memory is protectable: 2KB granularity
- All interrupt vectors and the clock are in first 512 bytes of memory
- I/O done via channel programs in memory, initiated with privileged instructions
- Dynamic address translation (virtual memory) added for Model 67

- The (Real) 360 Architecture**
- Two execution modes
 - supervisor and problem (user)
 - all sensitive instructions are privileged instructions
 - Memory is protectable: 2KB granularity
 - All interrupt vectors and the clock are in first 512 bytes of memory
 - I/O done via channel programs in memory, initiated with privileged instructions
 - Dynamic address translation (virtual memory) added for Model 67



- guest OS *source code* is *modified* (and recompiled)
- to simply mean "*VMM*" (even without paravirtualization)
- ◊ please note that some would use the term "*hypervisor*"
- *hypervisor* interface between virtual and real machines
- provides more convenient interfaces for virtualization
- virtual machine differs from real machine
- *Paravirtualization*
- fix the hardware so it's virtualizable
- *Hardware virtualization*
- no need to modify guest OS
- ◊ VMWare does this
- do so dynamically (i.e., *dynamic* binary rewriting)
- replace sensitive instructions with hypercalls
- rewrite kernel binaries of guest OSes
- *Binary rewriting*

What to Do?

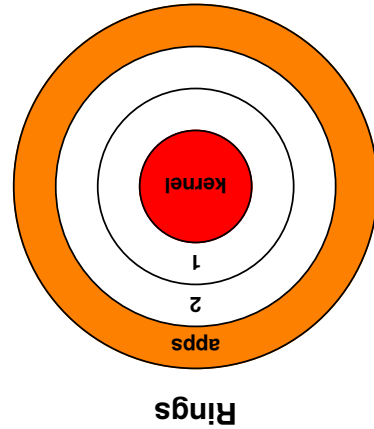


- *Privilege-mode code run via binary translator*
- guest OS is unmodified
- replaces sensitive instructions with hypercalls
- translated code is cached
- usually translated just once
- *VMWare*
- U.S. patent 6,397,242
- VirtualBox appears to do something similar to VMWare
- see <https://www.virtualbox.org/manual/ch10.html#idp58764736>
- for more details

Binary Rewriting



- An x86 processor can be in one of 4 *modes*

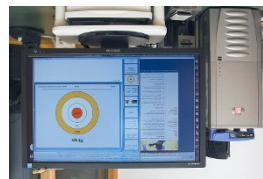


Rings



- *popf*
- pops flags (word) off stack, setting processor flags according to word's content
- sets all flags if in ring 0
- ◊ including interrupt-disable flag
- just some of them if in other rings
- ◊ ignores interrupt-disable flag
- bad news: if invoked in *user* mode, does *not* cause a *trap*
- therefore, this instruction will execute differently in the OS when it's running on top of a VM (as compared to running on a real machine)
- ◊ since the OS is running in user mode under the virtual memory scheme
- this is one of the major problem to virtualize x86 systems
- there is another major problem (related to device I/O)

