# CS 350 Operating Systems Course Notes

Fall 2013

David R. Cheriton School of Computer Science University of Waterloo

	What is an Operating System?	
• Three views of an o	operating system	
<b>Application View:</b>	what services does it provide?	
System View: what	at problems does it solve?	
Implementation V	<b>'iew:</b> how is it built?	
An operating system	n is part cop, part facilitator.	



# **Other Views of an Operating System**

System View: The OS manages the hardware resources of a computer system.

- Resources include processors, memory, disks and other storage devices, network interfaces, I/O devices such as keyboards, mice and monitors, and so on.
- The operating system allocates resources among running programs. It controls the sharing of resources among programs.
- The OS itself also uses resources, which it must share with application programs.

Implementation View: The OS is a concurrent, real-time program.

- Concurrency arises naturally in an OS when it supports concurrent applications, and because it must interact directly with the hardware.
- Hardware interactions also impose timing constraints.

CS350

**Operating Systems** 

Fall 2013

<ul> <li>Some terminology:</li> <li>kernel: The operating system kernel is the part of the operating system that responds to system calls, interrupts and exceptions.</li> <li>operating system: The operating system as a whole includes the kernel, and may include other related programs that provide services for applications. This may include things like: <ul> <li>utility programs</li> <li>command interpreters</li> <li>programming libraries</li> </ul> </li> </ul>
<ul> <li>kernel: The operating system kernel is the part of the operating system that responds to system calls, interrupts and exceptions.</li> <li>operating system: The operating system as a whole includes the kernel, and may include other related programs that provide services for applications. This may include things like: <ul> <li>utility programs</li> <li>command interpreters</li> <li>programming libraries</li> </ul> </li> </ul>
<ul> <li>operating system: The operating system as a whole includes the kernel, and may include other related programs that provide services for applications. This may include things like:</li> <li>utility programs</li> <li>command interpreters</li> <li>programming libraries</li> </ul>





	<b>Course Outline</b>	
• Introduction		
• Threads and Concurr	rency	
• Synchronization		
• Processes and the Ke	ernel	
• Virtual Memory		
• Scheduling		
• Devices and Device	Management	
• File Systems		
• Interprocess Commu	nication and Networking (time permitti	ng)

	<b>Review: Program Execution</b>	
• Registers		
- program counter,	stack pointer,	
• Memory		
– program code		
<ul> <li>program data</li> </ul>		
<ul> <li>program stack co</li> </ul>	ntaining procedure activation records	
• CPU		
– fetches and execu	tes instructions	

CS350	

Jpe ng Sys

				<b>Review: MIPS Register Usage</b>	
R0,	zero	=	##	zero (always returns 0)	
R1,	at	=	##	reserved for use by assembler	
R2,	v0	=	##	return value / system call number	
R3,	v1	=	##	return value	
R4,	a0	=	##	1st argument (to subroutine)	
R5,	al	=	##	2nd argument	
R6,	a2	=	##	3rd argument	
R7,	a3	=	##	4th argument	

# **Review: MIPS Register Usage**

R08-R15,	t0-t7	=	##	temps (not preserved by subroutines)
R24-R25,	t8-t9	=	##	temps (not preserved by subroutines)
			##	can be used without saving
R16-R23,	s0-s7	=	##	preserved by subroutines
			##	save before using,
			##	restore before return
R26-27,	k0-k1	=	##	reserved for interrupt handler
R28,	gp	=	##	global pointer
			##	(for easy access to some variables)
R29,	sp	=	##	stack pointer
R30,	s8/fp	=	##	9th subroutine reg / frame pointer
R31,	ra	=	##	return addr (used by jal)

CS350

Operating Systems

Fall 2013

	What is a Thread?
•	A thread represents the control state of an executing program.
)	A thread has an associated context (or state), which consists of
	<ul> <li>the processor's CPU state, including the values of the program counter (PC), the stack pointer, other registers, and the execution mode (privileged/non-privileged)</li> </ul>
	- a stack, which is located in the address space of the thread's process
	magine that you would like to suspend the program execution, and resume
i	t again later. Think of the thread context as the information you would
1	need in order to restart program execution from where it left off when it was suspended.





### **Example: Concurrent Mouse Simulations**

```
static void mouse_simulation(void * unusedpointer,
                              unsigned long mousenumber)
{
  int i; unsigned int bowl;
  for(i=0;i<NumLoops;i++) {</pre>
    /* for now, this mouse chooses a random bowl from
     * which to eat, and it is not synchronized with
     * other cats and mice
     */
    /* legal bowl numbers range from 1 to NumBowls */
    bowl = ((unsigned int)random() % NumBowls) + 1;
    mouse_eat(bowl);
  }
  /* indicate that this mouse is finished */
  V(CatMouseWait);
}
CS350
                            Operating Systems
```

```
<u>Implementing Threads</u>

a thread library is responsibile for implementing threads
the thread library stores threads' contexts (or pointers to the threads' contexts) when they are not running
the data structure used by the thread library to store a thread context is sometimes called a thread control block

In the OS/161 kernel's thread implementation, thread contexts are stored in thread structures.
```

7

Fall 2013



```
Threads and Concurrency
                                                               10
                   The OS/161 thread Structure
/* see kern/include/thread.h */
struct thread {
                            /* Name of this thread */
 char *t_name;
 const char *t_wchan_name; /* Wait channel name, if sleeping */
                            /* State this thread is in */
 threadstate_t t_state;
 /* Thread subsystem internal fields. */
 struct thread_machdep t_machdep; /* Any machine-dependent goo */
 struct threadlistnode t_listnode; /* run/sleep/zombie lists */
 void *t_stack;
                                 /* Kernel-level stack */
 struct switchframe *t_context; /* Register context (on stack) */
 struct cpu *t_cpu;
                                 /* CPU thread runs on */
 struct proc *t_proc;
                                /* Process thread belongs to */
 . . .
```





```
Dispatching on the MIPS (1 of 2)
```

```
/* See kern/arch/mips/thread/switch.S */
switchframe switch:
  /* a0: address of switchframe pointer of old thread. */
  /* al: address of switchframe pointer of new thread. */
   /* Allocate stack space for saving 10 registers. 10*4 = 40 */
   addi sp, sp, -40
   sw
        ra, 36(sp)
                     /* Save the registers */
        gp, 32(sp)
   sw
        s8, 28(sp)
   SW
        s6, 24(sp)
   sw
        s5, 20(sp)
   sw
        s4, 16(sp)
   sw
        s3, 12(sp)
   SW
   sw
        s2, 8(sp)
        s1, 4(sp)
   sw
        s0, 0(sp)
   sw
   /* Store the old stack pointer in the old thread */
        sp, 0(a0)
   sw
CS350
                            Operating Systems
                                                             Fall 2013
```

```
Threads and Concurrency
                                                                    14
                    Dispatching on the MIPS (2 of 2)
   /* Get the new stack pointer from the new thread */
   lw
        sp, 0(a1)
                  /* delay slot for load */
   nop
   /* Now, restore the registers */
   lw
        s0, 0(sp)
   lw
        s1, 4(sp)
        s2, 8(sp)
   lw
        s3, 12(sp)
   lw
   lw
        s4, 16(sp)
        s5, 20(sp)
   lw
        s6, 24(sp)
   lw
        s8, 28(sp)
   lw
   lw
        gp, 32(sp)
   lw
        ra, 36(sp)
                          /* delay slot for load */
   nop
   /* and return. */
   j ra
   addi sp, sp, 40
                          /* in delay slot */
   .end switchframe_switch
                                                                Fall 2013
```

```
CS350
```

# **Dispatching on the MIPS (Notes)**

- Not all of the registers are saved during a context switch
- This is because the context switch code is reached via a call to thread\_switch and by convention on the MIPS not all of the registers are required to be preserved across subroutine calls
- thus, after a call to switchframe\_switch returns, the caller (thread\_switch) does not expect all registers to have the same values as they had before the call - to save time, those registers are not preserved by the switch
- if the caller wants to reuse those registers it must save and restore them

CS350

**Operating Systems** 

Fall 2013

15

	Thread Library Interface	
• the thread library inte	rface allows program code to manipulate t	threads
• one key thread library	v interface function is <i>Yield()</i>	
• Yield() causes the call to choose some other causes a context swite	ling thread to stop and wait, and causes the waiting thread to run in its place. In other ch.	e thread library words, Yield()
<ul> <li>in addition to Yield services:</li> </ul>	(), thread libraries typically provide other	r thread-related
- create new thread		
– end (and destroy)	a thread	
- cause a thread to $b$	block (to be discussed later)	
		E 11 0010

# The OS/161 Thread Interface (incomplete)

```
19
```

### Creating Threads Using thread\_fork()

```
/* From kern/synchprobs/catmouse.c */
for (index = 0; index < NumMice; index++) {</pre>
  error = thread fork("mouse simulation thread",
    NULL, mouse simulation, NULL, index);
  if (error) {
    panic("mouse_simulation: thread_fork failed: %s\n",
     strerror(error));
  }
}
/* wait for all of the cats and mice to finish */
for(i=0;i<(NumCats+NumMice);i++) {</pre>
  P(CatMouseWait);
}
CS350
                          Operating Systems
                                                         Fall 2013
```

```
      Threads and Concurrency
      20

      Scheduling

      • scheduling means deciding which thread should run next

      • scheduling is implemented by a scheduler, which is part of the thread library

      • simple FIFO scheduling:

      - scheduler maintains a queue of threads, often called the ready queue

      - the first thread in the ready queue is the running thread

      - on a context switch the running thread is moved to the end of the ready queue, and new first thread is allowed to run

      - newly created threads are placed at the end of the ready queue

      • more on scheduling later ...
```

# Preemption

- Yield() allows programs to *voluntarily* pause their execution to allow another thread to run
- sometimes it is desirable to make a thread stop running even if it has not called Yield()
- this kind of *involuntary* context switch is called *preemption* of the running thread
- to implement preemption, the thread library must have a means of "getting control" (causing thread library code to be executed) even though the application has not called a thread library function
- this is normally accomplished using *interrupts*

CS350

Operating Systems

Fall 2013

nreads and Concurrency	22
Review: Interrupts	
• an interrupt is an event that occurs during the execu	ution of a program
• interrupts are caused by system devices (hardware), controller, a network interface	, e.g., a timer, a disk
• when an interrupt occurs, the hardware automatical location in memory	lly transfers control to a fixed
• at that memory location, the thread library places a <i>interrupt handler</i>	procedure called an
• the interrupt handler normally:	
1. saves the current thread context (in OS/161, this the current thread's stack)	s is saved in a <i>trap frame</i> on
2. determines which device caused the interrupt ar processing	nd performs device-specific
3. restores the saved thread context and resumes exwhere it left off at the time of the interrupt.	xecution in that context
S350 Operating Systems	Fall 2013

## **Round-Robin Scheduling**

- round-robin is one example of a preemptive scheduling policy
- round-robin scheduling is similar to FIFO scheduling, except that it is preemptive
- as in FIFO scheduling, there is a ready queue and the thread at the front of the ready queue runs
- unlike FIFO, a limit is placed on the amount of time that a thread can run before it is preempted
- the amount of time that a thread is allocated is called the scheduling *quantum*
- when the running thread's quantum expires, it is preempted and moved to the back of the ready queue. The thread at the front of the ready queue is dispatched and allowed to run.

CS350

Operating Systems

Fall 2013







	_
	Concurrency
• On multiproc processor.	essors, several threads can execute simultaneously, one on each
• On uniproces preemption an	sors, only one thread executes at a time. However, because of nd timesharing, threads appear to run concurrently.
Concurrency a	and synchronization are important even on uniprocessors.
Concurrency	and synchronization are important even on uniprocessors.
Concurrency	and synchronization are important even on uniprocessors.
Concurrency	and synchronization are important even on uniprocessors.
Concurrency	and synchronization are important even on uniprocessors.

	Thread Synchronization
• C	concurrent threads can interact with each other in a variety of ways:
-	<ul> <li>Threads share access, through the operating system, to system devices (more on this later)</li> </ul>
-	- Threads may share access to program data, e.g., global variables.
► A n v	common synchronization problem is to enforce <i>mutual exclusion</i> , which neans making sure that only one thread at a time uses a shared object, e.g., a ariable or a device.
• T se	The part of a program in which the shared object is accessed is called a <i>critical ection</i> .

# **Critical Section Example (Part 0)**

```
/* Note the use of volatile */
int volatile total = 0;
void add() {
                                void sub() {
   int i;
                                   int i;
   for (i=0; i<N; i++) {
                                   for (i=0; i<N; i++) {</pre>
     total++;
                                       total--;
   }
                                    }
}
                                }
```

If one thread executes add and another executes sub what is the value of total when they have finished?

CS350

**Operating Systems** 

Fall 2013

```
Synchronization
                                                              4
                  Critical Section Example (Part 0)
/* Note the use of volatile */
int volatile total = 0;
void add() {
                                void sub() {
   loadaddr R8 total
                                    loadaddr R10 total
   for (i=0; i<N; i++) {
                                   for (i=0; i<N; i++) {
      lw R9 0(R8)
                                        lw R11 0(R10)
      add R9 1
                                        sub R11 1
      sw R9 0(R8)
                                       sw R11 0(R10)
   }
                                    }
}
                                 }
                                                         Fall 2013
```

Synchronization 5 **Critical Section Example (Part 0)** Thread 1 Thread 2 loadaddr R8 total lw R9 0(R8) R9=0 add R9 1 R9=1 <INTERRUPT> loadaddr R10 total lw R11 0(R10) R11=0 sub R11 1 R11=-1 sw R11 0(R10) total=-1 <INTERRUPT> sw R9 0(R8) total=1 One possible order of execution. CS350 **Operating Systems** Fall 2013

	Critical S	ectio	n Example (Part 0)	
			Thread 2	
• R8	total			
R8)	R9=0			
	<interru< td=""><td>PT&gt;</td><td></td><td></td></interru<>	PT>		
			loadaddr R10 to	otal
			lw R11 0(R10)	R11=0
	<interru< td=""><td>PT&gt;</td><td></td><td></td></interru<>	PT>		
	R9=1			
R8)	total=1			
	<interru< td=""><td>PT&gt;</td><td></td><td></td></interru<>	PT>		
			sub R11 1	R11=-1
			sw R11 0(R10)	total=-1
	r R8 (R8) (R8)	Critical S R8 total (R8) R9=0 <interru (INTERRU R9=1 (R8) total=1 <interru< td=""><td>Critical Section R8 total (R8) R9=0 <interrupt> R9=1 (R8) total=1 <interrupt></interrupt></interrupt></td><td>Critical Section Example (Part 0) Thread 2 Thread 2 Thread 2 TR8 total TR8) R9=0 (INTERRUPT&gt; Ioadaddr R10 to Iw R11 0(R10) (INTERRUPT&gt; R9=1 TR8) total=1 (INTERRUPT&gt; Sub R11 1 Sw R11 0(R10)</td></interru<></interru 	Critical Section R8 total (R8) R9=0 <interrupt> R9=1 (R8) total=1 <interrupt></interrupt></interrupt>	Critical Section Example (Part 0) Thread 2 Thread 2 Thread 2 TR8 total TR8) R9=0 (INTERRUPT> Ioadaddr R10 to Iw R11 0(R10) (INTERRUPT> R9=1 TR8) total=1 (INTERRUPT> Sub R11 1 Sw R11 0(R10)

#### The use of volatile

```
/* What if we DO NOT use volatile */
int volatile total = 0;
void add() {
                                void sub() {
   loadaddr R8 total
                                    loadaddr R10 total
   lw R9 0(R8)
                                    lw R11 0(R10)
   for (i=0; i<N; i++) {</pre>
                                   for (i=0; i<N; i++) {
      add R9 1
                                       sub R11 1
   }
                                    }
   sw R9 0(R8)
                                    sw R11 0(R10)
                                }
}
```

Without volatile the compiler could optimize the code. If one thread executes add and another executes sub, what is the value of total when they have finished?

