

CSE 451: Operating Systems

Winter 2013

Hardware Trends and OS Support

Gary Kimura

Even coarse architectural trends impact tremendously the design of systems

- Processing power
 - doubling every 18 months (not Moore's Law)
 - 60% improvement each year
 - factor of 100 every decade

 - 1980: 1 MHz Apple II+ == \$2,000
 - 1980 also 1 MIPS VAX-11/780 == \$120,000
 - 2008: Intel Quad-Core 2.66GHz == \$900
 - 2009: err... make that \$700
 - 2012: err... make that \$155

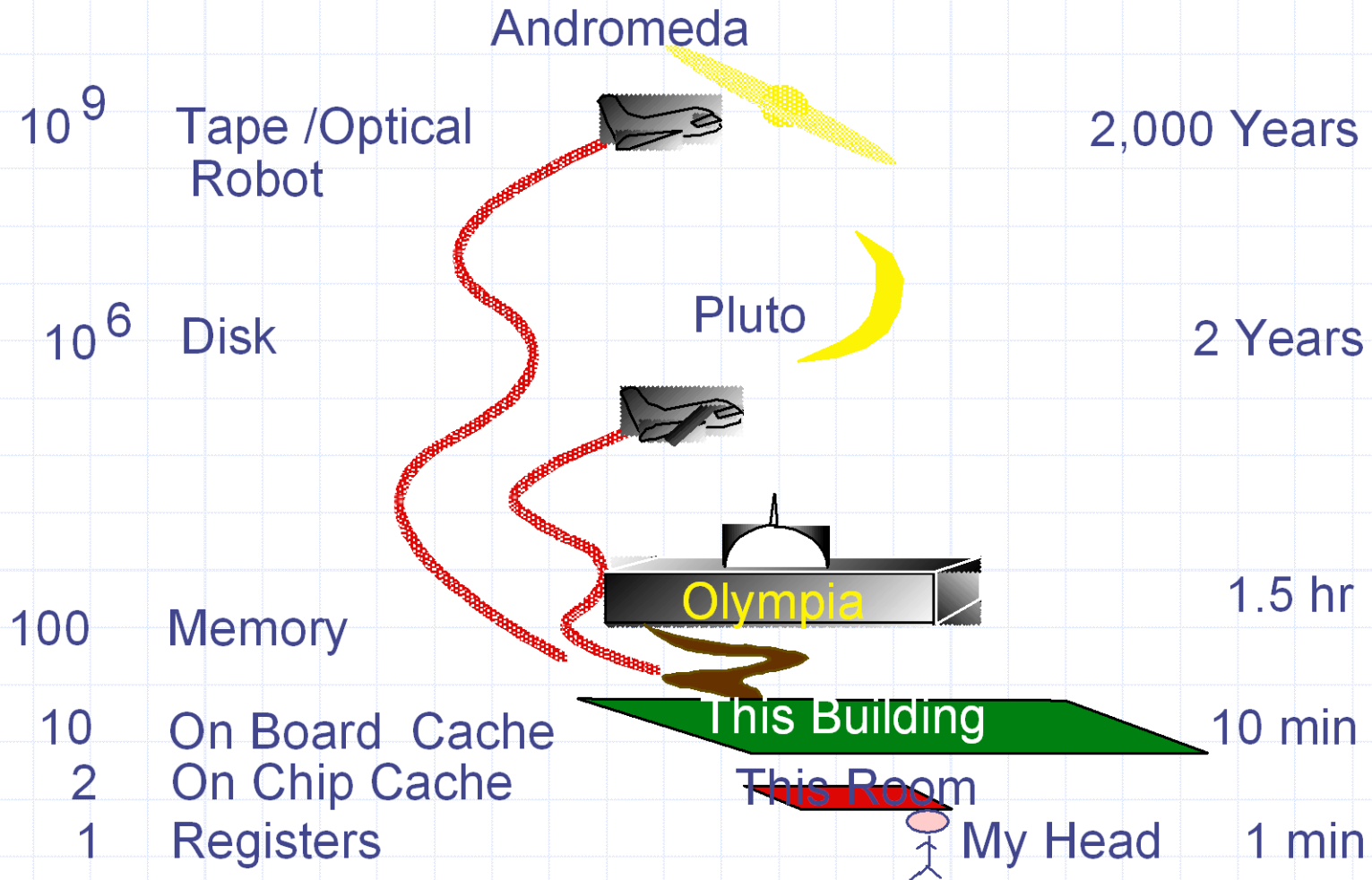


- Primary memory capacity
 - same story, same reason (Moore's Law)
 - 1972: 1MB = \$1,000,000
 - 1982: 4 MB for DECSYSTEM 20 - \$60,000
 - 2009: 2 GB for Dell Inspiron - \$75
 - 2010: 2 GB for Dell Inspiron - \$36

- Disk capacity, 1975-1989
 - doubled every 3+ years
 - 25% improvement each year
 - factor of 10 every decade
 - Still exponential, but far less rapid than processor performance
- Disk capacity since 1990
 - doubling every 12 months
 - 100% improvement each year
 - factor of 1000 every decade
 - 10x as fast as processor performance!
- Disk Performance has NOT kept up
 - Speed in 1983: 500KB/sec
 - Speed in 2011: ~100MB/sec (only 200x!)

- Optical network bandwidth today
 - *Doubling every 9 months*
 - 150% improvement each year
 - Factor of 10,000 every decade
 - 10x as fast as disk capacity!
 - 100x as fast as processor performance!!
- What are some of the implications of these trends?
 - Just one example: We have always designed systems so that they “spend” processing power in order to save “scarce” storage and bandwidth!

Storage Latency: How Far Away is the Data?



Lower-level architecture affects the OS even more dramatically

- The operating system supports **sharing** and **protection**
 - multiple applications can run concurrently, sharing resources
 - a buggy or malicious application can't nail other applications or the system
- There are many approaches to achieving this
- The architecture determines which approaches are viable (reasonably efficient, or even possible)