A Sensitive x86 Instruction



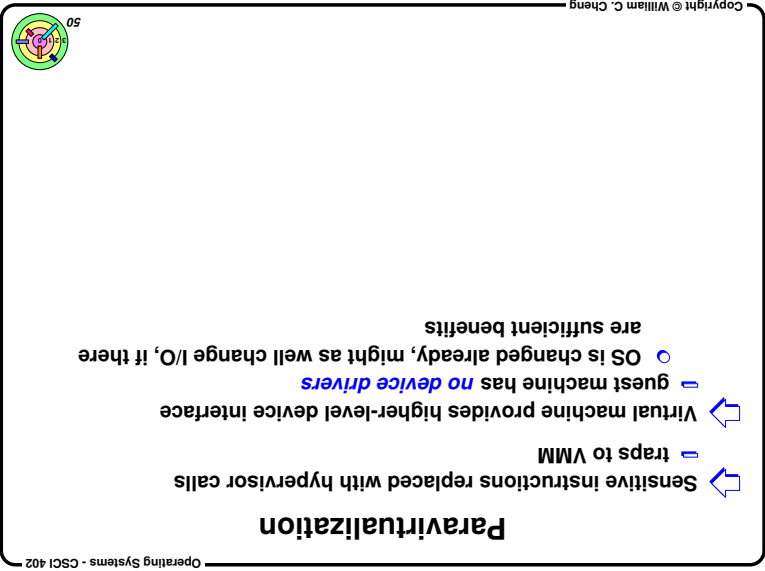
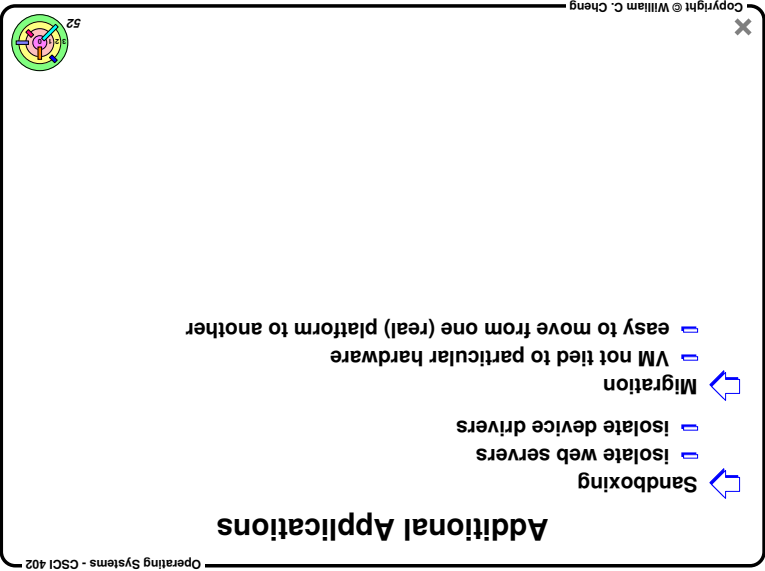
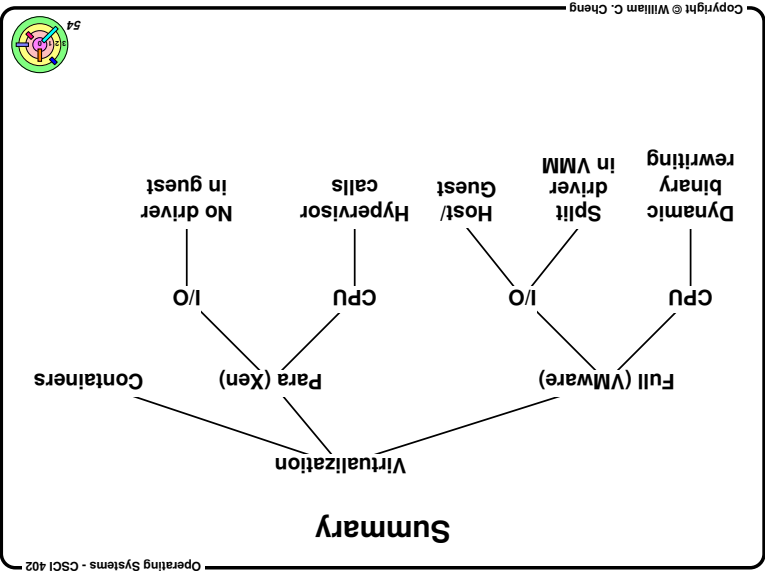