CS350

Operating Systems

Fall 2013

7

```
Synchronization
                                                               8
                        The use of volatile
/* What if we DO NOT use volatile */
int volatile total = 0;
void add() {
                                void sub() {
   loadaddr R8 total
                                    loadaddr R10 total
   lw R9 0(R8)
                                    lw R11 0(R10)
   add R9 N
                                    sub R11 N
   sw R9 0(R8)
                                    sw R11 0(R10)
                                }
}
```

The compiler could aggressively optimize the code., Volatile tells the compiler that the object may change for reasons which cannot be determined from the local code (e.g., due to interaction with a device or because of another thread).

CS350

#### The use of volatile

```
/* Note the use of volatile */
int volatile total = 0;
void add() {
                                 void sub() {
   loadaddr R8 total
                                    loadaddr R10 total
   for (i=0; i<N; i++) {</pre>
                                    for (i=0; i<N; i++) {</pre>
      lw R9 0(R8)
                                        lw R11 0(R10)
      add R9 1
                                        sub R11 1
      sw R9 0(R8)
                                        sw R11 0(R10)
   }
                                    }
}
                                 }
```

The volatile declaration forces the compiler to load and store the value on every use. Using volatile is necessary but not sufficient for correct behaviour. Mutual exclusion is also required to ensure a correct ordering of instructions.

CS350

Operating Systems

Fall 2013

```
10
Synchronization
              Ensuring Correctness with Multiple Threads
/* Note the use of volatile */
int volatile total = 0;
void add() {
                                    void sub() {
   int i;
                                        int i;
   for (i=0; i<N; i++) {</pre>
                                        for (i=0; i<N; i++) {</pre>
     Allow one thread to to execute and make others wait
         total++;
                                               total--;
     Permit one waiting thread to continue execution
   }
                                        }
}
                                    }
   Threads must enforce mutual exclusion.
```

# **Critical Section Example (Part 1)**

```
int list_remove_front(list *lp) {
    int num;
    list_element *element;
    assert(!is_empty(lp));
    element = lp->first;
    num = lp->first->item;
    if (lp->first == lp->last) {
        lp->first = lp->last = NULL;
    } else {
        lp->first = element->next;
    }
    lp->num_in_list--;
    free(element);
    return num;
}
```

The list\_remove\_front function is a critical section. It may not work properly if two threads call it at the same time on the same list. (Why?)

CS350

Operating Systems

Fall 2013

11

```
Synchronization
                                                                12
                  Critical Section Example (Part 2)
void list_append(list *lp, int new_item) {
   list_element *element = malloc(sizeof(list_element));
   element->item = new_item
   assert(!is_in_list(lp, new_item));
   if (is_empty(lp)) {
     lp->first = element; lp->last = element;
   } else {
     lp->last->next = element; lp->last = element;
   ł
   lp->num_in_list++;
}
   The list_append function is part of the same critical section as
   list_remove_front. It may not work properly if two threads call
   it at the same time, or if a thread calls it while another has called
   list_remove_front
```

CS350



<ul> <li>On a uniprocessor, only one thread at a time is actually running.</li> <li>If the running thread is executing a critical section, mutual exclusion may be violated if <ol> <li>the running thread is preempted (or voluntarily yields) while it is in the critical section, and</li> <li>the scheduler chooses a different thread to run, and this new thread enter the same critical section that the preempted thread was in</li> </ol> </li> </ul>		Disabling Interrupts
<ul> <li>If the running thread is executing a critical section, mutual exclusion may be violated if</li> <li>1. the running thread is preempted (or voluntarily yields) while it is in the critical section, and</li> <li>2. the scheduler chooses a different thread to run, and this new thread enter the same critical section that the preempted thread was in</li> </ul>	On a uniprocessor, or	nly one thread at a time is actually running.
<ol> <li>the running thread is preempted (or voluntarily yields) while it is in the critical section, and</li> <li>the scheduler chooses a different thread to run, and this new thread enter the same critical section that the preempted thread was in</li> </ol>	If the running thread violated if	is executing a critical section, mutual exclusion may be
2. the scheduler chooses a different thread to run, and this new thread enter the same critical section that the preempted thread was in	1. the running thread critical section, an	l is preempted (or voluntarily yields) while it is in the nd
	2. the scheduler choo the same critical s	oses a different thread to run, and this new thread enters section that the preempted thread was in
• Since preemption is caused by timer interrupts, mutual exclusion can be enforced by disabling timer interrupts before a thread enters the critical section and re-enabling them when the thread leaves the critical section.	Since preemption is c enforced by disabling and re-enabling them	caused by timer interrupts, mutual exclusion can be g timer interrupts before a thread enters the critical section when the thread leaves the critical section.

#### **Interrupts in OS/161**

This is one way that the OS/161 kernel enforces mutual exclusion on a single processor. There is a simple interface

- spl0() sets IPL to 0, enabling all interrupts.
- splhigh() sets IPL to the highest value, disabling all interrupts.
- splx(s) sets IPL to S, enabling whatever state S represents.

These are used by splx() and by the spinlock code.

- splraise(int oldipl, int newipl)
- spllower(int oldipl, int newipl)
- For splraise, NEWIPL > OLDIPL, and for spllower, NEWIPL < OLDIPL.

See kern/include/spl.h and kern/thread/spl.c

CS350

Operating Systems

Fall 2013

	<b>Pros and Cons of Disabling Interrupts</b>
• advantages	
– does no	ot require any hardware-specific synchronization instructions
– works f	or any number of concurrent threads
• disadvanta	ges:
– indiscri threater	minate: prevents all preemption, not just preemption that would n the critical section
<ul> <li>ignoring</li> <li>time. (Vinterrup</li> </ul>	g timer interrupts has side effects, e.g., kernel unaware of passage of Worse, OS/161's splhigh() disables <i>all</i> interrupts, not just timer ots.) Keep critical sections <i>short</i> to minimize these problems.
– wiii iio	temorce mutual exclusion on multiprocessors (why??)

# Peterson's Mutual Exclusion Algorithm

```
/* shared variables */
/* note: one flag array and turn variable */
/* for each critical section */
boolean volatile flag[2]; /* shared, initially false */
int volatile turn;
                             /* shared */
flag[i] = true; /* for one thread, i = 0 and j = 1 * /
                   /* for the other, i=1 and j=0 */
turn = j;
while (flag[j] && turn == j) { } /* busy wait */
 critical section
                      /* e.g., call to list_remove_front */
flag[i] = false;
  Ensures mutual exclusion and avoids starvation, but works only for two
  threads. (Why?)
CS350
                         Operating Systems
                                                        Fall 2013
```

Synch	Hardware-Specific Synchronization Instructions
•	a test-and-set instruction <i>atomically</i> sets the value of a specified memory location and either
	- places that memory location's <i>old</i> value into a register, or
	<ul> <li>checks a condition against the memory location's old value and records the result of the check in a register</li> </ul>
•	for presentation purposes, we will abstract such an instruction as a function TestAndSet(address,value), which takes a memory location (address) and a value as parameters. It atomically stores value at the memory location specified by address and returns the previous value stored at that address
•	Often only two values are used 0 and 1 so the value parameter is not used and a value of 1 is implied (e.g., in OS/161)

# A Spin Lock Using Test-And-Set in OS/161

- a test-and-set instruction can be used to enforce mutual exclusion
- for each critical section, define a shared variable

```
volatile spinlock_data_t lk_lock; /* initially 0 */
```

We will use the lock variable to keep track of whether there is a thread in the critical section, in which case the value of lk\_lock will be 1

• before a thread can enter the critical section, it does the following:

```
while (spinlock_data_testandset(&lk->lk_lock) != 0) {
   /* busy wait */
}
```

- if lk\_lock == 0 then it is set to 1 and the thread enters the critical section
- when the thread leaves the critical section, it does:

```
spinlock_data_set(&lk->lk_lock, 0);
```

CS350

Operating Systems

Fall 2013

	A Spin Lock Using Test-And-Set
this enforce	s mutual exclusion (why?), but starvation is a possibility
This constru	ict is sometimes known as a <i>spin lock</i> , since a thread "spins" in
multiproces	sors.

#### Spinlocks in OS/161

```
struct spinlock {
   volatile spinlock_data_t lk_lock; /* word for spin */
   struct cpu *lk_holder; /* CPU holding this lock */
};
void spinlock_init(struct spinlock *lk);
void spinlock_cleanup(struct spinlock *lk);
void spinlock_acquire(struct spinlock *lk);
void spinlock_release(struct spinlock *lk);
bool spinlock_do_i_hold(struct spinlock *lk);
```

CS350

Operating Systems

Fall 2013

```
Synchronization
                                                              22
                       Spinlocks in OS/161
spinlock_init(struct spinlock *lk)
{
  spinlock_data_set(&lk->lk_lock, 0);
  lk->lk_holder = NULL;
}
void spinlock_cleanup(struct spinlock *lk)
ł
  KASSERT(lk->lk_holder == NULL);
  KASSERT(spinlock_data_get(&lk->lk_lock) == 0);
}
void spinlock_data_set(volatile spinlock_data_t *sd,
  unsigned val)
{
  *sd = val;
CS350
                                                          Fall 2013
                          Operating Systems
```

### Spinlocks in OS/161

```
void spinlock_acquire(struct spinlock *lk)
ł
  struct cpu *mycpu;
  splraise(IPL_NONE, IPL_HIGH);
  /* this must work before curcpu initialization */
  if (CURCPU_EXISTS()) {
    mycpu = curcpu->c_self;
    if (lk->lk_holder == mycpu) {
      panic("Deadlock on spinlock %p\n", lk);
    }
  } else {
    mycpu = NULL;
  }
CS350
                         Operating Systems
                                                        Fall 2013
```

```
Synchronization 24
Spinlocks in OS/161
while (1) {
    /* Do test-test-and-set to reduce bus contention */
    if (spinlock_data_get(&lk->lk_lock) != 0) {
        continue;
        }
        if (spinlock_data_testandset(&lk->lk_lock) != 0) {
        continue;
        }
        break;
    }
    lk->lk_holder = mycpu;
}
```

#### Spinlocks in OS/161

```
void spinlock_release(struct spinlock *lk)
{
  /* this must work before curcpu initialization */
  if (CURCPU_EXISTS()) {
    KASSERT(lk->lk_holder == curcpu->c_self);
  }
  lk->lk holder = NULL;
  spinlock_data_set(&lk->lk_lock, 0);
  spllower(IPL_HIGH, IPL_NONE);
}
                                                        Fall 2013
CS350
                          Operating Systems
```

```
Synchronization
                                                                   26
                     Load-Link / Store-Conditional
Load-link returns the current value of a memory location, while a subsequent
store-conditional to the same memory location will store a new value only if no
updates have occurred to that location since the load-link.
spinlock_data_testandset(volatile spinlock_data_t *sd)
{
  spinlock_data_t x,y;
  /* Test-and-set using LL/SC.
   * Load the existing value into X, and use Y to store 1.
   * After the SC, Y contains 1 if the store succeeded,
   * 0 if it failed. On failure, return 1 to pretend
    * that the spinlock was already held.
    */
  y = 1;
                                                               Fall 2013
```

# Load-Link / Store-Conditional

CS350	C	Operating Systems Fall 2013
,		
}		
return x;		
}		
return 1;		
if (y == 0) {		
: "=r" (x), "+r"	(y)	: "r" (sd));
".set pop"	/ *	restore assembler mode */
"sc %1, 0(%2);"	/ *	*sd = y; y = success? */
"ll %0, 0(%2);"	/ *	x = *sd */
".set volatile;"	/ *	avoid unwanted optimization */
".set mips32;"	/ *	allow MIPS32 instructions */
".set push;"	/ *	save assembler mode */
asm volatile(		

Pros and Cons of Spinlocks
Pros:
<ul> <li>– can be efficient for short critical sections</li> </ul>
<ul> <li>using hardware specific synchronization instructions means it works on multiprocessors</li> </ul>
Cons:
- CPU is busy (nothing else runs) while waiting for lock
- starvation is possible
If critical section is short prefer spinlock.
If critical section is long prefer blocking lock.
Hybrid locks will spin for a period of time before blocking.
Question: How to determine how long to spin for hybrid lock?