 - includes instruction set (synchronization, I/O, ...)
 - Chip (local memory, sharing caches)
 - also hardware components like MMU or DMA controllers

- Architectural support can vastly simplify (or complicate!) OS tasks
 - e.g.: early PC operating systems (DOS, MacOS) lacked support for virtual memory, in part because at that time PCs lacked necessary hardware support
 - Apollo workstation used two CPUs as a bandaid for non-restartable instructions!
 - Until 2006, Intel-based PCs still lacked support for 64-bit addressing (which has been available for a decade on other platforms: MIPS, Alpha, IBM, etc...)
 - changing rapidly due to AMD's 64-bit architecture

Architectural Features affecting OS's

- (CSE351 review) These features were built primarily to support OS's:
 - timer (clock) operation
 - synchronization instructions (e.g., atomic test-and-set)
 - memory protection
 - I/O control operations
 - interrupts and exceptions
 - disabling hardware interrupts
 - protected modes of execution (kernel vs. user)
 - protected and privileged instructions
 - system calls (and software interrupts)

Privileged Instructions

- some instructions are restricted to the OS
 - known as **protected** or **privileged** instructions
- e.g., only the OS can:
 - directly access I/O devices (disks, network cards)
 - why?
 - manipulate memory state management
 - page table pointers, TLB loads, etc.
 - why?
 - manipulate special ‘mode bits’
 - interrupt priority level
 - why?
 - halt instruction
 - why?

OS Protection

- So how does the processor know if a protected instruction should be executed?
 - the architecture must support at least two modes of operation: **kernel** mode and **user** mode
 - VAX, x86 support 4 protection modes in hardware
 - Multics supported 4 modes in software
 - why more than 2?
 - mode is set by status bit in a protected processor register
 - user programs execute in user mode
 - OS executes in kernel mode (OS == kernel)
- Protected instructions can only be executed in the kernel mode
 - what happens if user mode executes a protected instruction?