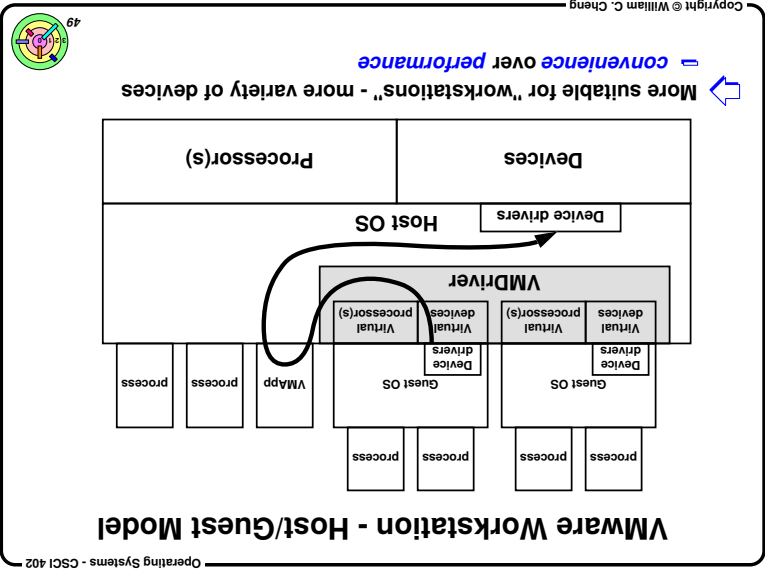
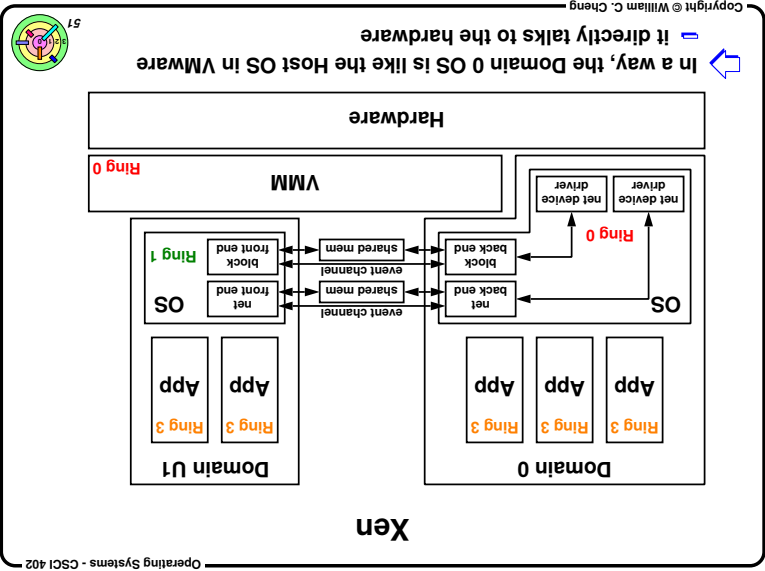
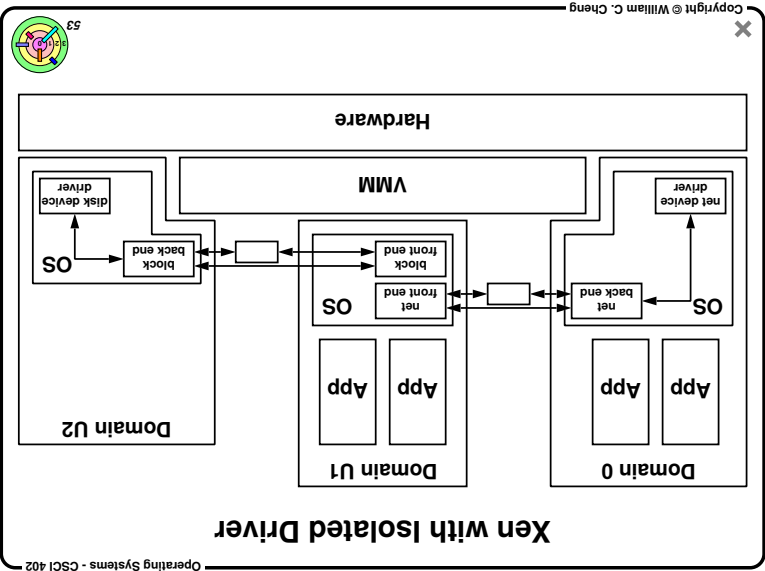
≠