CS350

Fall 2013

	Semaphores	
<ul> <li>A sema exclusion synchrophysical exclusion of the syn</li></ul>	phore is a synchronization primitive that can be used to enforce non requirements. It can also be used to solve other kinds of nization problems.	nutual
• A sema operation	phore is an object that has an integer value, and that supports two	)
<b>P:</b> if th wait	e semaphore value is greater than 0, decrement the value. Otherv until the value is greater than 0 and then decrement it.	vise,
V: incr	ement the value of the semaphore	
• Two kin	ids of semaphores:	
countir	g semaphores: can take on any non-negative value	
<b>binary</b> sem	<b>semaphores:</b> take on only the values 0 and 1. (V on a binary aphore with value 1 has no effect.)	
By defin	ition, the $P$ and $V$ operations of a semaphore are <i>atomic</i> .	
	Operating Systems	Eall 20

A Simple Example	e using Semaphores
void add() {	<pre>void sub() {</pre>
int i;	int i;
for (i=0; i <n; i++)="" td="" {<=""><td>for (i=0; i<n; i++)="" td="" {<=""></n;></td></n;>	for (i=0; i <n; i++)="" td="" {<=""></n;>
P(sem);	P(sem);
total++;	total;
V(sem);	V(sem);
}	}
}	}
What type of semaphore can be used	for sem?
CS250 Onorstin	g Systems Fall 201

```
Synchronization
                                                              31
                       OS/161 Semaphores
struct semaphore {
  char *sem_name;
  struct wchan *sem_wchan;
  struct spinlock sem_lock;
  volatile int sem_count;
};
struct semaphore *sem_create(const char *name,
  int initial_count);
void P(struct semaphore *s);
void V(struct semaphore *s);
void sem_destroy(struct semaphore *s);
   see kern/include/synch.h and kern/thread/synch.c
CS350
                          Operating Systems
                                                          Fall 2013
```

```
Synchronization y 32
Synchronization y 1
Struct semaphore *s;
s = sem_create("MySem1", 1); /* initial value is 1 */
P(s); /* do this before entering critical section */
critical section /* e.g., call to list_remove_front */
V(s); /* do this after leaving critical section */
```

```
OS/161 Semaphores: P() from kern/thread/synch.c
P(struct semaphore *sem)
ł
  KASSERT(sem != NULL);
  KASSERT(curthread->t_in_interrupt == false);
  spinlock_acquire(&sem->sem_lock);
    while (sem->sem_count == 0) {
      /* Note: we don't maintain strict FIFO ordering */
      wchan_lock(sem->sem_wchan);
      spinlock_release(&sem->sem_lock);
      wchan_sleep(sem->sem_wchan);
      spinlock_acquire(&sem->sem_lock);
    }
    KASSERT(sem->sem_count > 0);
    sem->sem_count--;
  spinlock_release(&sem->sem_lock);
}
CS350
                         Operating Systems
                                                        Fall 2013
```

```
Synchronization 34
OS/161 Semaphores: V() from kern/thread/synch.c
V(struct semaphore *sem)
{
KASSERT(sem != NULL);
spinlock_acquire(&sem->sem_lock);
sem->sem_count++;
KASSERT(sem->sem_count > 0);
wchan_wakeone(sem->sem_wchan);
spinlock_release(&sem->sem_lock);
}
```

	Thread Blocking
<ul> <li>Sometimes previous sl semaphore</li> </ul>	a thread will need to wait for an event. One example is on the ide: a thread that attempts a $P()$ operation on a zero-valued must wait until the semaphore's value becomes positive.
• other exam	ples that we will see later on:
– wait fo	r data from a (relatively) slow device
– wait fo	r input from a keyboard
– wait fo	r busy device to become idle
• In these cir anything u	cumstances, we do not want the thread to run, since it cannot do seful.
• To handle	this, the thread scheduler can <i>block</i> threads.

CS350

Operating Systems

Fall 2013
## **Thread Blocking in OS/161**

- wchan\_sleep() is much like thread\_yield(). The calling thread is voluntarily giving up the CPU, so the scheduler chooses a new thread to run, the state of the running thread is saved and the new thread is dispatched. However:
  - after a thread\_yield(), the calling thread is *ready* to run again as soon as it is chosen by the scheduler
  - after a wchan\_sleep(), the calling thread is *blocked*, and must not be scheduled to run again until after it has been explicitly unblocked by a call to wchan\_wakeone() or wchan\_wakeall().

CS350

Operating Systems

Fall 2013



Pr	oducer/Consumer Synchronization	
• suppose we have thr remove items from t	reads that add items to a list (producers) a he list (consumers)	and threads that
• suppose we want to instead they must wa	ensure that consumers do not consume i ait until the list has something in it	f the list is empty -
• this requires synchro	onization between consumers and produc	cers
<ul> <li>semaphores can pro- slide</li> </ul>	vide the necessary synchronization, as sh	own on the next

Producer/Consumer Synchronization using Semaphores	
struct semaphore *s; s = sem_create("Items", 0); /* initial value is 0 ;	- /
Producer's Pseudo-code: add item to the list (call list_append()) V(s);	
Consumer's Pseudo-code:	
remove item from the list (call list_remove_front	())



Synchronization

## OS/161 Locks

• OS/161 also uses a synchronization primitive called a *lock*. Locks are intended to be used to enforce mutual exclusion.

```
struct lock *mylock = lock_create("LockName");
```

lock\_aquire(mylock);

```
critical section /* e.g., call to list_remove_front */
lock_release(mylock);
```

- A lock is similar to a binary semaphore with an initial value of 1. However, locks also enforce an additional constraint: the thread that releases a lock must be the same thread that most recently acquired it.
- The system enforces this additional constraint to help ensure that locks are used as intended.

CS350

Operating Systems

Fall 2013

```
§ynchrinter jet in either of read (shared) or
write (exclusive) mode
• In OS/161 reader/writer locks might look like this:
struct rwlock *rwlock = rw_lock_create("RWLock");
rwlock_aquire(rwlock, READ_MODE);
can only read shared resources
/* access is shared by readers */
rwlock_release(rwlock);
rwlock_aquire(rwlock, WRITE_MODE);
can read and write shared resources
/* access is exclusive to only one writer */
rwlock_release(rwlock);
```

·	Critical Section Requirements	
• Mutual exclusi thread can exec	<b>on</b> : While one thread is executing in the critic ute in that critical section.	cal section no othe
• <b>Progress</b> : The t section.	hread in the critical section will eventually le	ave the critical
• Bounded waiti	<b>ng</b> : Any thread will wait for a bounded amou the critical section.	nt of time before it

	Performance Issues
Ð	<b>Overhead</b> : the memory and CPU resources used when acquiring and releasing access to critical sections
	Contention: competition for access to the critical section
	Granularity: the amount of code executed while in a critical section
 	Why are these important issues?

Synchronization 47 Lock Overhead, Contention and Granularity (Option 1) void add() { void sub() { int i; int i; for (i=0; i<N; i++) {</pre> for (i=0; i<N; i++) {</pre> P / Acquire P / Acquire total++; total--; V / Release V / Release } } } } Should one use P()/V(), spinlock\_acquire()/spinlock\_release() or lock\_acquire()/lock\_release? CS350 **Operating Systems** Fall 2013

```
48
Synchronization
          Lock Overhead, Contention and Granularity (Option 2)
void add() {
                                     void sub() {
   int i;
                                         int i;
                                         P / Acquire
   P / Acquire
      for (i=0; i<N; i++) {</pre>
                                           for (i=0; i<N; i++) {
           total++;
                                                total--;
      }
                                           }
   V / Release
                                         V / Release
}
                                     }
```

Which option is better Option 1 (previous slide) or 2 (this slide)? Why? Does the choice of where to do synchronization influence the choice of which mechanism to use for synchronization?

	<b>Condition Variables</b>	
• OS/161 supports <i>variables</i>	another common synchronization primitiv	ve: condition
• each condition v variables are onl <i>lock</i>	ariable is intended to work together with a y used <i>from within the critical section that</i>	lock: condition <i>is protected by the</i>
• three operations	are possible on a condition variable:	
wait: This cause with the conc lock.	es the calling thread to block, and it release lition variable. Once the thread is unblock	es the lock associated ed it reacquires the
signal: If thread those threads	ls are blocked on the signaled condition va is unblocked.	riable, then one of
<b>broadcast:</b> Like condition var	e signal, but unblocks all threads that are b iable.	locked on the

ynchronization	50	
Using Condition Variables		
• Condition variables get their name because they allow threads to w arbitrary conditions to become true inside of a critical section.	ait for	
• Normally, each condition variable corresponds to a particular cond of interest to an application. For example, in the bounded buffer producer/consumer example on the following slides, the two condi	ition that is tions are:	
<ul> <li><math>count &gt; 0</math> (condition variable notempty)</li> <li><math>count &lt; N</math> (condition variable notfull)</li> </ul>		
• when a condition is not true, a thread can wait on the correspondit variable until it becomes true	ing condition	
• when a thread detects that a condition is true, it uses signal or b to notify any threads that may be waiting	roadcast	
Note that signalling (or broadcasting to) a condition variable that waiters has <i>no effect</i> . Signals do not accumulate.	at has no	

#### Waiting on Condition Variables

- when a blocked thread is unblocked (by signal or broadcast), it reacquires the lock before returning from the wait call
- a thread is in the critical section when it calls wait, and it will be in the critical section when wait returns. However, in between the call and the return, while the caller is blocked, the caller is out of the critical section, and other threads may enter.
- In particular, the thread that calls signal (or broadcast) to wake up the waiting thread will itself be in the critical section when it signals. The waiting thread will have to wait (at least) until the signaller releases the lock before it can unblock and return from the wait call.

This describes Mesa-style condition variables, which are used in OS/161. There are alternative condition variable semantics (Hoare semantics), which differ from the semantics described here.

CS350

Operating Systems

Fall 2013

```
Synchronization
                                                           52
          Bounded Buffer Producer Using Condition Variables
                           /* must initially be 0 */
int volatile count = 0;
struct lock *mutex;
                         /* for mutual exclusion */
struct cv *notfull, *notempty; /* condition variables */
/* Initialization Note: the lock and cv's must be created
 * using lock_create() and cv_create() before Produce()
 * and Consume() are called */
Produce(itemType item) {
  lock_acquire(mutex);
  while (count == N) {
     cv_wait(notfull, mutex);
  }
  add item to buffer (call list_append())
  count = count + 1;
  cv_signal(notempty, mutex);
  lock_release(mutex);
}
```

Synchronization

# **Bounded Buffer Consumer Using Condition Variables**

```
itemType Consume() {
  lock_acquire(mutex);
  while (count == 0) {
     cv_wait(notempty, mutex);
  }
  remove item from buffer (call list_remove_front())
  count = count - 1;
  cv_signal(notfull, mutex);
  lock_release(mutex);
  return(item);
```

Both Produce() and Consume() call cv\_wait() inside of a while loop. Why?

CS350

}

**Operating Systems** 

Fall 2013



#### Monitors in OS/161

- The C language, in which OS/161 is written, does not support monitors.
- However, programming convention and OS/161 locks and condition variables can be used to provide monitor-like behavior for shared kernel data structures:
  - define a C structure to implement the object's data elements
  - define a set of C functions to manipulate that structure (these are the object "methods")
  - ensure that only those functions directly manipulate the structure
  - create an OS/161 lock to enforce mutual exclusion
  - ensure that each access method acquires the lock when it starts and releases the lock when it finishes
  - if desired, define one or more condition variables and use them within the methods.

CS350

Operating Systems

Fall 2013

Deadlocks		
Suppose there are two threads and two locks, lounlocked.	ckA and lockB, both initially	
Suppose the following sequence of events occurs	5	
1. Thread 1 does lock_acquire(lockA).		
2. Thread 2 does lock_acquire(lockB).		
3. Thread 1 does lock_acquire(lockB) and held by thread 2.	nd blocks, because lockB is	
<ol> <li>Thread 2 does lock_acquire(lockA) and held by thread 1.</li> </ol>	nd blocks, because lockA is	
These two threads are <i>deadlocked</i> - neither threating will not resolve the deadlock. The threads are	d can make progress. Wait-	

	Deadlocks (Another Simple Example)	
• Sup	pose a machine has 64 MB of memory. The following sequence ours.	of events
1.	Thread A starts, requests $30 \text{ MB}$ of memory.	
2.	Thread $B$ starts, also requests 30 MB of memory.	
3.	Thread $A$ requests an additional 8 MB of memory. The kernel blo $A$ since there is only 4 MB of available memory.	cks thread
4.	Thread $B$ requests an additional 5 MB of memory. The kernel blo $B$ since there is not enough memory available.	cks threa
The	se two threads are deadlocked.	
\$350	Operating Systems	Fall 20







## **Deadlock Detection and Recovery**

- main idea: the system maintains the resource allocation graph and tests it to determine whether there is a deadlock. If there is, the system must recover from the deadlock situation.
- deadlock recovery is usually accomplished by terminating one or more of the threads involved in the deadlock
- when to test for deadlocks? Can test on every blocked resource request, or can simply test periodically. Deadlocks persist, so periodic detection will not "miss" them.

Deadlock detection and deadlock recovery are both costly. This approach makes sense only if deadlocks are expected to be infrequent.

CS350

**Operating Systems** 

Detecting Deadlock in a Resource Allocation Graph		
• System State Notation		
– $D_i$ : demand vector	r for thread $T_i$	
– $A_i$ : current allocat	ion vector for thread $T_i$	
– $U$ : unallocated (av	ailable) resource vector	
• Additional Algorithm	Notation:	
- R: scratch resource	e vector	
– $f_i$ : algorithm is fin	ished with thread $T_i$ ? (boolean)	
8250	One section - Sectores	E-11 201/

Synchronization

63

```
/* initialization */
R = U
for all i, f_i = \text{false}
/* can each thread finish? */
while \exists i (\neg f_i \land (D_i \leq R)) \{
R = R + A_i
f_i = \text{true}
}
/* if not, there is a deadlock */
if \exists i (\neg f_i) then report deadlock
else report no deadlock
```

CS350

Operating Systems

**Detecting Deadlock (cont'd)** 

Fall 2013







## What is a Process?

Answer 1: a process is an abstraction of a program in execution

Answer 2: a process consists of

- an *address space*, which represents the memory that holds the program's code and data
- a *thread* of execution (possibly several threads)
- other resources associated with the running program. For example:
  - open files
  - sockets
  - attributes, such as a name (process identifier)
  - ...

A process with one thread is a *sequential* process. A process with more than one thread is a *concurrent* process.

CS350

**Operating Systems** 

Fall 2013

Multiprogramming		
multiprog	gramming means having multiple processes existing at the same time	
most mod	lern, general purpose operating systems support multiprogramming	
all proces	ses share the available hardware resources, with the sharing ed by the operating system:	
– Each The C	process uses some of the available memory to hold its address space. OS decides which memory and how much memory each process gets	
<ul> <li>OS ca proce</li> </ul>	In coordinate shared access to devices (keyboards, disks), since sses use these devices indirectly, by making system calls.	
<ul> <li>Proce</li> <li>the op</li> </ul>	sses <i>timeshare</i> the processor(s). Again, timesharing is controlled by perating system.	
OS ensur	es that processes are isolated from one another. Interprocess	
communi processes	cation should be possible, but only at the explicit request of the involved.	