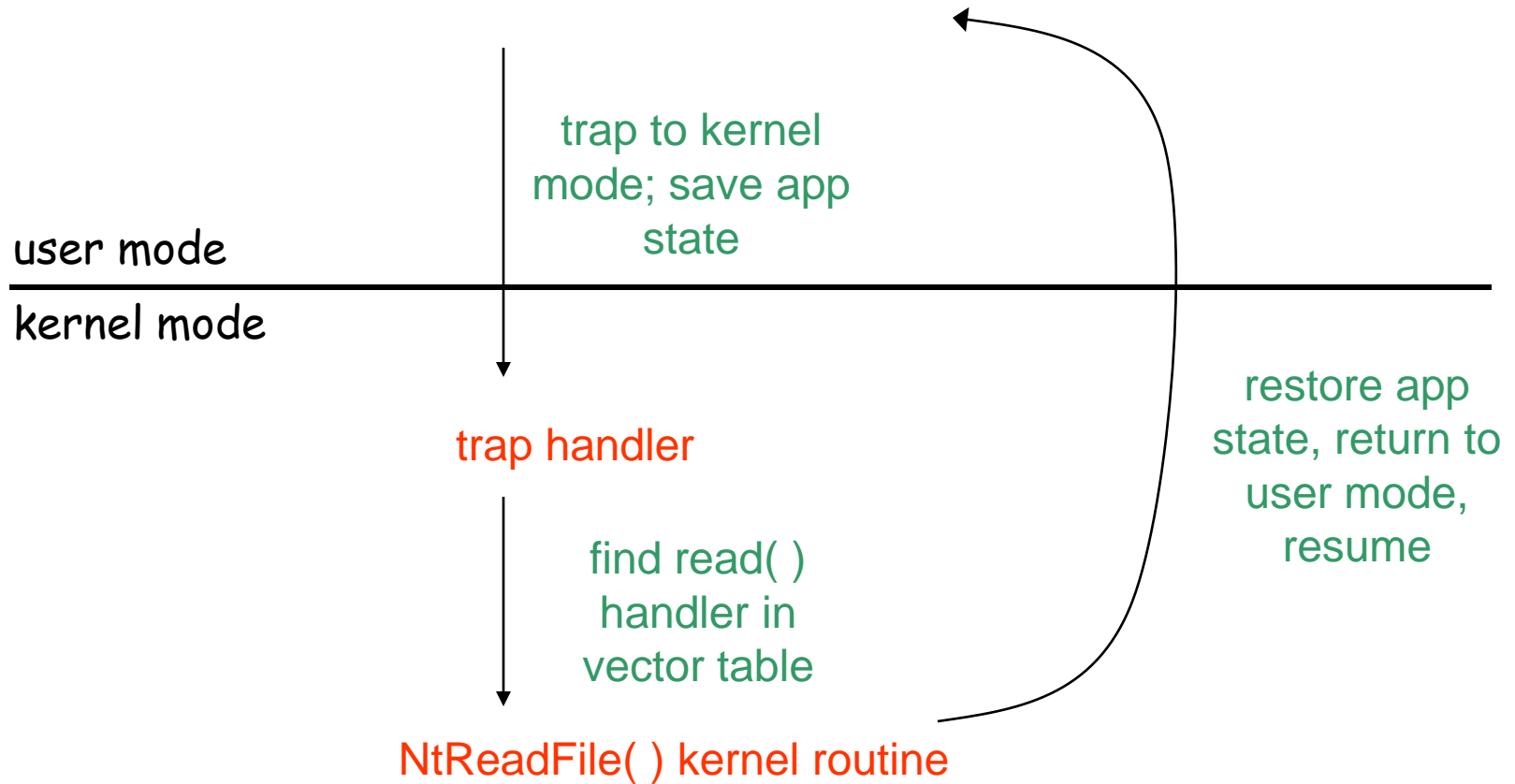
Crossing Protection Boundaries

- So how do user programs do something privileged?
 - e.g., how can you write to a disk if you can't do I/O instructions?
- User programs must call an OS procedure
 - OS defines a sequence of **system calls**
 - how does the user-mode to kernel-mode transition happen?
- There must be a system call instruction
 - Looks a lot like a protected procedure call

