





Locality In A Memory-Reference Pattern

Kubiatowicz CS162 © UCB Fall 2020

- Program Memory Access Patterns have temporal and spatial locality
 - Group of Pages accessed along a given time slice called the "Working Set"
 - Working Set defines minimum number of pages for process to behave well
- Not enough memory for Working Set \Rightarrow Thrashing
 - Better to swap out process?





- If Δ too small will not encompass entire locality
- if Δ too large will encompass several localities
- if Δ = $\infty \Longrightarrow$ will encompass entire program
- D = Σ|WSi| ≡ total demand frames
- if D > m \Rightarrow Thrashing
 - Policy: if D > m, then suspend/swap out processes

Kubiatowicz CS162 © UCB Fall 2020

- This can improve overall system behavior by a lot!

Lec 17.11

What about Compulsory Misses? Linux Memory Details? Memory management in Linux considerably more complex than the Recall that compulsory misses are misses that occur the first time that a examples we have been discussing page is seen Memory Zones: physical memory categories - Pages that are touched for the first time - ZONE DMA: < 16MB memory, DMAable on ISA bus - Pages that are touched after process is swapped out/swapped back in - ZONE NORMAL: $16MB \rightarrow 896MB$ (mapped at 0xC000000) Clustering: - ZONE HIGHMEM: Everything else (> 896MB) - On a page-fault, bring in multiple pages "around" the faulting page • Each zone has 1 freelist, 2 LRU lists (Active/Inactive) - Since efficiency of disk reads increases with sequential reads, makes Many different types of allocation sense to read several sequential pages - SLAB allocators, per-page allocators, mapped/unmapped Working Set Tracking: Many different types of allocated memory: - Use algorithm to try to track working set of application - Anonymous memory (not backed by a file, heap/stack) - When swapping process back in, swap in working set - Mapped memory (backed by a file) Allocation priorities - Is blocking allowed/etc 10/26/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 17.13 10/26/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 17.14

Linux Virtual memory map (Pre-Meltdown)



10/26/20

Pre-Meltdown Virtual Map (Details)

- · Kernel memory not generally visible to user
 - Exception: special VDSO (virtual dynamically linked shared objects) facility that maps kernel code into user space to aid in system calls (and to provide certain actual system calls such as gettimeofday())
- · Every physical page described by a "page" structure
 - Collected together in lower physical memory
 - Can be accessed in kernel virtual space
 - Linked together in various "LRU" lists
- · For 32-bit virtual memory architectures:
- When physical memory < 896MB
 - » All physical memory mapped at 0xC0000000
- When physical memory >= 896MB
 - » Not all physical memory mapped in kernel space all the time
 - » Can be temporarily mapped with addresses > 0xCC000000
- · For 64-bit virtual memory architectures:
 - All physical memory mapped above 0xFFFF80000000000

Kubiatowicz CS162 © UCB Fall 2020



Requirements of I/O

- · So far in CS 162, we have studied:
 - Abstractions: the APIs provided by the OS to applications running in a process
 - Synchronization/Scheduling: How to manage the CPU
- What about I/O?
 - Without I/O, computers are useless (disembodied brains?)
 - But... thousands of devices, each slightly different
 » How can we standardize the interfaces to these devices?
 - Devices unreliable: media failures and transmission errors
 » How can we make them reliable???
 - Devices unpredictable and/or slow
 - » How can we manage them if we don't know what they will do or how they will perform?



Recall: OS Basics: I/O



Lec 17.19

10/26/20

Recall: Range of Timescales

Jeff Dean: "Numbers **Everyone Should** Know"

L1 cache reference	0.	.5 ns
Branch mispredict	5	ns
L2 cache reference	7	ns
Mutex lock/unlock	25	ns
Main memory reference	100	ns
Compress 1K bytes with Zippy	3,000	ns
Send 2K bytes over 1 Gbps network	20,000	ns
Read 1 MB sequentially from memory	250,000	ns
Round trip within same datacenter	500,000	ns
Disk seek	10,000,000	ns
Read 1 MB sequentially from disk	20,000,000	ns
Send packet CA->Netherlands->CA	150,000,000	ns

10/26/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 17.21

10/26/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 17.22

In a Picture



- I/O devices you recognize are supported by I/O Controllers
- Processors accesses them by reading and writing IO registers as if they were memory
 - -Write commands and arguments, read status and results

Lec 17.23

Example: Device Transfer Rates in Mb/s (Sun Enterprise 6000)

- Device rates vary over 12 orders of magnitude!!!
- · System must be able to handle this wide range
 - Better not have high overhead/byte for fast devices
 - Better not waste time waiting for slow devices





Modern I/O Systems



Lec 17.24



PCI Bus Evolution



PCI started life out as a bus

- · But a parallel bus has many limitations
 - Multiplexing address/data for many requests
 - Slowest devices must be able to tell what's happening (e.g., for arbitration)
 - Bus speed is set to that of the slowest device

PCI Express "Bus"

- No longer a parallel bus
- Really a collection of fast serial channels or "lanes"
- Devices can use as many as they need to achieve a desired bandwidth
- · Slow devices don't have to share with fast ones
- One of the successes of device abstraction in Linux was the ability to migrate from PCI to PCI Express
 - The physical interconnect changed completely, but the old API still worked