	The OS Kernel	
• The kernel is a pr	ogram. It has code and data like any other	program.
• Usually kernel co do not	de runs in a privileged execution mode, w	hile other programs
C\$350	Operating Systems	Fall 201



	Kernel Privilege, Kernel Protection	
• What does it mea	an to run in privileged mode?	
• Kernel uses privi	lege to	
<ul> <li>– control hardw</li> </ul>	vare	
<ul> <li>protect and is</li> </ul>	olate itself from processes	
• privileges vary fr	om platform to platform, but may include:	
– ability to exec	cute special instructions (like halt)	
– ability to mar	nipulate processor state (like execution mod	le)
– ability to acco	ess memory addresses that can't be accesse	d otherwise
• kernel ensures th	at it is <i>isolated</i> from processes. No process	can execute or
change kernel co mechanisms like	de, or read or write kernel data, except thro system calls.	ough controlled
<u></u>	Operating Systems	Fall 201

	System Calls
•	System calls are an interface between processes and the kernel.
•	A process uses system calls to request operating system services.
•	From point of view of the process, these services are used to manipulate the
	abstractions that are part of its execution environment. For example, a process might use a system call to
	– open a file
	– send a message over a pipe
	– create another process
	- increase the size of its address space

Processes	and	the	Kernel
-----------	-----	-----	--------

# How System Calls Work

- The hardware provides a mechanism that a running program can use to cause a system call. Often, it is a special instruction, e.g., the MIPS syscall instruction.
- What happens on a system call:
  - the processor is switched to system (privileged) execution mode
  - key parts of the current thread context, such as the program counter, are saved

- the program counter is set to a fixed (determined by the hardware) memory address, which is within the kernel's address space

CS350

Operating Systems

Fall 2013

	System Call Execution and Return
•	Once a system call occurs, the calling thread will be executing a system call handler, which is part of the kernel, in privileged mode.
•	The kernel's handler determines which service the calling process wanted, and performs that service.
)	When the kernel is finished, it returns from the system call. This means:
	<ul> <li>restore the key parts of the thread context that were saved when the system call was made</li> </ul>
	- switch the processor back to unprivileged (user) execution mode
)	Now the thread is executing the calling process' program again, picking up where it left off when it made the system call.
	A system call causes a thread to stop executing application code and to start executing kernel code in privileged mode. The system call return switches the thread back to executing application code in unprivileged mode.



OS/161 close System Call Descript	tion
Library: standard C library (libc)	
Synopsis:	
<pre>#include <unistd.h> int close(int fd);</unistd.h></pre>	
<b>Description:</b> The file handle fd is closed	
<b>Return Values:</b> On success, close returns 0. On error, -1 is set according to the error encountered.	is returned and errno
Errors:	
<b>EBADF:</b> fd is not a valid file handle	
<b>EIO:</b> A hard I/O error occurred	

```
An Example System Call: A Tiny OS/161 Application that Uses close
/* Program: user/uw-testbin/syscall.c */
#include <unistd.h>
#include <errno.h>
int
main()
{
  int x;
  x = close(999);
  if (x < 0) {
    return errno;
  }
  return x;
}
CS350
                            Operating Systems
```

Processes and the Kernel 12 Disassembly listing of user/uw-testbin/syscall 00400050 <main>: 400050: 27bdffe8 addiu sp, sp, -24 400054: afbf0010 sw ra,16(sp) 400058: 0c100077 jal 4001dc <close> 40005c: 240403e7 li a0,999 400060: 04410003 bgez v0,400070 <main+0x20> 400064: 00000000 nop 400068: 3c021000 lui v0,0x1000 40006c: 8c420000 lw v0, 0(v0)400070: 8fbf0010 ra,16(sp) lw 400074: 00000000 nop 400078: 03e00008 jr ra 40007c: 27bd0018 addiu sp, sp, 24 The above can be obtained by disassembling the compiled syscall executable file with cs350-objdump -d

11

# System Call Wrapper Functions from the Standard Library

```
...
004001dc <close>:
    4001dc: 08100030    j 4000c0 <___syscall>
    4001e0: 24020031    li    v0,49
004001e4 <read>:
    4001e4: 08100030    j 4000c0 <___syscall>
    4001e8: 24020032    li    v0,50
...
```

The above is disassembled code from the standard C library (libc), which is linked with user/uw-testbin/syscall.o.

CS350

Operating Systems

Fall 2013



#### **OS/161 System Call Code Definitions**

/* Contains a number for	ever	y more-or-less standard 😽	<b>k</b> /
/* Unix system call (you	will	implement some subset). *	r /
#define SYS_close	49		
#define SYS_read	50		
#define SYS_pread	51		
//#define SYS_readv	52	/* won't be implementing *	<b>k</b> /
//#define SYS_preadv	53	/* won't be implementing *	<b>k</b> /
#define SYS_getdirentry	54		
#define SYS_write	55		

This comes from kern/include/kern/syscall.h. The files in kern/include/kern define things (like system call codes) that must be known by both the kernel and applications.

CS350

Operating Systems

Fall 2013

```
Processes and the Kernel
                                                                   16
             The OS/161 System Call and Return Processing
004000c0 <___syscall>:
  4000c0: 000000c
                        syscall
  4000c4: 10e00005
                               a3,4000dc <__syscall+0x1c>
                        beqz
  4000c8: 0000000
                        nop
  4000cc: 3c011000
                        lui at,0x1000
  4000d0: ac220000
                             v0,0(at)
                        sw
  4000d4: 2403ffff
                             v1,-1
                        li
  4000d8: 2402ffff
                        li
                             v0,-1
  4000dc: 03e00008
                        jr
                             ra
  4000e0: 00000000
                        nop
   The system call and return processing, from the standard C library. Like the
   rest of the library, this is unprivileged, user-level code.
```

#### **OS/161 MIPS Exception Handler**

```
common_exception:
```

```
/* Coming from user mode - find kernel stack */
mfc0 kl, c0_context /* we keep the CPU number here */
srl kl, kl, CTX_PTBASESHIFT /* shift to get the CPU number */
sll kl, kl, 2 /* shift back to make array index */
lui k0, %hi(cpustacks) /* get base address of cpustacks[] */
addu k0, k0, kl /* index it */
move kl, sp /* Save previous stack pointer */
b 2f /* Skip to common code */
lw sp, %lo(cpustacks)(k0) /* Load kernel sp (in delay slot) */
```

CS350

Operating Systems

```
Fall 2013
```

```
Processes and the Kernel
                                                                  18
                   OS/161 MIPS Exception Handler
1:
  /* Coming from kernel mode - just save previous stuff */
 move k1, sp
                  /* Save previous stack in k1 (delay slot) */
2:
  /* At this point:
   * Interrupts are off. (The processor did this for us.)
   * k0 contains the value for curthread, to go into s7.
   * k1 contains the old stack pointer.
   * sp points into the kernel stack.
   * All other registers are untouched.
   */
   When the syscall instruction occurs, the MIPS transfers control to ad-
   dress 0x8000080.
                       This kernel exception handler lives there.
                                                              See
   kern/arch/mips/locore/exception-mips1.S
```



	OS/161 MIPS Exception Handler (cont'd)
The	e common_exception code then does the following:
1.	allocates a <i>trap frame</i> on the thread's kernel stack and saves the user-level application's complete processor state (all registers except k0 and k1) into the trap frame.
2.	calls the mips_trap function to continue processing the exception.
3.	when mips_trap returns, restores the application processor state from the trap frame to the registers
4.	issues MIPS jr and rfe (restore from exception) instructions to return control to the application code. The jr instruction takes control back to the location specified by the application program counter when the syscall occurred (i.e., exception PC) and the rfe (which happens in the delay slot of the jr) restores the processor to unprivileged mode





ł

#### **OS/161 System Call Handler**

```
syscall(struct trapframe *tf)
   callno = tf->tf_v0; retval = 0;
   switch (callno) {
     case SYS_reboot:
       err = sys_reboot(tf->tf_a0);
       break;
     case SYS
               _time:
       err = sys___time((userptr_t)tf->tf_a0,
         (userptr_t)tf->tf_a1);
       break;
     /* Add stuff here */
     default:
       kprintf("Unknown syscall %d\n", callno);
       err = ENOSYS;
       break;
   }
```

syscall checks system call code and invokes a handler for the indicated system call. See kern/arch/mips/syscall/syscall.c

CS350

**Operating Systems** 

```
Fall 2013
```

23

```
Processes and the Kernel
                                                                     24
               OS/161 MIPS System Call Return Handling
  if (err) {
    tf->tf_v0 = err;
    tf -> tf_a3 = 1;
                          /* signal an error */
  } else {
    /* Success. */
    tf->tf v0 = retval;
    tf \rightarrow tf_a3 = 0;
                          /* signal no error */
  }
  /* Advance the PC, to avoid the syscall again. */
  tf \rightarrow tf_epc += 4;
  /* Make sure the syscall code didn't forget to lower spl */
  KASSERT(curthread->t_curspl == 0);
  /* ...or leak any spinlocks */
  KASSERT(curthread->t_iplhigh_count == 0);
}
```

syscall must ensure that the kernel adheres to the system call return convention.

CS350

- Exceptions are another way that control is transferred from a process to the kernel.
- Exceptions are conditions that occur during the execution of an instruction by a process. For example, arithmetic overflows, illegal instructions, or page faults (to be discussed later).
- Exceptions are detected by the hardware.
- When an exception is detected, the hardware transfers control to a specific address.
- Normally, a kernel exception handler is located at that address.

Exception handling is similar to, but not identical to, system call handling. (What is different?)

CS350

Operating Systems

Fall 2013

			MIPS Exceptions
EX TRO	0	/ *	Interrupt */
EX MOD	1	/*	TLB Modify (write to read-only page) *
EX TLBL	2	/*	TLB miss on load */
EX TLBS	3	/*	TLB miss on store */
EX ADEL	4	/*	Address error on load */
EX_ADES	5	/*	Address error on store */
EX_IBE	6	/ *	Bus error on instruction fetch */
EX_DBE	7	/ *	Bus error on data load *or* store */
EX_SYS	8	/ *	Syscall */
EX_BP	9	/ *	Breakpoint */
EX_RI	10	/ *	Reserved (illegal) instruction */
EX_CPU	11	/ *	Coprocessor unusable */
EX_OVF	12	/ *	Arithmetic overflow */

## **Interrupts (Revisited)**

- Interrupts are a third mechanism by which control may be transferred to the kernel
- Interrupts are similar to exceptions. However, they are caused by hardware devices, not by the execution of a program. For example:
  - a network interface may generate an interrupt when a network packet arrives
  - a disk controller may generate an interrupt to indicate that it has finished writing data to the disk
  - a timer may generate an interrupt to indicate that time has passed
- Interrupt handling is similar to exception handling current execution context is saved, and control is transferred to a kernel interrupt handler at a fixed address.

CS350

**Operating Systems** 

Fall 2013

Interrupts,	Exceptions, and System Calls: Summar	ry
• interrupts, exceptions is transferred from an	and system calls are three mechanisms by application program to the kernel	which control
• when these events occ and transfers control to be located	o a predefined location, at which a kernel	rivileged mode handler should
• the handler saves the a executed on the CPU, control is returned to t	application thread context so that the kerne and restores the application thread contex the application	el code can be t just before
		_

	Implementation of Processes	
• The kernel main data structure of	tains information about all of the processes i ten called the process table.	n the system in a
• Per-process info	rmation may include:	
<ul> <li>process ident</li> </ul>	tifier and owner	
– the address s	pace for the process	
<ul> <li>threads below</li> </ul>	nging to the process	
<ul> <li>lists of resou</li> </ul>	rces allocated to the process, such as open fil	les
<ul> <li>accounting in</li> </ul>	nformation	
CS350	Operating Systems	Fall 201

Processes and the Ke	rnel	30
	OS/161 Process	
/* From ke	rn/include/proc.h */	
struct pro	c {	
char *p_	name; /* Name of this process */	
struct s	pinlock p_lock; /* Lock for this structure	*/
struct t	hreadarray p_threads; /* Threads in process	s */
struct a struct v	ddrspace *p_addrspace; /* virtual address s node *p_cwd; /* current working directory ;	space */ */
/* add m };	ore material here as needed */	

#### **OS/161 Process**

```
/* From kern/include/proc.h */
/* Create a fresh process for use by runprogram() */
struct proc *proc_create_runprogram(const char *name);
/* Destroy a process */
void proc_destroy(struct proc *proc);
/* Attach a thread to a process */
/* Must not already have a process */
int proc_addthread(struct proc *proc, struct thread *t);
/* Detach a thread from its process */
void proc_remthread(struct thread *t);
...
CS350 Operating Systems Fall 2013
```

	Implementing Timesharing	
• whenever a syste from the running	em call, exception, or interrupt occurs, control is transferre g program to the kernel	d
• at these points, the running process'	he kernel has the ability to cause a context switch from the thread to another process' thread	
• notice that these executing kernel	context switches always occur while a process' thread is code	
By switching fro nel timeshares th	m one process's thread to another process's thread, the ker e processor among multiple processes.	<u>-</u>
350	Operating Systems Fa	11 201











Creationfork,execvfork,execvDestruction_exit,kill_exitSynchronizationwait,waitpid,pause,waitpidAttribute Mgmtgetpid,getuid,nice,getrusage,getpid		Linux	OS/161
Destruction_exit,kill_exitSynchronizationwait,waitpid,pause,waitpidAttribute Mgmtgetpid,getuid,nice,getrusage,getpid	Creation	fork,execv	fork,execv
Synchronizationwait,waitpid,pause,waitpidAttribute Mgmtgetpid,getuid,nice,getrusage,getpid	Destruction	_exit,kill	_exit
Attribute Mgmt getpid,getuid,nice,getrusage, getpid	Synchronization	wait,waitpid,pause,	waitpid
	Attribute Mgmt	getpid,getuid,nice,getrusage,	getpid

## The fork, \_exit, getpid and waitpid system calls

```
main()
{
   rc = fork(); /* returns 0 to child, pid to parent */
   if (rc == 0) {
      my_pid = getpid();
      x = child_code();
      exit(x);
   } else {
     child_pid = rc;
     parent_code();
     child_exit = waitpid(child_pid);
     parent_pid = getpid();
   }
}
CS350
                          Operating Systems
                                                         Fall 2013
```

```
Processes and the Kernel
                                                            40
                      The execv system call
int main()
{
  int rc = 0;
  char *args[4];
  args[0] = (char *) "/testbin/argtest";
  args[1] = (char *) "first";
  args[2] = (char *) "second";
  args[3] = 0;
  rc = execv("/testbin/argtest", args);
  printf("If you see this execv failed\n");
  printf("rc = %d errno = %d\n", rc, errno);
  exit(0);
}
```

```
CS350
```

# The Process Model

- Although the general operations supported by the process interface are straightforward, there are some less obvious aspects of process behaviour that must be defined by an operating system.
  - **Process Initialization:** When a new process is created, how is it initialized? What is in the address space? What is the initial thread context? Does it have any other resources?
  - **Multithreading:** Are concurrent processes supported, or is each process limited to a single thread?
  - **Inter-Process Relationships:** Are there relationships among processes, e.g, parent/child? If so, what do these relationships mean?