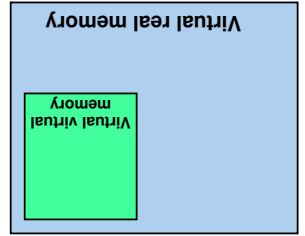


How They Are Different

- Two execution modes
- supervisor and problem (user)
- *all* sensitive instructions are privileged instructions
- Memory is protectable:
- 2k-byte granularity
- All interrupt vectors and the clock are in first 512 bytes of memory
- I/O done via *channel programs*
- in memory, initiated with privileged instructions
- Dynamic address translation (virtual memory) added for *Model 67*
- *Intel x86*
- Four execution modes
- rings 0 through 3
- *not all* sensitive instructions are privileged instructions
- Memory is protectable:
- segment system + virtual memory
- Special register points to interrupt table
- I/O done via *memory-mapped* addresses
- Virtual memory is standard

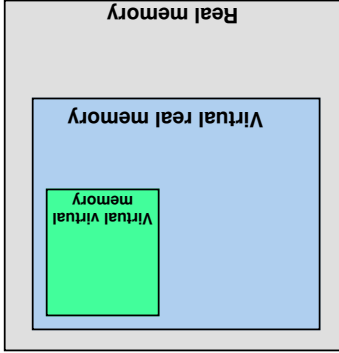






- A user process thinks it's
accessing virtual memory
— but it's really dealing with
virtual real memory
- The OS in a VM thinks it's
managing real memory
— but it's really dealing with
virtual real memory

Virtual Machines Meet Virtual Memory



- A user process thinks it's accessing virtual memory
- but it's really dealing with *virtual virtual memory*
- The OS in a VM thinks it's managing real memory
- but it's really dealing with *virtual real memory*
- VMM needs to manage real memory
- how can we *virtualize* *virtual memory*?

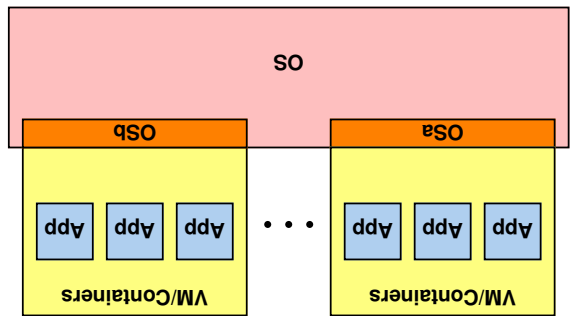
Virtual Machines Meet Virtual Memory

7.2.6 Virtualizing Virtual Memory



- ➡ A user process thinks it's accessing virtual memory
- ➡ but it's really dealing with virtual virtual memory

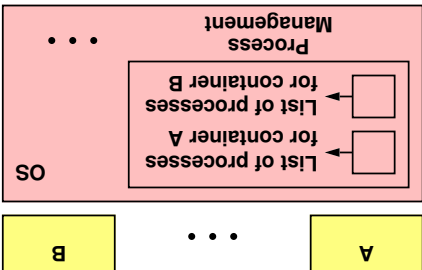
Virtual Machines Meet Virtual Memory

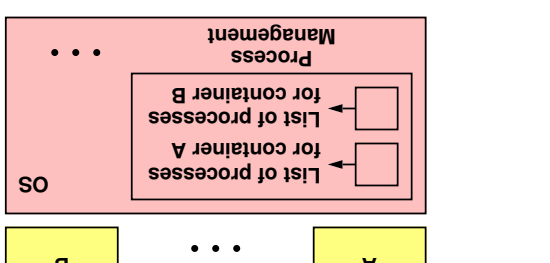


- the OS provides the abstraction that each container runs on top of a separate OS (but there is really only one OS)
- e.g., OpenVZ, Linux Containers (LXC), Docker