- The magic instruction atomically:
 - Saves the current IP
 - Sets the execution mode to kernel
 - Sets the IP to a handler address
- With that, it's a lot like a local procedure call
 - Caller puts arguments in a place callee expects (registers or stack)
 - One of the args is a system call number, indicating which OS function to invoke
 - Callee (OS) saves caller's state (registers, other control state) so it can use the CPU
 - OS function code runs
 - OS must verify caller's arguments (e.g., pointers)
 - OS returns using a special instruction
 - Automatically sets PC to return address and sets execution mode to user

A Kernel Crossing Illustrated

App: ReadFile(Handle, Buffer, Count, &BytesRead, Overlapped)



System Call Issues

- What would be wrong if the instruction was like a regular subroutine call that specified the address of the routine in kernel mode to execute?
- What would happen if kernel didn't save state?
- Why must the kernel verify arguments?
- How can you reference kernel objects as arguments or results to/from system calls?

Interrupts and Exceptions

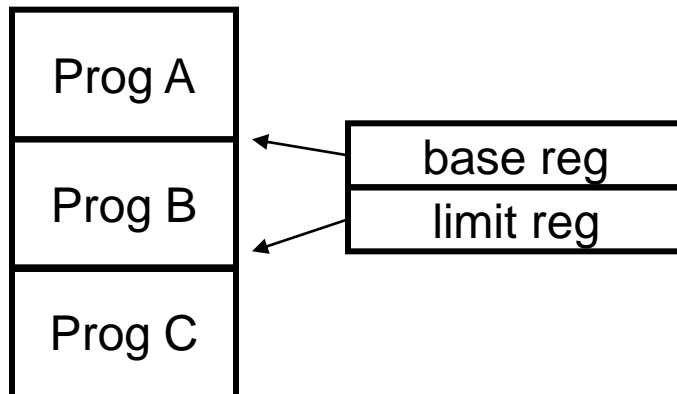
- Two main types of events: **interrupts** and **exceptions**
 - exceptions are caused by software executing instructions
 - e.g. the x86 'int' instruction
 - e.g. a page fault, write to a read-only page
 - an expected (i.e., related to the instructions in your code) exception is a “**trap**”, an unexpected exception is a “**fault**”
 - interrupts are caused by hardware devices
 - e.g. device finishes I/O
 - e.g. timer fires

Exceptions

- Hardware must detect special conditions
 - page fault, write to a read-only page, overflow, divide by zero, trace trap, odd address trap, privileged instruction trap, syscall...
- Must transfer control to handler within the OS
- Hardware must save state on fault
 - faulting process can be restarted afterwards
- VM traps are useful
 - Modern operating systems use VM traps for many functions: debugging, distributed VM, garbage collection, copy-on-write...
- Exceptions are a performance optimization
 - Lazy, let the hardware detect things
 - Expensive, conditions could be detected by inserting extra instructions in the code

Memory Protection

- OS must protect user programs from each other
 - maliciousness, ineptitude
- OS must also protect itself from user programs
 - integrity and security
 - what about protecting user programs from OS?
- Simplest scheme: **base** and **limit** registers
 - are these protected?



base and limit registers
are loaded by OS before
starting program

More sophisticated memory protection

- coming later in the course
- virtual memory
 - paging, segmentation
 - page tables, page table pointers
 - translation lookaside buffers (TLBs)
 - page fault handling

OS control flow

- After the OS has booted, all entry to the kernel happens as the result of an **event**
 - event immediately stops current execution
 - changes mode to kernel mode, event handler is called
 - fundamentally, the OS doesn't actually DO anything, just reacts!
- Kernel defines handlers for each event type
 - specific types are defined by the architecture
 - e.g.: timer event, I/O interrupt, system call trap
 - when the processor receives an event of a given type, it
 - transfers control to handler within the OS
 - handler saves program state (PC, regs, etc.)
 - handler functionality is invoked
 - handler restores program state, returns to program