Lec 17.27



Port-Mapped I/O in Pintos Speaker Driver

Pintos: devices	/speaker.c		Pintos: threads/io.h
13 /* Sets the PC speaker to emit a	tone at the given FREQUENCY, in		
14 Hz. */		7	/* Reads and returns a byte from PORT. */
15 void		8	<pre>static inline uint8_t</pre>
<pre>16 speaker_on (int frequency)</pre>		9	<pre>inb (uint16_t port)</pre>
17 (10	{
<pre>18 if (frequency >= 20 && frequer</pre>	cy <= 20000)	11	/* See [TA32-v2a] "IN", */
19 { 20 (* Cat the tiles along 1 d	hadden and an all an all an all and all all all all all	12	viet0 t deter
20 /- Set the timer channel t	the given ERECHENCY then	12	uinto_t uata;
22 connect the timer chan	el output to the speaker. */	13	asm volatile ("inb %w1, %b0" : "=a" (data) : "Nd" (port
23 enum intr level old level	= intr disable ():	14	return data;
<pre>24 pit_configure_channel (2,</pre>	3, frequency);	15	}
25 outb (SPEAKER_PORT_GATE, 1	nb (SPEAKER_PORT_GATE) SPEAKER_GATE_ENABLE);		
26 intr_set_level (old_level)	3		
27 }			
28 else		64	/* Writes byte DATA to PORT. */
29 {		65	static inline void
30 /* FREQUENCY is outside the second sec	e range of normal human hearing.	66	outh (wint16 t port wint8 t data)
31 Just turn off the speak	er. */	67	((anero_e pore, aneo_e aaea)
32 speaker_ott ();		67	1
24		68	/* See [IA32-v2b] "OUT". */
35		69	asm volatile ("outb %b0, %w1" : : "a" (data), "Nd" (
36 /* Turn off the PC speaker, by c	isconnecting the timer channel's	70	}
37 output from the speaker. */			
38 void			
<pre>39 speaker_off (void)</pre>			
40 (
41 onum_intr_level old_levelic 42 outb (SPEAKER_PORT_GATE, inb (<pre>tr_disable (); SPEAKER_PORT_GATE) & ~SPEAKER_GATE_ENABLE);</pre>		
<pre>43 intr_set_level (old_level);</pre>			
	14 11 1 1 00100 011		

Example: Memory-Mapped Display Controller

· Memory-Mapped: 0x80020000 - Hardware maps control registers and display memory into Graphics physical address space Comman Queue » Addresses set by HW jumpers or at boot time 0x80010000 - Simply writing to display memory (also called the "frame Display Memory buffer") changes image on screen » Addr: 0x8000F000 - 0x8000FFFF 0x8000F000 - Writing graphics description to cmd gueue » Say enter a set of triangles describing some scene Command 0x0007F004 » Addr: 0x80010000 - 0x8001FFFF Status 0x0007F000 - Writing to the command register may cause on-board graphics hardware to do something » Say render the above scene Physical » Addr: 0x0007F004 Address Can protect with address translation Space 10/26/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 17.32

There's more than just a CPU in there!

10/26/20

10/26/20

Chip-scale Features of 2015 x86 (Sky Lake)



Transferring Data To/From Controller

Programmed I/O:

- Each byte transferred via processor in/out or load/store
- Pro: Simple hardware, easy to program
- Con: Consumes processor cycles proportional to data size

Direct Memory Access:

- Give controller access to memory bus
- Ask it to transfer

(from OSC book):

10/26/20

- data blocks to/from memory directly
- DMA controller 2 device driver tells transfers bytes to buffer X, increasing disk controller to transfer C bytes memory address from disk to buffer cache and decreasing C at address X until C = 0DMA/bus/ when C = 0, DMA interrupts CPU to signal memory^x buffer interrupt - 2 memory bus controller transfer completion Sample interaction with DMA controller . disk controller initiates DMA transfer IDE disk controller disk controller sends each byte to DMA disk disk controller disk) disk

device driver is told

to transfer disk data

to buffer at address)

Transferring Data To/From Controller

Programmed I/O:

- Each byte transferred via processor in/out or load/store
- Pro: Simple hardware, easy to program
- Con: Consumes processor cycles proportional to data size

Direct Memory Access:

- Give controller access to memory bus
- Ask it to transfer data blocks to/from memory directly
- Sample interaction with DMA controller (from OSC book):



Kubiatowicz CS162 © UCB Fall 2020

Lec 17.38

I/O Device Notifying the OS

Kubiatowicz CS162 © UCB Fall 2020

- The OS needs to know when:
 - The I/O device has completed an operation
 - The I/O operation has encountered an error
- I/O Interrupt:
 - Device generates an interrupt whenever it needs service
- Pro: handles unpredictable events well
- Con: interrupts relatively high overhead
- Polling:
 - OS periodically checks a device-specific status register
 - » I/O device puts completion information in status register
 - Pro: low overhead
- Con: may waste many cycles on polling if infrequent or unpredictable I/O operations
- Actual devices combine both polling and interrupts
 - For instance High-bandwidth network adapter:
 - » Interrupt for first incoming packet
 - » Poll for following packets until hardware queues are empty

Kernel Device Structure



10/26/20

Lec 17.37

10/26/20

10/26/20