CS350

Operating Systems

Fall 2013


	Simple Address Translation: Dynamic Relocation
•	hardware provides a <i>memory management unit</i> which includes a <i>relocation</i> register
•	at run-time, the contents of the relocation register are added to each virtual address to determine the corresponding physical address
•	the OS maintains a separate relocation register value for each process, and ensures that relocation register is reset on each context switch
	Properties
	<ul> <li>each virtual address space corresponds to a contiguous range of physical addresses</li> </ul>
	- OS must allocate/deallocate variable-sized chunks of physical memory
	<ul> <li>potential for <i>external fragmentation</i> of physical memory: wasted, unallocated space</li> </ul>







- Each virtual address space is divided into fixed-size chunks called pages
- The physical address space is divided into *frames*. Frame size matches page size.
- OS maintains a *page table* for each process. Page table specifies the frame in which each of the process's pages is located.
- At run time, MMU translates virtual addresses to physical using the page table of the running process.
- Properties
  - simple physical memory management
  - potential for *internal fragmentation* of physical memory: wasted, allocated space
  - virtual address space need not be physically contiguous in physical space after translation.

Operating Systems

Fall 2013









translation speed: A	
It must be fast.	ldress translation happens very frequently. (How frequently?
sparseness: Many pro their code and dat	grams will only need a small part of the available space for a.
the kernel: Each proc the kernel? In whi	ess has a virtual address space in which to run. What about ch address space does it run?

### **Speed of Address Translation**

- Execution of each machine instruction may involve one, two or more memory operations
  - one to fetch instruction
  - one or more for instruction operands
- Address translation through a page table adds one extra memory operation (for page table entry lookup) for each memory operation performed during instruction execution
  - Simple address translation through a page table can cut instruction execution rate in half.
  - More complex translation schemes (e.g., multi-level paging) are even more expensive.
- Solution: include a Translation Lookaside Buffer (TLB) in the MMU
  - TLB is a fast, fully associative address translation cache
  - TLB hit avoids page table lookup

CS350

Operating Systems

Fall 2013

irtual Memory	12
TLB	
• Each entry in the TLB contains a (page number, frame number	er) pair.
• If address translation can be accomplished using a TLB entry page table is avoided.	, access to the
• Otherwise, translate through the page table, and add the result the TLB, replacing an existing entry if necessary. In a <i>hardwa</i> TLB, this is done by the MMU. In a <i>software controlled</i> TLB kernel.	ting translation to <i>are controlled</i> , it is done by the
• TLB lookup is much faster than a memory access. TLB is an memory - page numbers of all entries are checked simultaneo. However, the TLB is typically small (typically hundreds, e.g. entries).	associative pusly for a match. 128, or 256
• If the MMU cannot distinguish TLB entries from different ad the kernel must clear or invalidate the TLB. (Why?)	dress spaces, then

### The MIPS R3000 TLB

- The MIPS has a software-controlled TLB that can hold 64 entries.
- Each TLB entry includes a virtual page number, a physical frame number, an address space identifier (not used by OS/161), and several flags (valid, read-only).
- OS/161 provides low-level functions for managing the TLB:

TLB\_Write: modify a specified TLB entry

TLB\_Random: modify a random TLB entry

TLB\_Read: read a specified TLB entry

TLB\_Probe: look for a page number in the TLB

• If the MMU cannot translate a virtual address using the TLB it raises an exception, which must be handled by OS/161.

See kern/arch/mips/include/tlb.h

CS350

Operating Systems

Fall 2013

TLB Shootdown				
If one a processor cha mappings of that add	anges the virtual-to-physical mapping of an address, ress in other processors' TLBs would no longer be valid.			
The changing process their TLB.	sor tells the other processors to invalidate that mapping in			
This is called a "TLB (eliminating) entries i	shootdown". The processor is shooting down in other TLBs that are no longer valid.			
In OS/161 is it pos	sible to have the same virtual address stored in multiple			





#### Segmentation

- Often, programs (like sort) need several virtual address segments, e.g, for code, data, and stack.
- One way to support this is to turn *segments* into first-class citizens, understood by the application and directly supported by the OS and the MMU.
- Instead of providing a single virtual address space to each process, the OS provides multiple virtual segments. Each segment is like a separate virtual address space, with addresses that start at zero.
- With segmentation, a virtual address can be thought of as having two parts: (segment ID, address within segment)
- Each segment:
  - can grow (or shrink) independently of the other segments, up to some maximum size
  - has its own attributes, e.g, read-only protection

CS350

Operating Systems

Fall 2013













```
Virtual Memory
                                                                    22
                    OS/161 Address Spaces: dumbvm
 • OS/161 starts with a very simple virtual memory implementation
 • virtual address spaces are described by addrspace objects, which record the
   mappings from virtual to physical addresses
struct addrspace {
#if OPT_DUMBVM
  vaddr_t as_vbasel; /* base virtual address of code segment */
  paddr_t as_pbase1; /* base physical address of code segment */
  size_t as_npages1; /* size (in pages) of code segment */
  vaddr_t as_vbase2; /* base virtual address of data segment */
  paddr_t as_pbase2; /* base physical address of data segment */
  size_t as_npages2; /* size (in pages) of data segment */
  paddr_t as_stackpbase; /* base physical address of stack */
#else
  /* Put stuff here for your VM system */
#endif
};
   This amounts to a slightly generalized version of simple dynamic relocation,
   with three bases rather than one.
```

Fall 2013

### Address Translation Under dumbvm

- the MIPS MMU tries to translate each virtual address using the entries in the TLB
- If there is no valid entry for the page the MMU is trying to translate, the MMU generates a TLB fault (called an *address exception*)
- The vm\_fault function (see kern/arch/mips/vm/dumbvm.c) handles this exception for the OS/161 kernel. It uses information from the current process' addrspace to construct and load a TLB entry for the page.
- On return from exception, the MIPS retries the instruction that caused the page fault. This time, it may succeed.

vm\_fault is not very sophisticated. If the TLB fills up, OS/161 will crash!

CS350

Operating Systems

	Shared Virtual Memory	
• virtual memory sha	aring allows parts of two or more address s	spaces to overlap
• shared virtual mem	nory is:	
<ul> <li>a way to use ph</li> <li>can be shared b</li> </ul>	hysical memory more efficiently, e.g., one only several processes	copy of a program
– a mechanism fo	or interprocess communication	
• sharing is accompl the same physical a	ished by mapping virtual addresses from s address	everal processes to
• unit of sharing can	be a page or a segment	
	On anting Suptana	





An	Addr	ess S	pace	for	the	Kernel
----	------	-------	------	-----	-----	--------

- Each process has its own address space. What about the kernel?
- Three possibilities:

Kernel in physical space: disable address translation in privileged system execution mode, enable it in unprivileged mode

Kernel in separate virtual address space: need a way to change address translation (e.g., switch page tables) when moving between privileged and unprivileged code

Kernel mapped into portion of address space of every process: OS/161,

Linux, and other operating systems use this approach

- memory protection mechanism is used to isolate the kernel from applications
- one advantage of this approach: application virtual addresses (e.g., system call parameters) are easy for the kernel to use

CS350

**Operating Systems** 

Fall 2013







Virtual Memory		3
	ELF Files	
• ELF files contain kernel when it is	n address space segment descriptions, which loading a new address space	n are useful to the
• the ELF file iden	tifies the (virtual) address of the program's	first instruction
<ul> <li>the ELF file also symbol tables) the tools used to built</li> </ul>	contains lots of other information (e.g., sec nat is useful to compilers, linkers, debugger ld programs	ction descriptors, s, loaders and other
CS350	Operating Systems	Fall 201

	Address Space Segments in ELF Files
•	The ELF file contains a header describing the segments and segment <i>images</i> .
•	Each ELF segment describes a contiguous region of the virtual address space.
•	The header includes an entry for each segment which describes:
	- the virtual address of the start of the segment
	- the length of the segment in the virtual address space
	- the location of the start of the segment image in the ELF file (if present)
	- the length of the segment image in the ELF file (if present)
•	the image is an exact copy of the binary data that should be loaded into the specified portion of the virtual address space
•	the image may be smaller than the address space segment, in which case the res of the address space segment is expected to be zero-filled
=	To initialize an address space, the kernel copies images from the ELF file to the specified portions of the virtual address space

	ELF Files and OS/161	
• OS/161's dumb segments:	vm implementation assumes that an ELF file	e contains two
– a text segmer	nt, containing the program code and any rea	d-only data
– a data segme	ent, containing any other global program dat	ta
• the ELF file doe	es not describe the stack (why not?)	
<ul> <li>dumbvm creates</li> <li>virtual address (</li> </ul>	s a <i>stack segment</i> for each process. It is 12 p Dx7fffffff	bages long, ending
Look at kern/s	syscall/loadelf.c to see how OS/16	l loads segments

Virtual Memory		34			
ELF Sections and Segments					
• In the ELF file,	a program's code and data are grouped togeth	her into sections,			
based on their p	roperties. Some sections:				
.text: program	code				
.rodata: read-o	only global data				
.data: initialize	ed global data				
.bss: uninitializ	ed global data (Block Started by Symbol)				
.sbss: small un	initialized global data				
• not all of these	sections are present in every ELF file				
• normally					
- the .text;	and .rodata sections together form the text	t segment			
- the .data,	.bss and .sbss sections together form the	e data segement			
• space for <i>local</i>	program variables is allocated on the stack wh	hen the program			
runs					
CS350	Operating Systems	Fall 2013			

```
The user/uw-testbin/segments.c Example Program (1 of 2)
#include <unistd.h>
#define N
             (200)
int x = 0xdeadbeef;
int t1;
int t2;
int t3;
int array[4096];
char const *str = "Hello World\n";
const int z = 0xabcddcba;
struct example {
  int ypos;
  int xpos;
};
CS350
                          Operating Systems
```

```
Virtual Memory 36
The user/uw-testbin/segments.c Example Program (2 of 2)
int
main()
{
    int count = 0;
    const int value = 1;
    t1 = N;
    t2 = 2;
    count = x + t1;
    t2 = z + t2 + value;
    reboot(RB_POWEROFF);
    return 0; /* avoid compiler warnings */
}
```

35

# **ELF Sections for the Example Program**

Section Headers:						
[Nr] Name	Туре	Addr	Off	Size	Flg	
[ 0]	NULL	00000000	000000	000000		
[ 1] .text	PROGBITS	00400000	010000	000200	AX	
[ 2] .rodata	PROGBITS	00400200	010200	000020	A	
[ 3] .reginfo	MIPS_REGINFO	00400220	010220	000018	A	
[ 4] .data	PROGBITS	10000000	020000	000010	WA	
[ 5] .sbss	NOBITS	10000010	020010	000014	WAp	
[ 6] .bss	NOBITS	1000030	020010	004000	WA	
 Flags: W (write),	A (alloc), X (ex	cecute), p	p (proce	essor sp	pecific)	
## Size = number of bytes (e.g., .text is 0x200 = 512 bytes ## Off = offset into the ELF file ## Addr = virtual address						

The cs350-readelf program can be used to inspect OS/161 MIPS ELF files: cs350-readelf -a segments

CS350

Operating Systems

Fall 2013

Virtual Memor	ту						38		
ELF Segments for the Example Program									
Program Type REGINFO LOAD LOAD	Headers: Offset 0x010220 0x010000 0x020000	VirtAddr 0x00400220 0x00400000 0x10000000	PhysAddr 0x00400220 0x00400000 0x10000000	FileSiz 0x00018 0x00238 0x00010	MemSiz 0x00018 0x00238 0x04030	Flg R R E RW	Align 0x4 0x10000 0x10000		
<ul> <li>segment info, like section info, can be inspected using the cs350-readelf program</li> <li>the REGINFO section is not used</li> </ul>									
• the fir	st LOAD se	gment includes	s the .text and .	rodata sect	ions				
• the se	cond LOAD	segment inclu	des .data, .sbss	, and .bss					

#### Contents of the Example Program's .text Section

```
Contents of section .text:
 400000 3c1c1001 279c8000 2408fff8 03a8e824 <....$
. . .
## Decoding 3c1c1001 to determine instruction
## 0x3c1c1001 = binary 1111000001110000010000000000
## instr
         | rs
                 | rt
                        immediate
## 6 bits | 5 bits | 5 bits |
                            16 bits
## 001111 | 00000 | 11100 |
                           0001 0000 0000 0001
## LUI
         0
                 | reg 28|
                            0x1001
## LUI
         unused reg 28
                            0x1001
## Load upper immediate into rt (register target)
## lui gp, 0x1001
  The cs350-objdump program can be used to inspect OS/161 MIPS ELF
  file section contents: cs350-objdump -s segments
```

CS350

Operating Systems

Fall 2013

```
ytuny y
```

#### Contents of the Example Program's .data Section

```
Contents of section .data:
  10000000 deadbeef 00400210 00000000 00000000 ....@....@.....
## Size = 0x10 bytes = 16 bytes (padding for alignment)
## int x = deadbeef (4 bytes)
## char const *str = "Hello World\n"; (4 bytes)
## address of str = 0x1000004
## value stored in str = 0x00400210.
## NOTE: this is the address of the start
## of the string literal in the .rodata section
```

The .data section contains the initialized global variables str and x.