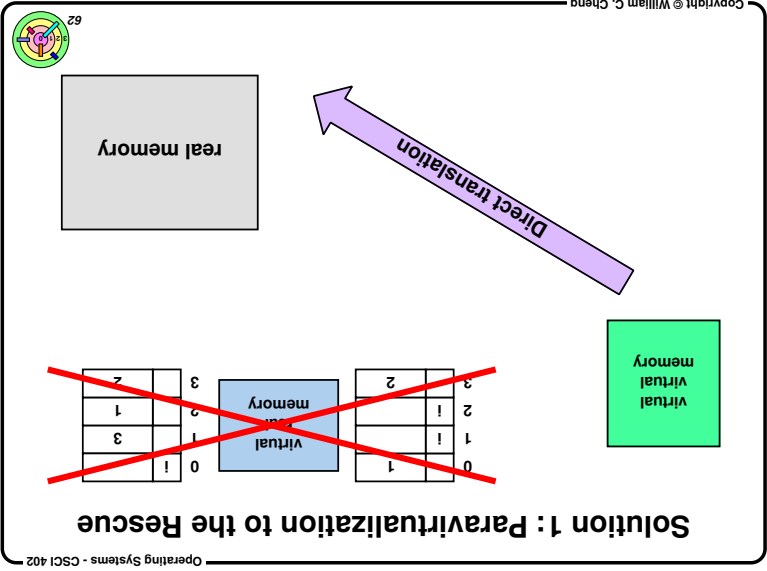
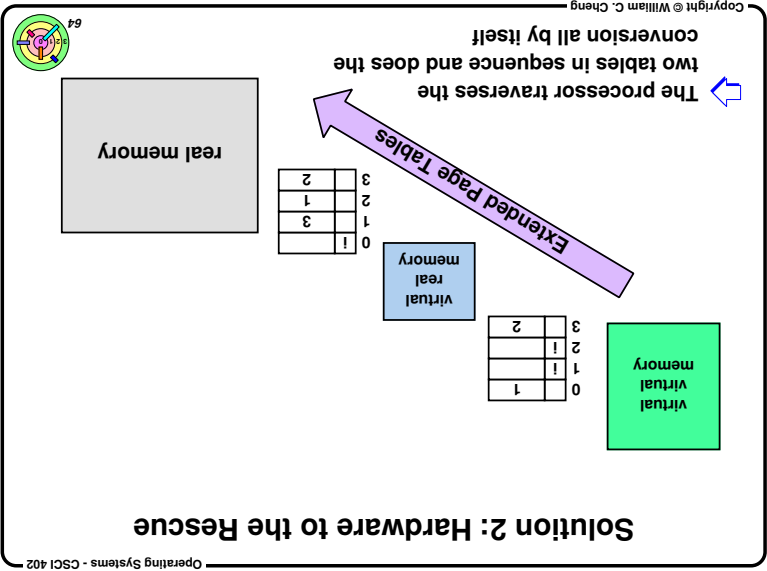
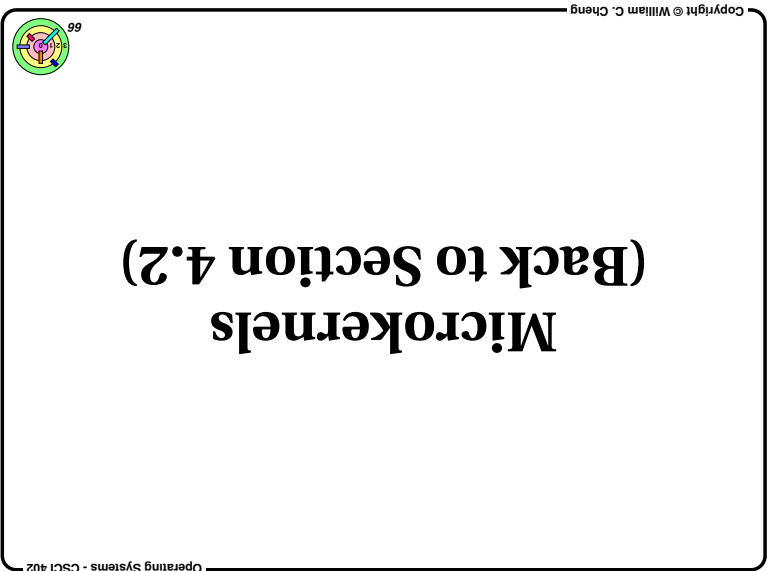
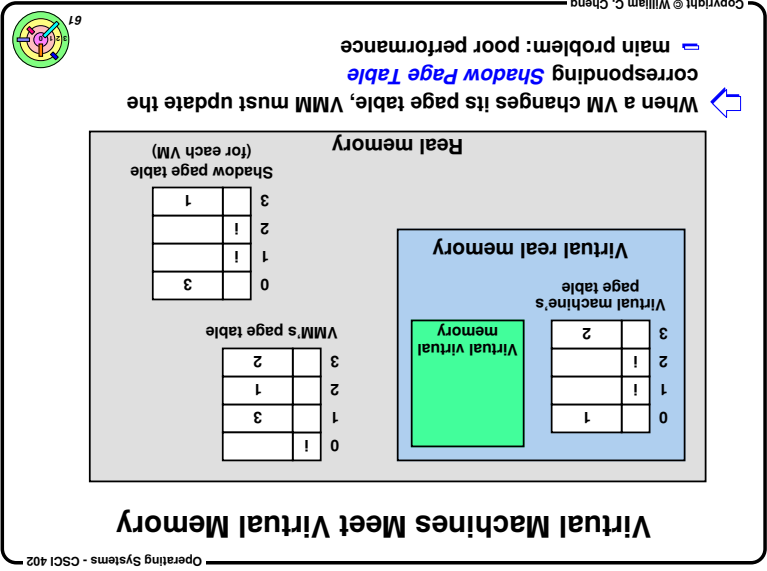
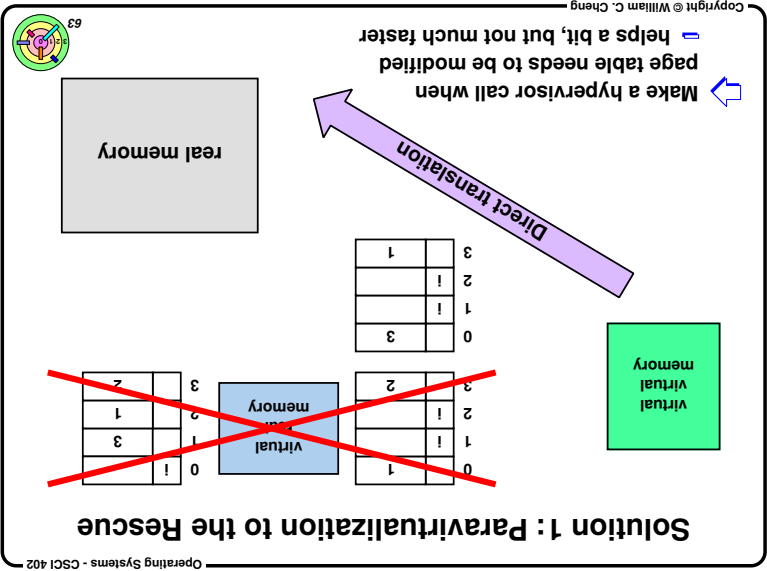
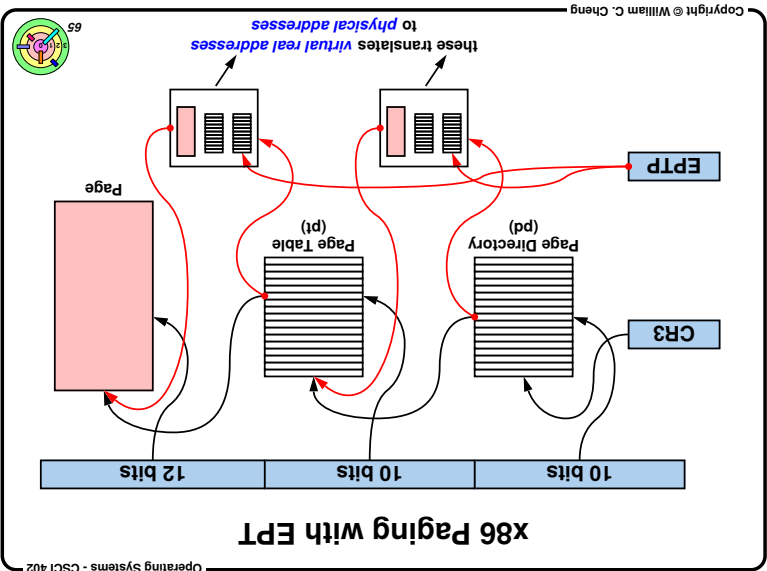
One More Kind Of Virtualization