I/O Control

- Issues:
 - how does the kernel start an I/O?
 - special I/O instructions
 - memory-mapped I/O
 - how does the kernel notice an I/O has finished?
 - polling
 - interrupts
- Interrupts are basis for **asynchronous** I/O
 - device performs an operation asynchronously to CPU
 - device sends an interrupt signal on bus when done
 - in memory, a **vector table** contains list of addresses of kernel routines to handle various interrupt types
 - who populates the vector table, and when?
 - CPU switches to address indicated by vector specified by interrupt signal

I/O Control (continued)

device interrupts



CPU stops current operation, switches to kernel mode, and saves current PC and other state on kernel stack



CPU fetches proper vector from vector table and branches to that address (to routine to handle interrupt)



interrupt routine examines device database and performs action required by interrupt



handler completes operation, restores saved (interrupted state) and returns to user mode (or calls scheduler to switch to another program)

Timers

- How can the OS prevent runaway user programs from hogging the CPU (infinite loops?)
 - use a hardware timer that generates a periodic interrupt
 - before it transfers to a user program, the OS loads the timer with a time to interrupt
 - “quantum”: how big should it be set?
 - when timer fires, an interrupt transfers control back to OS
 - at which point OS must decide which program to schedule next
 - very interesting policy question: we’ll dedicate a class to it
- Should the timer be privileged?
 - for reading or for writing?

Synchronization

- Interrupts cause a wrinkle:
 - may occur any time, causing code to execute that interferes with code that was interrupted
 - OS must be able to **synchronize** concurrent processes
- Synchronization:
 - guarantee that short instruction sequences (e.g., read-modify-write) execute atomically
 - one method: turn off interrupts before the sequence, execute it, then re-enable interrupts
 - architecture must support disabling interrupts
 - in user-mode???
 - multi-core?!?
 - another method: have special complex atomic instructions
 - read-modify-write
 - test-and-set
 - load-linked store-conditional

“Concurrent programming”

- Management of **concurrency** and **asynchronous** events is biggest difference between “**systems programming**” and “**traditional application programming**”
 - modern “event-oriented” application programming is a middle ground (CSE333)
- Arises from the architecture
- Can be sugar-coated, but cannot be totally abstracted away
- Huge intellectual challenge
 - Unlike vulnerabilities due to buffer overruns, which are just sloppy programming

Architectures are still evolving

- New features are still being introduced to meet modern demands
 - Support for virtual machine monitors
 - Hardware transaction support (to simplify parallel programming)
 - Support for security (encryption, trusted modes)
 - Increasingly sophisticated video / graphics
 - Other stuff that hasn't been invented yet...
- In current technology transistors are free – CPU makers are looking for new ways to use transistors to make their chips more desirable
- Intel's big challenge: finding applications that require new hardware support, so that you will want to upgrade to a new computer to run them

Some questions

- Why wouldn't you want a user program to be able to access an I/O device (e.g., the disk) directly?
- OK, so what keeps this from happening? What prevents user programs from directly accessing the disk?
- So, how does a user program cause disk I/O to occur?
- What prevents a user program from scribbling on the memory of another user program?
- What prevents a user program from scribbling on the memory of the operating system?
- What prevents a user program from running away with the CPU?

More Kernel/User Mode Switching Modes

- From user-mode to kernel
 - Interrupts
 - Triggered by timer and I/O devices
- Exceptions
 - Triggered by unexpected program behavior
 - Or malicious behavior!
- System calls (aka protected procedure call)
 - Request by program for kernel to do some operation on its behalf
 - Only limited # of very carefully coded entry points