CS350

Operating Systems

Fall 2013

41

```
42
Virtual Memory
       Contents of the Example Program's .bss and .sbss Sections
. . .
10000000 D x
10000004 D str
                     ## S indicates sbss section
10000010 S t3
10000014 S t2
10000018 S t1
1000001c S errno
10000020 S __argv
                     ## B indicates bss section
10000030 B array
10004030 A _end
10008000 A _gp
```

The t1, t2, and t3 variables are in the .sbss section. The array variable is in the .bss section. There are no values for these variables in the ELF file, as they are uninitialized. The cs350-nm program can be used to inspect symbols defined in ELF files: cs350-nm -n <filename>, in this case cs350-nm -n segments.



Virtual Memory	44
<b>Exploiting Secondary Storage</b>	
Goals:	
• Allow virtual address spaces that are larger than the physical address	s space.
• Allow greater multiprogramming levels by using less of the available memory for each process.	le (primary)
Method:	
• Allow pages (or segments) from the virtual address space to be store secondary memory, as well as primary memory.	ed in
• Move pages (or segments) between secondary and primary memory are in primary memory when they are needed.	so that they



Large Virtual Address Spaces	
Virtual memory allows for very large virtual address spaces, and ve virtual address spaces require large page tables.	ery large
example: $2^{48}$ byte virtual address space, 8 Kbyte ( $2^{13}$ byte) pages, table entries means	4 byte page
$\frac{2^{48}}{2^{13}}2^2 = 2^{37}$ bytes per page table	
page tables for large address spaces may be very large, and	
- they must be in memory, and	
<ul> <li>they must be physically contiguous</li> </ul>	
some solutions:	
– multi-level page tables - page the page tables	
– inverted page tables	





### **Paging Policies**

## When to Page?:

Demand paging brings pages into memory when they are used. Alternatively, the OS can attempt to guess which pages will be used, and *prefetch* them.

## What to Replace?:

Unless there are unused frames, one page must be replaced for each page that is loaded into memory. A replacement policy specifies how to determine which page to replace.

Similar issues arise if (pure) segmentation is used, only the unit of data transfer is segments rather than pages. Since segments may vary in size, segmentation also requires a *placement policy*, which specifies where, in memory, a newly-fetched segment should be placed.

CS350

**Operating Systems** 

Fall 2013

Global vs. Local Page Replacement						
• When the system's should the replacer	page reference string is generated by ment policy take this into account?	nore than one process,				
<b>Global Policy:</b> A the process to w may replace a p	global policy is applied to all in-memory which each one "belongs". A page reque age that belongs another process, Y.	ry pages, regardless of ested by process X				
<b>Local Policy:</b> Und processes accor policy is then ap requested by pr	ler a local policy, the available frames a ding to some memory allocation policy oplied separately to each process's allo ocess X replaces another page that "bel	are allocated to 7. A replacement cated space. A page longs" to process X.				
2220		E 11 2012				



A Simple Replacement Policy: FIFO									
• the FIFO po	olicy: repla	ce the pag	ge that l	nas bee	en in me	emory	the le	onges	st
• a three from	ne evemple								

Num	1	2	3	4	5	6	7	8	9	10	11	12
Refs	a	b	c	d	a	b	e	a	b	c	d	e
Frame 1	a	a	a	d	d	d	e	e	e	e	e	e
Frame 2		b	b	b	a	a	a	a	a	c	с	c
Frame 3			с	c	с	b	b	b	b	b	d	d
Fault ?	X	X	X	X	X	X	X			Х	X	

# **Optimal Page Replacement**

- There is an optimal page replacement policy for demand paging.
- The OPT policy: replace the page that will not be referenced for the longest time.

Num	1	2	3	4	5	6	7	8	9	10	11	12
Refs	a	b	c	d	a	b	e	a	b	с	d	e
Frame 1	a	a	a	a	a	a	a	a	a	с	с	c
Frame 2		b	b	b	b	b	b	b	b	b	d	d
Frame 3			c	d	d	d	e	e	e	e	e	e
Fault?	х	х	X	х			x			Х	Х	

• OPT requires knowledge of the future.

CS350

**Operating Systems** 

<b>Other Replacement Policies</b>	
• FIFO is simple, but it does not consider:	
Frequency of Use: how often a page has been used?	
Recency of Use: when was a page last used?	
Cleanliness: has the page been changed while it is in memory?	
• The <i>principle of locality</i> suggests that usage ought to be considered in a replacement decision.	
• Cleanliness may be worth considering for performance reasons.	

	Locality	
<ul> <li>Locality is a proper property of program</li> </ul>	rty of the page reference string. In other w ms themselves.	ords, it is a
• Temporal locality s used again.	says that pages that have been used recently	y are likely to be
• Spatial locality say be used next.	vs that pages "close" to those that have bee	n used are likely to
In practice, page re	ference strings exhibit strong locality. Wh	y?
CS350	Operating Systems	Fall 201

- Counting references to pages can be used as the basis for page replacement decisions.
- Example: LFU (Least Frequently Used) Replace the page with the smallest reference count.
- Any frequency-based policy requires a reference counting mechanism, e.g., MMU increments a counter each time an in-memory page is referenced.
- Pure frequency-based policies have several potential drawbacks:
  - Old references are never forgotten. This can be addressed by periodically reducing the reference count of every in-memory page.
  - Freshly loaded pages have small reference counts and are likely victims ignores temporal locality.

Operating Systems

			Lea	ast F	Rece	ntly	Use	d: L	RU				
ie three-	frai	me e	exam	ple:									
Nui	m	1	2	3	4	5	6	7	8	9	10	11	12
Ret	fs	a	b	c	d	a	b	e	а	b	c	d	e
Frame	1	a	а	а	d	d	d	e	e	e	c	с	c
Frame	2		b	b	b	a	a	a	а	а	a	d	d
Frame	3			c	c	c	b	b	b	b	b	b	e
Fault	?	x	x	x	x	х	x	x			x	Х	x



What if t	he MMU Does Not Provide a "Use" B	it?
• the kernel can emulat	e the "use" bit, at the cost of extra excep	ptions
1. When a page is lo been loaded) and	aded into memory, mark it as <i>invalid</i> (e set its simulated "use" bit to false.	ven though it as
2. If a program atten	npts to access the page, an exception wi	ll occur.
3. In its exception hat and marks the pag	undler, the OS sets the page's simulated ge <i>valid</i> so that further accesses do not c	"use" bit to "true" ause exceptions.
• This technique requir page:	es that the OS maintain extra bits of inf	ormation for each
1. the simulated "use	e" bit	
2. an "in memory" b	it to indicate whether the page is in mer	nory
2. un minemory o	it to indicate whether the page is in mor	nory

CS350

Operating Systems

	Page Cleanliness: the "Modified" Bit
• A page is <i>modi</i> loaded into me	<i>fied</i> (sometimes called dirty) if it has been changed since it was nory.
• A modified pag	e is more costly to replace than a clean page. (Why?)
• The MMU iden when the conte	tifies modified pages by setting a <i>modified bit</i> in the TLB entry nts of the page change.
• Operating syste	m clears the modified bit when it cleans the page
• The modified b	it potentially has two roles:
- Indicates w	nich pages need to be cleaned.
– Can be used	to influence the replacement policy.
MIPS TLB entr	ies do not include a modified bit.

What if the MMU I	Does Not Provide a "Modified" Bit?
• Can emulate it in similar fash	ion to the "use" bit
1. When a page is loaded int actually writeable) and se	o memory, mark it as <i>read-only</i> (even if it is t is t is simulated "modified" bit to false.
2. If a program attempts to n	nodify the page, a protection exception will occur.
3. In its exception handler, if the page's simulated "moo writeable.	the page is supposed to be writeable, the OS sets dified" bit to "true" and marks the page as
• This technique requires that the each page:	he OS maintain two extra bits of information for
1. the simulated "modified"	bit
2. a "writeable" bit to indica	te whether the page is supposed to be writeable

Operating Systems

Enhanced Second Chance Replacement Algorithm					
• Classify pages according to their use and modified bits:					
(0,0): not recently used, clean.					
(0,1): not recently used, modified.					
(1,0): recently used, clean					
(1,1): recently used, modified					
Algorithm:					
1. Sweep once looking for $(0,0)$ page. Don't clear use b	oits while looking.				
2. If none found, look for (0,1) page, this time clearing frames.	"use" bits for bypassed				
3. If step 2 fails, all use bits will be zero, repeat from st (guaranteed to find a page).	ep 1				

#### Page Cleaning

- A modified page must be cleaned before it can be replaced, otherwise changes on that page will be lost.
- *Cleaning* a page means copying the page to secondary storage.
- Cleaning is distinct from replacement.
- Page cleaning may be *synchronous* or *asynchronous*:

**synchronous cleaning:** happens at the time the page is replaced, during page fault handling. Page is first cleaned by copying it to secondary storage. Then a new page is brought in to replace it.

- **asynchronous cleaning:** happens before a page is replaced, so that page fault handling can be faster.
  - asynchronous cleaning may be implemented by dedicated OS *page cleaning threads* that sweep through the in-memory pages cleaning modified pages that they encounter.

CS350

CS350

Operating Systems

Fall 2013

Fall 2013

				DC.	lauy	511	non	laiy					
FIFO re	placement	, 4 f	rame	es									
	Num	1	2	3	4	5	6	7	8	9	10	11	12
	Refs	a	b	с	d	a	b	e	a	b	c	d	e
	Frame 1	a	a	a	a	a	a	e	e	e	e	d	d
	Frame 2		b	b	b	b	b	b	a	a	а	a	e
	Frame 3			с	c	с	c	c	c	b	b	b	b
	Frame 4				d	d	d	d	d	d	c	с	c
Γ	Fault?	X	x	x	x			X	X	x	х	X	x
FIFO ex faults.	ample on	Slid	e 52	with	n san	ne re	efere	nce	strin	g ha	d 3 fr	ames	and or

Operating Systems

# Stack Policies

67

Fall 2013

- Let B(m, t) represent the set of pages in the system with m frames of memory, at time t, under some given replacement policy, for some given reference string.
- A replacement policy is called a *stack policy* if, for all reference strings, all m and all t:

$$B(m,t) \subseteq B(m+1,t)$$

- If a replacement algorithm imposes a total order, independent of the number of frames (i.e., memory size), on the pages and it replaces the largest (or smallest) page according to that order, then it satisfies the definition of a stack policy.
- Examples: LRU is a stack algorithm. FIFO and CLOCK are not stack algorithms. (Why?)

Stack algorithms do not suffer from Belady's anomaly.

CS350

Operating Systems

```
Virtual Memory
                                                                                   68
                                    Prefetching
  • Prefetching means moving virtual pages into memory before they are needed,
    i.e., before a page fault results.
  • The goal of prefetching is latency hiding: do the work of bringing a page into
    memory in advance, not while a process is waiting.
  • To prefetch, the operating system must guess which pages will be needed.
  • Hazards of prefetching:
     - guessing wrong means the work that was done to prefetch the page was
        wasted
     - guessing wrong means that some other potentially useful page has been
        replaced by a page that is not used
  • most common form of prefetching is simple sequential prefetching: if a process
    uses page x, prefetch page x + 1.

    sequential prefetching exploits spatial locality of reference

                                                                              Fall 2013
CS350
                                    Operating Systems
```



	How Much Physical Memory Does a Process Need?				
•	Principle of locality suggests that some portions of the process's virtual address space are more likely to be referenced than others.				
•	A refinement of this principle is the <i>working set model</i> of process reference behaviour.				
•	According to the working set model, at any given time some portion of a program's address space will be heavily used and the remainder will not be. The heavily used portion of the address space is called the <i>working set</i> of the process.				
•	The working set of a process may change over time.				
•	The resident set of a process is the set of pages that are located in memory.				
=	According to the working set model, if a process's resident set includes its working set it will rarely page fault				

**Resident Set Sizes (Example)** PID VSZ RSS COMMAND 805 13940 5956 /usr/bin/qnome-session 831 848 /usr/bin/ssh-agent 2620 834 7936 5832 /usr/lib/gconf2/gconfd-2 11 838 6964 2292 gnome-smproxy 840 14720 5008 gnome-settings-daemon 848 8412 3888 sawfish 851 34980 7544 nautilus 853 19804 14208 gnome-panel 9656 2672 gpilotd 857 867 4608 1252 gnome-name-service

CS350

Operating Systems

Fall 2013


Virtual Memory		7
	Page Fault Frequency	
• A more direct w <i>fault frequencie</i>	yay to allocate memory to processes is to meas s - the number of page faults they generate p	asure their <i>page</i> er unit time.
• If a process's pa low, it may be a	age fault frequency is too high, it needs more ble to surrender memory.	memory. If it is
• The working set sharp "knee".	t model suggests that a page fault frequency	plot should have a
C\$350	Operating Systems	Fall 201





Virtual Memory	76
Swapping Out Processes	
• Swapping a process out means removing all of its pages from memory, or marking them so that they will be removed by the normal page replacement process. Suspending a process ensures that it is not runnable while it is swappe out.	ed
<ul> <li>Which process(es) to suspend?</li> <li>low priority processes</li> <li>blocked processes</li> <li>large processes (lots of space freed) or small processes (easier to reload)</li> </ul>	
• There must also be a policy for making suspended processes ready when system load has decreased.	m