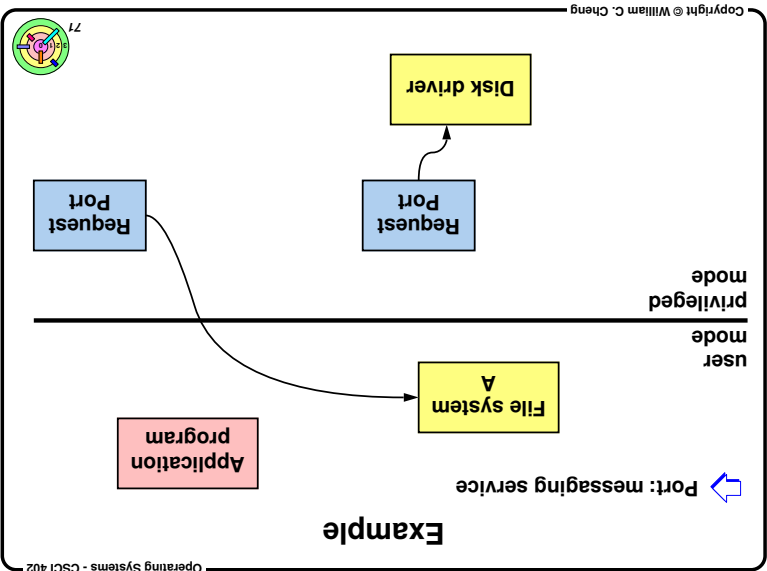
- Containerized OS (or OS Containers)
- not covered in textbook