Mode Switch

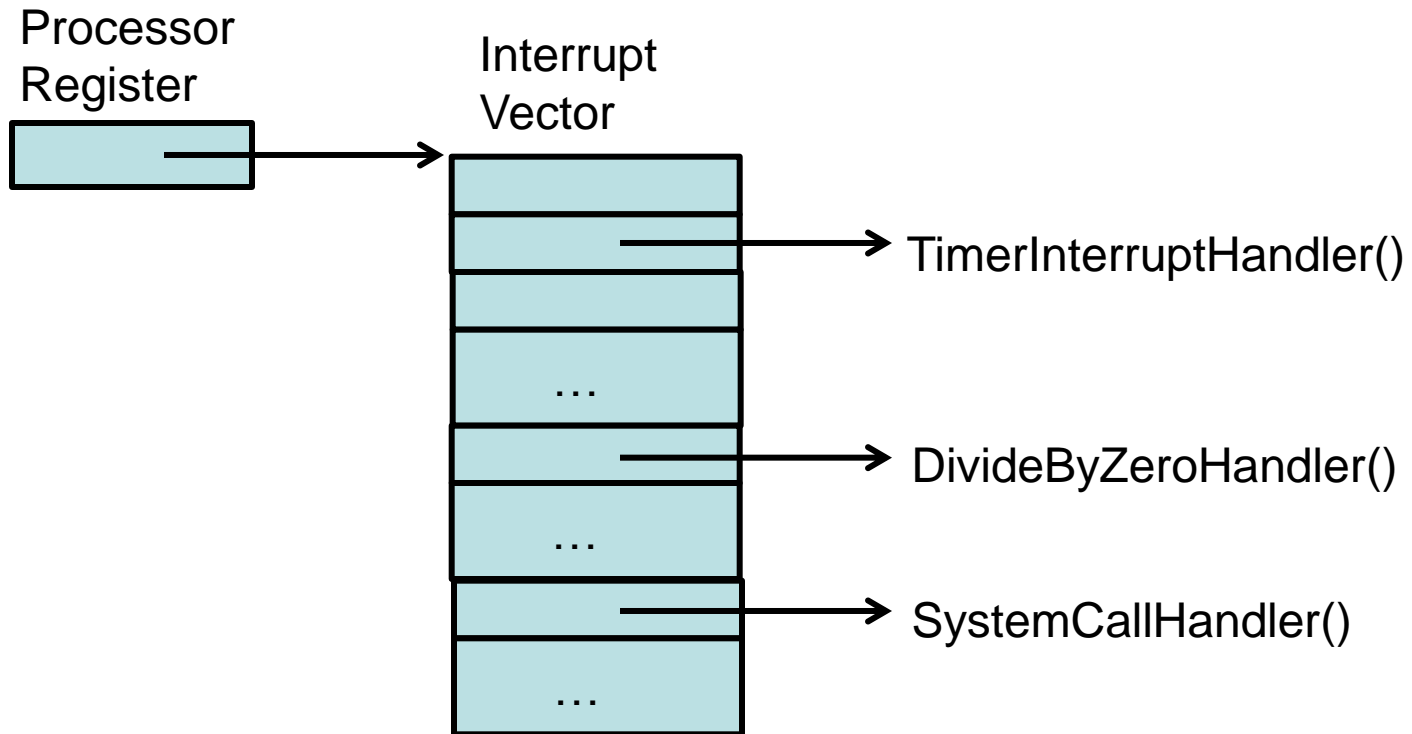
- From kernel-mode to user
 - New process/new thread start
 - Jump to first instruction in program/thread
 - Return from interrupt, exception, system call
 - Resume suspended execution
 - Process/thread context switch
 - Resume some other process
 - User-level upcall
 - Asynchronous notification to user program

How do we take interrupts safely?

- Interrupt vector
 - Limited number of entry points into kernel
- Kernel Interrupt stack
 - Handler works regardless of state of user code
- Interrupt masking
 - Handler is non-blocking
- Atomic transfer of control
 - Single instruction to change:
 - Program counter
 - Stack pointer
 - Memory Protection
 - Kernel/user mode
- Transparent restartable execution
 - User program does not know interrupt occurred

Interrupt Vector

- Table set up by OS kernel; pointers to code to run on different events



Interrupt Stack

- Per-processor, located in kernel (not user) memory
 - Usually a thread had both: kernel and user stack
- Why can't interrupt handler run on the stack of the interrupted user process?