	The Neture of Dreamon Executions	
	The Nature of Frogram Executions	
• A running thr <i>bursts</i>	read can be modeled as alternating series of C.	<i>PU bursts</i> and <i>I/O</i>
– during a C	CPU burst, a thread is executing instructions	
<ul> <li>during an and is not</li> </ul>	I/O burst, a thread is waiting for an I/O opera executing instructions	tion to be performed
	Operating Systems	Eall 201
CS350	Operating Systems	1 all 201
CS350	Operating Systems	
CS350	Operating Systems	
CS350 Processor Scheduling	Operating Systems	
CS350 Processor Scheduling	Preemptive vs. Non-Preemptive	
CS350 Processor Scheduling	Preemptive vs. Non-Preemptive	
Processor Scheduling     A non-preemp	Preemptive vs. Non-Preemptive ptive scheduler runs only when the running th pugh its own actions. e.g.,	read gives up the
<ul> <li>Processor Scheduling</li> <li>A non-preemp processor thro – the thread</li> </ul>	Preemptive vs. Non-Preemptive <i>ptive</i> scheduler runs only when the running th pugh its own actions, e.g., terminates	read gives up the
<ul> <li>Processor Scheduling</li> <li>A non-preemp processor throon processor thro</li></ul>	Preemptive vs. Non-Preemptive <i>ptive</i> scheduler runs only when the running th pugh its own actions, e.g., terminates blocks because of an I/O or synchronization	read gives up the
<ul> <li>CS350</li> <li>Processor Scheduling</li> <li>A non-preemp processor throon processor throon the thread - the thre</li></ul>	Preemptive vs. Non-Preemptive <i>ptive</i> scheduler runs only when the running th pugh its own actions, e.g., terminates blocks because of an I/O or synchronization of performs a Yield system call (if one is provid	read gives up the operation
<ul> <li>CS350</li> <li>Processor Scheduling</li> <li>A non-preemp processor throon processor throon the thread of the thread of the thread system)</li> </ul>	Preemptive vs. Non-Preemptive <i>ptive</i> scheduler runs only when the running th pugh its own actions, e.g., terminates blocks because of an I/O or synchronization of performs a Yield system call (if one is provid	read gives up the operation led by the operating
<ul> <li>CS350</li> <li>Processor Scheduling</li> <li>A non-preemp processor throad - the thread - the thread - the thread system)</li> <li>A preemptive</li> </ul>	Preemptive vs. Non-Preemptive <i>ptive</i> scheduler runs only when the running th pugh its own actions, e.g., terminates blocks because of an I/O or synchronization of performs a Yield system call (if one is provid scheduler may, in addition, force a running th	read gives up the operation led by the operating mead to stop running
<ul> <li>A non-preemp processor Scheduling</li> <li>A non-preemp processor throad - the thread - the thread - the thread system)</li> <li>A preemptive - typically, a</li> </ul>	Preemptive vs. Non-Preemptive <i>ptive</i> scheduler runs only when the running th pugh its own actions, e.g., terminates blocks because of an I/O or synchronization of performs a Yield system call (if one is provid scheduler may, in addition, force a running th a preemptive scheduler will be invoked period	read gives up the operation led by the operating pread to stop running lically by a timer

- a running thread that is preempted is moved to the ready state

FCFS and Round-Robin Scheduling			
First-Come, Firs	st-Served (FCFS):		
• non-pree	mptive - each thread runs until it blocks or termin	ates	
• FIFO rea	dy queue		
Round-Robin:			
• preempti	ve version of FCFS		
<ul> <li>running t blocked</li> </ul>	hread is preempted after a fixed time quantum, if	it has not already	
• preempte	ed thread goes to the end of the FIFO ready queue	:	

Shortest Job First (SJF) Scheduling		
• n	non-preemptive	
• r t	eady threads are scheduled according to the length of their next CPU burst - hread with the shortest burst goes first	
• 5	SJF minimizes average waiting time, but can lead to starvation	
• 5	SJF requires knowledge of CPU burst lengths	
	<ul> <li>Simplest approach is to estimate next burst length of each thread based on previous burst length(s). For example, exponential average considers all previous burst lengths, but weights recent ones most heavily:</li> </ul>	
	$B_{i+1} = \alpha b_i + (1 - \alpha)B_i$	
	where $B_i$ is the predicted length of the <i>i</i> th CPU burst, and $b_i$ is its actual length, and $0 \le \alpha \le 1$ .	
• S c	Shortest Remaining Time First is a preemptive variant of SJF. Preemption may occur when a new thread enters the ready queue.	









Processor Scheduling

### **Highest Response Ratio Next**

- non-preemptive
- response ratio is defined for each ready thread as:

$$\frac{w+b}{b}$$

where b is the estimated CPU burst time and w is the actual waiting time

- scheduler chooses the thread with the highest response ratio (choose smallest *b* in case of a tie)
- HRRN is an example of a heuristic that blends SJF and FCFS

CS350

Operating Systems

Fall 2013



	Prioritization	
• a scheduler r	nay be asked to take process or thread prioriti	es into account
• for example,	priorities could be based on	
– user class	sification	
– applicatio	on classification	
– applicatio	on specification	
(e.g., Lin	ux setpriority/sched_setschedule	er)
• scheduler car	n:	
<ul> <li>always cł</li> </ul>	noose higher priority threads over lower prior	ity threads
– use any s	cheduling heuristic to schedule threads of equ	al priority
• low priority t mechanism f	threads risk starvation. If this is not desired, s or elevating the priority of low priority thread	cheduler must have a ls that have waited a
long time		
CS350	Operating Systems	Fall 201

Multilevel Feedback Queues	
gives priority to interactive threads (those with short CPU bursts)	
scheduler maintains several ready queues	
• scheduler never chooses a thread in ready queue $i$ if there are threads in ready queue $j < i$ .	n any
• threads in ready queue $i$ use quantum $q_i$ , and $q_i < q_j$ if $i < j$	
• newly ready threads go into ready queue 0	
a level $i$ thread that is preempted goes into the level $i + 1$ ready queue	









	<b>Devices and Device Controllers</b>	
• network interface		
• graphics adapter		
• secondary storage (c	lisks, SSD, flash) and storage controllers	
• serial (e.g., mouse, l	keyboard)	
• sound		
• co-processors		
•		
		E 11 201









0			6	
Example: LAMEbus timer device registers				
Offset	Size	Type	Description	
0	1	status	current time (seconds)	
4	4	status	current time (nanoseconds)	
8	4	command	restart-on-expiry (auto-restart countdown?)	
12	4	status and command	interrupt (reading clears)	
16	4	status and command	countdown time (microseconds)	
20	4	command	speaker (causes beeps)	

Sys/161 uses memory-mapping. Each device's registers are mapped into the *physical address space* of the MIPS processor.

CS350

I/O

Example: L	AMEbus	disk	controller
------------	--------	------	------------

Offset	Size	Туре	Description
0	4	status	number of sectors
4	4	status and command	status
8	4	command	sector number
12	4	status	rotational speed (RPM)
32768	512	data	transfer buffer

CS350

Operating Systems

Fall 2013

7



I/O

### **Device Control Example: Controlling the Timer**

```
/* Registers (offsets within the device slot) */
#define LT_REG_SEC
                     0 /* time of day: seconds */
#define LT_REG_NSEC 4 /* time of day: nanoseconds */
#define LT_REG_ROE
                     8 /* Restart On countdown-timer Expiry flag
#define LT_REG_IRQ
                     12 /* Interrupt status register */
#define LT_REG_COUNT 16 /* Time for countdown timer (usec) */
#define LT_REG_SPKR 20 /* Beep control */
/* Get the number of seconds from the lamebus timer */
/* lt->lt_buspos is the slot number of the target device */
secs = bus_read_register(lt->lt_bus, lt->lt_buspos,
    LT_REG_SEC);
/* Get the timer to beep. Doesn't matter what value is sent */
bus_write_register(lt->lt_bus, lt->lt_buspos,
    LT_REG_SPKR, 440);
CS350
```

**Operating Systems** 

Fall 2013

```
I/O
                                                           10
            Device Control Example: Address Calculations
/* LAMEbus mapping size per slot */
#define LB_SLOT_SIZE
                               65536
#define MIPS_KSEG1
                    0xa0000000
#define LB_BASEADDR
                      (MIPS_KSEG1 + 0x1fe00000)
/* Compute the virtual address of the specified offset */
/* into the specified device slot */
void *
lamebus_map_area(struct lamebus_softc *bus, int slot,
                 u_int32_t offset)
{
    u_int32_t address;
    (void)bus;
                 // not needed
    assert(slot>=0 && slot<LB_NSLOTS);</pre>
    address = LB_BASEADDR + slot*LB_SLOT_SIZE + offset;
    return (void *)address;
ł
```

I/O

# **Device Control Example: Commanding the Device**

```
/* FROM: kern/arch/mips/mips/lamebus_mips.c */
/* Read 32-bit register from a LAMEbus device. */
u_int32_t
lamebus_read_register(struct lamebus_softc *bus,
    int slot, u_int32_t offset)
{
    u_int32_t *ptr = lamebus_map_area(bus, slot, offset);
    return *ptr;
}
/* Write a 32-bit register of a LAMEbus device. */
void
lamebus write register(struct lamebus softc *bus,
    int slot, u_int32_t offset, u_int32_t val)
{
    u_int32_t *ptr = lamebus_map_area(bus, slot, offset);
    *ptr = val;
}
CS350
                         Operating Systems
                                                       Fall 2013
```

I/O 12 **Device Data Transfer** • Sometimes, a device operation will involve a large chunk of data - much larger than can be moved with a single instruction. Example: reading a block of data from a disk. • Devices may have data buffers for such data - but how to get the data between the device and memory? • If the data buffer is memory-mapped, the kernel can move the data iteratively, one word at a time. This is called *program-controlled I/O*. • Program controlled I/O is simple, but it means that the CPU is *busy executing* kernel code while the data is being transferred. • The alternative is called Direct Memory Access (DMA). During a DMA data transfer, the CPU is not busy and is free to do something else, e.g., run an application. Sys/161 LAMEbus devices do program-controlled I/O.













### **Rotational Latency and Transfer Time**

- rotational latency depends on the rotational speed of the disk
- if the disk spins at  $\omega$  rotations per second:

$$0 \le t_{rot} \le \frac{1}{\omega}$$

• expected rotational latency:

$$\bar{t}_{rot} = \frac{1}{2\omega}$$

- transfer time depends on the rotational speed and on the amount of data transferred
- if k sectors are to be transferred and there are T sectors per track:

$$t_{transfer} = \frac{k}{T\omega}$$

CS350

Operating Systems

Fall 2013



Performance Implications of Disk Characteristics		
• larger transfers is, the cost (tim	s to/from a disk device are <i>more efficient</i> than ne) per byte is smaller for larger transfers. (V	n smaller ones. That Vhy?)
• sequential I/O	is faster than non-sequential I/O	
– sequential	I/O operations eliminate the need for (most)	seeks
<ul> <li>disks use o</li> <li>sequential</li> </ul>	ther techniques, like <i>track buffering</i> , to reduct I/O even more	ee the cost of
 C\$350	Operating Systems	Fall 201















### **Files and File Systems**

- files: persistent, named data objects
  - data consists of a sequence of numbered bytes
  - alternatively, a file may have some internal structure, e.g., a file may consist of sequence of numbered records
  - file may change size over time
  - file has associated meta-data (attributes), in addition to the file name
    - \* examples: owner, access controls, file type, creation and access timestamps
- file system: a collection of files which share a common name space
  - allows files to be created, destroyed, renamed, ...

CS350

Operating Systems

Fall 2013

	File Interface
• open, close	
– open returns a	file identifier (or handle or descriptor), which is used in
subsequent ope	erations to identify the file. (Why is this done?)
• read, write	
- must specify w	which file to read, which part of the file to read, and where to
put the data that	at has been read (similar for write).
<ul> <li>– often, file posit</li> </ul>	tion is implicit (why?)
• seek	
• get/set file attribute	es, e.g., Unix fstat, chmod





File Systems

# Sequential File Reading Example (Unix)

```
char buf[512];
int i;
int f = open("myfile",O_RDONLY);
for(i=0; i<100; i++) {
  read(f,(void *)buf,512);
}
close(f);
```

Read the first 100 \* 512 bytes of a file, 512 bytes at a time.

CS350

Operating Systems

Fall 2013

5

```
Ele Systems

File Systems

File Reading Example Using Seek (Unix)

char buf[512];
int i;
int f = open("myfile", O_RDONLY);
for(i=1; i<=100; i++) {
    lseek(f,(100-i)*512,SEEK_SET);
    read(f,(void *)buf,512);
}
close(f);
</pre>
```

Read the first 100 \* 512 bytes of a file, 512 bytes at a time, in reverse order.

File Systems

# File Reading Example Using Positioned Read

```
char buf[512];
int i;
int f = open("myfile",O_RDONLY);
for(i=0; i<100; i+=2) {
    pread(f,(void *)buf,512,i*512);
}
close(f);
```

Read every second 512 byte chunk of a file, until 50 have been read.

CS350

Operating Systems

Fall 2013

<pre>generic interface: vaddr ← mmap(file descriptor,fileoffset,length) munmap(vaddr,length) mmap call returns the virtual address to which the file is mapped munmap call unmaps mapped files within the specified virtual address range Memory-mapping is an alternative to the read/write file interface.</pre>		Memory-Mapped Files
<pre>vaddr</pre>	generic interface	:
mmap call returns the virtual address to which the file is mapped munmap call unmaps mapped files within the specified virtual address range Memory-mapping is an alternative to the read/write file interface.	$vaddr \leftarrow mma$ munmap(vadd	ap(file descriptor,fileoffset,length) r,length)
munmap call unmaps mapped files within the specified virtual address range Memory-mapping is an alternative to the read/write file interface.	mmap call return	s the virtual address to which the file is mapped
Memory-mapping is an alternative to the read/write file interface.	munmap call un	maps mapped files within the specified virtual address range
	Memory-mappin	g is an alternative to the read/write file interface.







	File Names
• application-visibl	e objects (e.g., files, directories) are given names
• the file system is	responsible for associating names with objects
• the namespace is	typically structured, often as a tree or a DAG
• namespace struct manage informati	ure provides a way for users and applications to organize and ion
• in a structured na describe a path fr	mespace, objects may be identified by <i>pathnames</i> , which om a root object to the object being identified, e.g.:
/home/	user/courses/cs350/notes/filesys.pdf





## Symbolic Links

- a *Symbolic link*, or *soft link*, is an association between two names in the file namespace. Think of it is a way of defining a synonym for a filename.
  - symlink(oldpath, newpath) creates a symbolic link from newpath to oldpath, i.e., newpath becomes a synonym for oldpath.
- symbolic links relate filenames to filenames, while hard links relate filenames to underlying file objects!

• referential integrity is *not* preserved for symbolic links, e.g., the system call above can succeed even if there is no object named oldpath

CS350

Operating Systems

Fall 2013

	UNIX/Linux Link Example (1 of 3)	
<pre>% cat &gt; file1 This is file1 <cntl-d> % ls -li 685844 -rw % ln file1 li % ln -s file1 % ln not-here ln: not-here: % ln -s not-h</cntl-d></pre>	1 user group 15 2008-08-20 file1 nk1 sym1 link2 No such file or directory ere sym2	
Files, hard lir	nks, and soft/symbolic links.	

File Systems

#### UNIX/Linux Link Example (2 of 3)

% ls -li 685844 -rw------ 2 user group 15 2008-08-20 file1 685844 -rw------ 2 user group 15 2008-08-20 link1 685845 lrwxrwxrwx 1 user group 5 2008-08-20 sym1 -> file1 685846 lrwxrwxrwx 1 user group 8 2008-08-20 sym2 -> not-here % cat file1 This is file1. % cat link1 This is file1. % cat sym1 This is file1. % cat sym2 cat: sym2: No such file or directory % /bin/rm file1

Accessing and manipulating files, hard links, and soft/symbolic links.