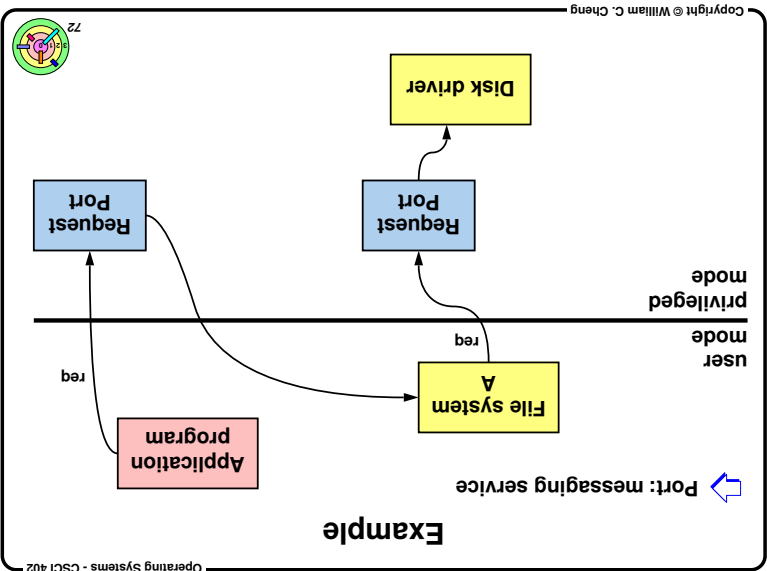
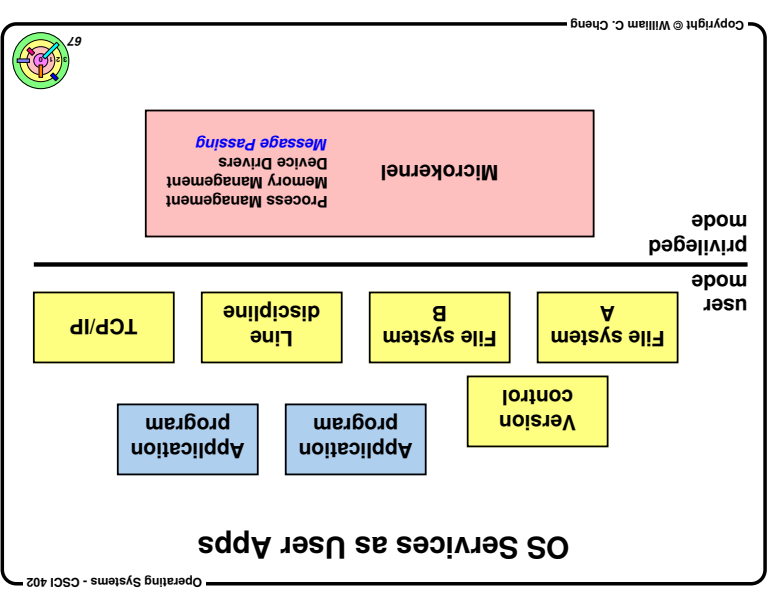
- ## Containerized OS
- Within the OS, the management of resources for each container is separated
- e.g., processes for container A is kept separate from processes for container B
- 
- The diagram illustrates a containerized operating system architecture. At the top, a large pink rectangle represents the 'OS'. Inside this rectangle, at the bottom, are two yellow boxes labeled 'A' and 'B', representing separate containers. Above these containers, within the OS box, is a smaller pink rectangle labeled 'Process Management'. This box contains two separate lists: 'List of processes for container A' and 'List of processes for container B', each with an arrow pointing to a small square box. The text 'OS' is written in the bottom right corner of the large pink rectangle. To the right of the diagram, there are two blue arrows pointing right, followed by the text 'Others may consider this "virtual machine", but we shouldn't because "guest OS" does not run in user space'.
- Others may consider this "virtual machine", but we shouldn't because "guest OS" does not run in user space

Containerized OS





- Copyright © William C. Cheng
- ### Implementation Issues
- How are modules linked together?
 - e.g., how would you implement `read()`/`write()`?
 - can't use system calls any more!
 - e.g., which file system supports `read()`/`write()`?
 - How is data moved around efficiently?
- 69



- Copyright © William C. Cheng
- ### Mach
- Developed at CMU, then Utah
 - Early versions shared kernel with Unix
 - basis of NeXT OS
 - Later versions still shared kernel with Unix
 - basis of OSF/1
 - basis of Macintosh OS X
 - Even later versions actually functioned as working microkernel
 - basis of GNU/HURD project
 - HURD: HIRD of Unix-replacing daemons
 - HIRD: HURD of interfaces representing depth
- 70

- Copyright © William C. Cheng
- ### Why?
- It's cool ...
 - Assume that OS coders are incompetent, malicious, or both ...
 - OS components run as protected user-level applications
 - easier to add, modify, and extend user-level components than kernel components
 - Extensibility
- 68



RPC

- Ports used to implement *remote procedure calls*
 - communication across process boundaries
 - if procedures are on same machine ...
 - local RPC



Successful Microkernel Systems

- ...
- ?
- ?



Mach Ports (1)

- Communication construct
 - client create *response port* and *capability* (like a key) to send data through it
 - include capability in the request message to server



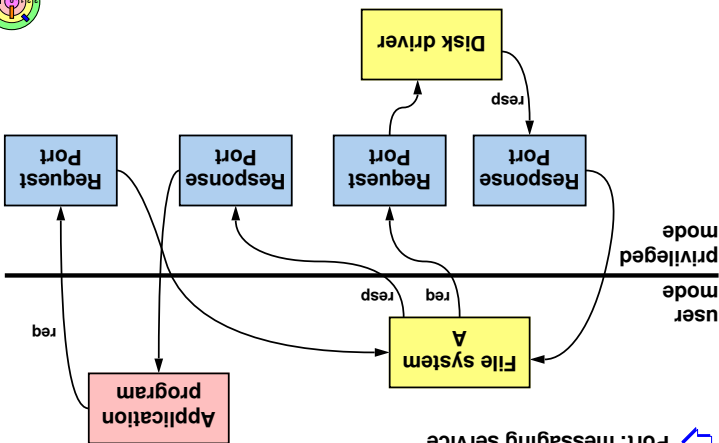
Mach Ports (1)

- Communication construct
 - client create *response port* and *capability* (like a key) to send data through it
 - include capability in the request message to server



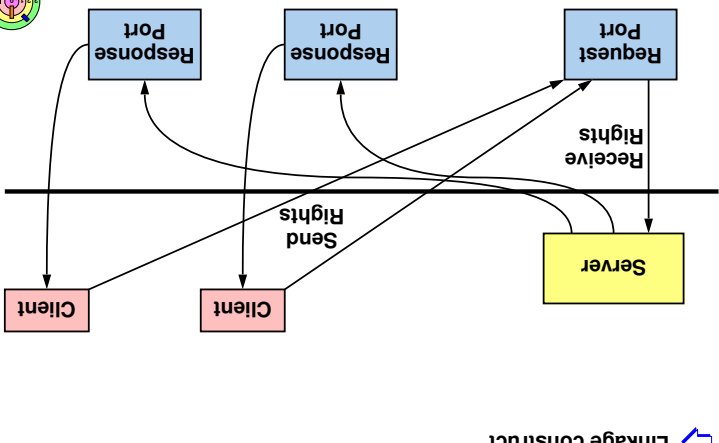
Example

- Port: messaging service



Mach Ports (1)

- ← Linkage construct



Attempts

- Windows NT 3.1
 - = graphics subsystem ran as user-level process
 - = moved to kernel in 4.0 for performance reasons
- Macintosh OS X
 - = based on Mach
 - = all services in kernel for performance reasons
- HURD
 - = based on Mach
 - = services implemented as user processes
 - = no one uses it, for performance reasons ...