CS350

Operating Systems

Fall 2013

17

```
File Systems
                                                                18
                  UNIX/Linux Link Example (3 of 3)
% ls -li
685844 -rw----- 1 user group 15 2008-08-20 link1
685845 lrwxrwxrwx 1 user group 5 2008-08-20 sym1 -> file1
685846 lrwxrwxrwx 1 user group 8 2008-08-20 sym2 -> not-here
% cat link1
This is file1.
% cat sym1
cat: sym1: No such file or directory
% cat > file1
This is a brand new file1.
<cntl-d>
% ls -li
685847 -rw----- 1 user group 27 2008-08-20 file1
685844 -rw----- 1 user group 15 2008-08-20 link1
685845 lrwxrwxrwx 1 user group 5 2008-08-20 sym1 -> file1
685846 lrwxrwxrwx 1 user group 8 2008-08-20 sym2 -> not-here
% cat link1
This is file1.
% cat sym1
This is a brand new file1.
   Different behaviour for hard links and soft/symbolic links.
```

Fall 2013





## Links and Multiple File Systems

- a hard link associates a name in the file system namespace with a file in that file system
- typically, hard links cannot cross file system boundaries
- for example, even after the mount operation illustrated on the previous slide, link(/x/a/x/g,/z/d) would result in an error, because the new link, which is in the root file system refers to an object in file system X
- soft links do not have this limitation
- for example, after the mount operation illustrated on the previous slide:
  - symlink(/x/a/x/g,/z/d) would succeed
  - open(/z/d) would succeed, with the effect of opening /z/a/x/g.
- even if the symlink operation were to occur *before* the mount command, it would succeed

CS350

Operating Systems

Fall 2013

File System Implementation	
• space management	
• file indexing (how to locate file data and meta-data)	
• directories	
• links	
• buffering, in-memory data structures	
• fault tolerance	

	S	pace A	Alloca	ation	anc	d L	ayo	ut					
• space may b	e allocated	in fixe	d-size	e chu	ınks,	, or	in c	chui	ıks	of v	vary	ing	size
• fixed-size cl	nunks: simp	le spac	ce ma	nage	men	nt, b	out i	nter	nal	frag	gme	entat	tion
• variable-size	e chunks: ex	ternal	fragn	nenta	atior	ı							
													]
		fixed-s	size al	locati	ion								
		iixed i	5120 ui	locut									
													]
		variabl	e–size	allo	catio	n							
• <i>lavout</i> matte	ers! Try to 1	av a fil	e out	seau	enti	allv	7 Of	• in 1	laro	e se	ane	entia	l extents
that can be	read and wr	itten ef	ficien	itly.	UIII)	uny	, 01	111	urg		que	/11/10	a extents
6250			0										E 11 00

	File Indexing
in	general, a file will require more than one chunk of allocated space
th	is is especially true because files can grow
ho	ow to find all of a file's data?
ch	naining:
	<ul> <li>each chunk includes a pointer to the next chunk</li> </ul>
	- OK for sequential access, poor for random access
ex	ternal chaining: DOS file allocation table (FAT), for example
	- like chaining, but the chain is kept in an external structure
pe	er-file index: Unix i-node, for example
	- for each file, maintain a table of pointers to the file's blocks or extents








File Systems		29
	Example: Unix i-nodes	
• an i-node is a r	ecord describing a file	
• each i-node is u location on the	uniquely identified by an i-number, which dete disk	ermines its physical
• an i-node is a fi	ixed size record containing:	
file attribute v	alues	
– file type		
– file owne	r and group	
<ul> <li>access co</li> </ul>	ontrols	
- creation,	reference and update timestamps	
– file size		
direct block p	<b>binters:</b> approximately 10 of these	
single indirect	block pointer	
double indirec	t block pointer	
triple indirect	block pointer	
2S350	Operating Systems	Fall 2013



#### **Directories**

- A directory consists of a set of entries, where each entry is a record that includes:
  - a file name (component of a path name)
  - the internal file identifier (e.g., i-number) of the file
- A directory can be implemented as a special type of file. The directory entries are the contents of the file.
- The file system should not allow directory files to be directly written by application programs. Instead, the directory is updated by the file system as files are created and destroyed

CS350

**Operating Systems** 

Fall 2013





### **Implementing Soft Links**

- soft links can be implemented as a special type of file
- for example, consider:

symlink(/y/k/g,/z/m)

- to implement this:
  - create a new symlink file
  - add a new entry in directory /z
    - \* file name in new entry is m
    - \* i-number in the new entry is the i-number of the new symlink file
  - store the pathname string "/y/k/g" as the contents of the new symlink file
- change the behaviour of the open system call so that when the symlink file is encountered during open(/z/m), the file /y/k/g will be opened instead.

CS350

Operating Systems

Fall 2013



### **Problems Caused by Failures**

- a single logical file system operation may require several disk I/O operations
- example: deleting a file
  - remove entry from directory
  - remove file index (i-node) from i-node table
  - mark file's data blocks free in free space index
- what if, because of a failure, some but not all of these changes are reflected on the disk?

CS350

Operating Systems

Fall 2013



Ir	terprocess Communication Mechanisms	
• shared storage		
– These mechan	nisms have already been covered. examples:	
* shared virtu	ual memory	
* shared files	3	
<ul> <li>processes musical</li> </ul>	st agree on a name (e.g., a file name, or a shared virtua	.1
memory key)	in order to establish communication	
• message based		
– signals		
– sockets		
– pipes		
CS350	Operating Systems	Fall 2013



Pro	operties of Message Passing Mechanism	S
Addressing: how to id	dentify where a message should go	
Directionality:		
• simplex (one-v	way)	
• duplex (two-w	ray)	
• half-duplex (tw	wo-way, but only one way at a time)	
Message Boundaries:	:	
datagram model:	: message boundaries	
stream model: no	o boundaries	

<ul> <li>Connections: need to connect before communicating?</li> <li>in connection-oriented models, recipient is specified at time of connection, not by individual send operations. All messages sent over a connection have the same recipient.</li> <li>in connectionless models, recipient is specified as a parameter to each send operation.</li> <li>Reliability: <ul> <li>can messages get lost?</li> <li>can messages get reordered?</li> <li>can messages get damaged?</li> </ul> </li> </ul>		Properties of Message Passing Mechanisms (cont'd)
<ul> <li>in connection-oriented models, recipient is specified at time of connection, not by individual send operations. All messages sent over a connection have the same recipient.</li> <li>in connectionless models, recipient is specified as a parameter to each send operation.</li> <li>Reliability: <ul> <li>can messages get lost?</li> <li>can messages get reordered?</li> <li>can messages get damaged?</li> </ul> </li> </ul>	Cor	mections: need to connect before communicating?
<ul> <li>in connectionless models, recipient is specified as a parameter to each send operation.</li> <li>Reliability: <ul> <li>can messages get lost?</li> <li>can messages get reordered?</li> <li>can messages get damaged?</li> </ul> </li> </ul>		• in connection-oriented models, recipient is specified at time of connection, not by individual send operations. All messages sent over a connection have the same recipient.
Reliability: • can messages get lost? • can messages get reordered? • can messages get damaged?		• in connectionless models, recipient is specified as a parameter to each send operation.
<ul> <li>can messages get lost?</li> <li>can messages get reordered?</li> <li>can messages get damaged?</li> </ul>	Reli	iability:
<ul><li> can messages get reordered?</li><li> can messages get damaged?</li></ul>		• can messages get lost?
• can messages get damaged?		• can messages get reordered?
		• can messages get damaged?

Interprocess	Communication
--------------	---------------

	Sockets	
• a socket is a comr	nunication end-point	
• if two processes a	re to communicate, each process must crea	te its own socket
• two common type	es of sockets	
stream sockets: under the strea	support connection-oriented, reliable, duple am model (no message boundaries)	ex communication
datagram socket communicatio	s: support connectionless, best-effort (unre n under the datagram model (message bour	liable), duplex
• both types of sock	xets also support a variety of address domai	ns, e.g.,
Unix domain: us same machine	seful for communication between processes	running on the
INET domain: u	seful for communication between process	running on
different mach	ines that can communicate using IP protoc	ols.
350	Operating Systems	Fall 201

Interprocess Communication		6
Usi	ng Datagram Sockets (Receiver)	
s = socket(address	Type, SOCK_DGRAM);	
<pre>bind(s,address);</pre>		
<pre>recvfrom(s,buf,buf</pre>	Length,sourceAddress);	
close(s);		
<ul><li>socket creates a soci</li><li>bind assigns an addression</li></ul>	ket ess to the socket	
• recvfrom receives a	message from the socket	
- buf is a buffer to h	hold the incoming message	
- sourceAddress	s is a buffer to hold the address of the	message sender
• both buf and source	eAddress are filled by the recvfro	om call
C\$350	Operating Systems	Fall 2013

# Using Datagram Sockets (Sender)

<pre>s = socket(addressType, SOCK_DGRAM); sendto(s,buf,msgLength,targetAddress)</pre>
<pre>close(s);</pre>
• socket creates a socket
• sendto sends a message using the socket
- buf is a buffer that contains the message to be sent
- msgLength indicates the length of the message in the buffer
- targetAddress is the address of the socket to which the message is to
be delivered

CS350

Operating Systems

Fall 2013

	More on Datagram Sockets
sendto a	nd recvfrom calls <i>may</i> block
- recvf specifie	from blocks if there are no messages to be received from the ed socket
– sendt messag	o blocks if the system has no more room to buffer undelivered
• datagram s	ocket communications are (in general) unreliable
– messag	es (datagrams) may be lost
– messag	es may be reordered
• The sendir How does	ng process must know the address of the receive process's socket. it know this?

A Socket Address Convention

9

Service	Port	Description	
echo	 7/udp		
systat	11/tcp		
netstat	15/tcp		
chargen	19/udp		
ftp	21/tcp		
ssh	22/tcp	# SSH Remote Login Prot	cocol
telnet	23/tcp		
smtp	25/tcp		
time	37/udp		
gopher	70/tcp	# Internet Gopher	
finger	79/tcp		
WWW	80/tcp	# WorldWideWeb HTTP	
pop2	109/tcp	# POP version 2	
imap2	143/tcp	# IMAP	
CS350	OI	perating Systems	Fall 2013

```
Interprocess Communication
                                                                      10
                  Using Stream Sockets (Passive Process)
s = socket(addressType, SOCK_STREAM);
bind(s,address);
listen(s,backlog);
ns = accept(s,sourceAddress);
recv(ns,buf,bufLength);
send(ns,buf,bufLength);
. . .
close(ns); // close accepted connection
close(s);
             // don't accept more connections
 • listen specifies the number of connection requests for this socket that will be
   queued by the kernel
 • accept accepts a connection request and creates a new socket (ns)
 • recv receives up to bufLength bytes of data from the connection
 • send sends bufLength bytes of data over the connection.
CS350
                                                                  Fall 2013
                              Operating Systems
```

No	tes on Using Stream Sockets (Passive Process)	
• accept create	es a new socket (ns) for the new connection	
<ul> <li>sourceAddr socket that has</li> </ul>	mess is an address buffer. accept fills it with the made the connection request	address of the
• additional cont the original soo	nection requests can be accepted using more accepted (s)	pt calls on
• accept block	as if there are no pending connection requests	
• connection is d	luplex (both send and recv can be used)	
02250	Operating Systems	Eall 201

dress
nd)
pt call
_





## **One-way Child/Parent Communication Using a Simplex Pipe**

```
int fd[2];
char m[] = "message for parent";
char y[100];
pipe(fd); // create pipe
pid = fork(); // create child process
if (pid == 0) {
  // child executes this
  close(fd[0]); // close read end of pipe
  write(fd[1],m,19);
  . . .
} else {
  // parent executes this
  close(fd[1]); // close write end of pipe
  read(fd[0],y,19);
  . . .
}
CS350
                         Operating Systems
```



15

Fall 2013





	Examples of Other Interprocess Communication Mechanisms
nam	ed pipe:
	• similar to pipes, but with an associated name (usually a file name)
	• name allows arbitrary processes to communicate by opening the same named pipe
	• must be explicitly deleted, unlike an unnamed pipe
mess	sage queue:
	• like a named pipe, except that there are message boundaries
	• msgsend call sends a message into the queue, msgrecv call receives the next message from the queue

CS350

Operating Systems

Fall 2013

	Signals	
• signals permit asyr	nchronous one-way communication	
– from a process	to another process, or to a group of proc	esses, via the kernel
– from the kernel	l to a process, or to a group of processes	
• there are many typ	es of signals	
• the arrival of a sign receiving process	nal may cause the execution of a <i>signal h</i>	<i>andler</i> in the
• there may be a diff	ferent handler for each type of signal	

Examples	of	Signal	Types
----------	----	--------	-------

Signal	Value	Action	Comment
SIGINT	2	Term	Interrupt from keyboard
SIGILL	4	Core	Illegal Instruction
SIGKILL	9	Term	Kill signal
SIGCHLD	20,17,18	Ign	Child stopped or terminated
SIGBUS	10,7,10	Core	Bus error
SIGXCPU	24,24,30	Core	CPU time limit exceeded
SIGSTOP	17,19,23	Stop	Stop process

CS350

Operating Systems

Fall 2013



### **Implementing IPC**

- application processes use descriptors (identifiers) provided by the kernel to refer to specific sockets and pipes, as well as files and other objects
- kernel *descriptor tables* (or other similar mechanism) are used to associate descriptors with kernel data structures that implement IPC objects
- kernel provides bounded buffer space for data that has been sent using an IPC mechanism, but that has not yet been received
  - for IPC objects, like pipes, buffering is usually on a per object basis
  - IPC end points, like sockets, buffering is associated with each endpoint









### **IP Routing Table Example**

• Routing table for zonker.uwaterloo.ca, which is on three networks, and has IP addresses 129.97.74.66, 172.16.162.1, and 192.168.148.1 (one per network):

Destination	Gateway	Interface
172.16.162.*	-	vmnet1
129.97.74.*	-	eth0
192.168.148.*	-	vmnet8
default	129.97.74.1	eth0

• routing table key:

destination: ultimate destination of packet

**gateway:** next hop towards destination (or "-" if destination is directly reachable)

interface: which network interface to use to send this packet

CS350

Operating Systems

Fall 2013











Additional Notes:

Additional Notes:

Additional Notes:

Additional Notes: