

CS162 Operating Systems and Systems Programming Lecture 9

Synchronization 4: Monitors and Readers/Writers (Con't), Process Structure, Device Drivers

September 28th, 2020
Prof. John Kubiatowicz
<http://cs162.eecs.Berkeley.edu>

Recall: Atomic Read-Modify-Write

```

• test&set (&address) {           /* most architectures */
    result = M[address];         // return result from "address" and
    M[address] = 1;              // set value at "address" to 1
    return result;
}

• swap (&address, register) {    /* x86 */
    temp = M[address];           // swap register's value to
    M[address] = register;       // value at "address"
    register = temp;
}

• compare&swap (&address, reg1, reg2) { /* x86 (returns old value), 68000 */
    if (reg1 == M[address]) {    // If memory still == reg1,
        M[address] = reg2;      // then put reg2 => memory
        return success;
    } else {                    // Otherwise do not change memory
        return failure;
    }
}

• load-linked&store-conditional(&address) { /* R4000, alpha */
    loop:
        ll r1, M[address];
        movi r2, 1;             // Can do arbitrary computation
        sc r2, M[address];
        beqz r2, loop;
}

```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.2

Recall: Better Locks using test&set

- Can we build test&set locks without busy-waiting?
 - Mostly. Idea: only busy-wait to atomically check lock value

– `int guard = 0; // Global Variable!`
 – `int mylock = FREE; // Interface: acquire(&mylock);`
 `// release(&mylock);`



```

acquire(int *thelock) {
    // Short busy-wait time
    while (test&set(guard));
    if (*thelock == BUSY) {
        put thread on wait queue;
        go to sleep() & guard = 0;
        // guard == 0 on wakeup!
    } else {
        *thelock = BUSY;
        guard = 0;
    }
}

release(int *thelock) {
    // Short busy-wait time
    while (test&set(guard));
    if anyone on wait queue {
        take thread off wait queue
        Place on ready queue;
    } else {
        *thelock = FREE;
    }
    guard = 0;
}

```

- Note: sleep has to be sure to reset the guard variable
 - Why can't we do it just before or just after the sleep?

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.3

Recall: Linux futex: Fast Userspace Mutex

```

#include <linux/futex.h>
#include <sys/time.h>

int futex(int *uaddr, int futex_op, int val,
          const struct timespec *timeout);

```

`uaddr` points to a 32-bit value in user space

`futex_op`

- FUTEX_WAIT – if `val == *uaddr` sleep till FUTEX_WAIT
 - » **Atomic** check that condition still holds after we disable interrupts (in kernel!)
- FUTEX_WAKE – wake up at most `val` waiting threads
- FUTEX_FD, FUTEX_WAKE_OP, FUTEX_CMP_QUEUE: More interesting operations!

`timeout`

- ptr to a `timespec` structure that specifies a timeout for the op

- Interface to the kernel `sleep()` functionality!
 - Let thread put themselves to sleep – conditionally!
- futex is not exposed in libc; it is used within the implementation of pthreads**
 - Can be used to implement locks, semaphores, monitors, etc...

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.4

Recall: Lock Using Atomic Instructions and Futex

- Three (3) states:
 - **UNLOCKED**: No one has lock
 - **LOCKED**: One thread has lock
 - **CONTESTED**: Possibly more than one (with someone sleeping)
- Clean interface!
- Lock grabbed cleanly by either
 - `compare_and_swap()`
 - First `swap()`
- No overhead if uncontested!
- Could build semaphores in a similar way!

```
typedef enum { UNLOCKED, LOCKED, CONTESTED } Lock;
Lock mylock = UNLOCKED; // Interface: acquire(&mylock);
                          //                      release(&mylock);

acquire(Lock *thelock) {
    // If unlocked, grab lock!
    if (compare_and_swap(thelock, UNLOCKED, LOCKED))
        return;

    // Keep trying to grab lock, sleep in futex
    while (swap(mylock, CONTESTED) != UNLOCKED)
        // Sleep unless someone releases hear!
        futex(thelock, FUTEX_WAIT, CONTESTED);
}

release(Lock *thelock) {
    // If someone sleeping,
    if (swap(thelock, UNLOCKED) == CONTESTED)
        futex(thelock, FUTEX_WAKE, 1);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.5

Recall: Monitors and Condition Variables

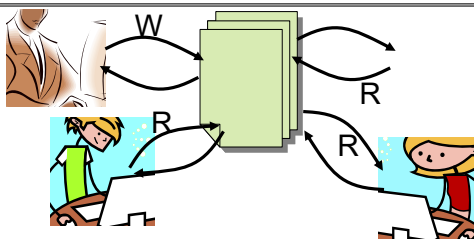
- **Monitor**: a lock and zero or more condition variables for managing concurrent access to shared data
 - Use of Monitors is a programming paradigm
 - Some languages like Java provide monitors in the language
- **Condition Variable**: a queue of threads waiting for something *inside* a critical section
 - Key idea: allow sleeping inside critical section by atomically releasing lock at time we go to sleep
 - Contrast to semaphores: Can't wait inside critical section
- Operations:
 - `Wait(&lock)`: Atomically release lock and go to sleep. Re-acquire lock later, before returning.
 - `Signal()`: Wake up one waiter, if any
 - `Broadcast()`: Wake up all waiters
- Rule: **Must hold lock when doing condition variable ops!**

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.6

Recall: Readers/Writers Problem



- Motivation: Consider a shared database
 - Two classes of users:
 - » Readers – never modify database
 - » Writers – read and modify database
 - Is using a single lock on the whole database sufficient?
 - » Like to have many readers at the same time
 - » Only one writer at a time

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.7

Recall: Structure of Mesa Monitor Program

- Monitors represent the synchronization logic of the program
 - Wait if necessary
 - Signal when change something so any waiting threads can proceed
- Basic structure of mesa monitor-based program:


```
lock
while (need to wait) {
    condvar.wait();
}
unlock

do something so no need to wait

lock

condvar.signal();

unlock
```

 - Check and/or update state variables (for the first while loop)
 - Wait if necessary (for the first while loop)
 - Check and/or update state variables (for the signal call)

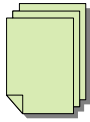
9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.8

Recall: Basic Readers/Writers Solution

- Correctness Constraints:
 - Readers can access database when no writers
 - Writers can access database when no readers or writers
 - Only one thread manipulates state variables at a time
- Basic structure of a solution:
 - Reader()
 - Wait until no writers
 - Access data base
 - Check out – wake up a waiting writer
 - Writer()
 - Wait until no active readers or writers
 - Access database
 - Check out – wake up waiting readers or writer
 - State variables (Protected by a lock called “lock”):
 - » int AR: Number of active readers; initially = 0
 - » int WR: Number of waiting readers; initially = 0
 - » int AW: Number of active writers; initially = 0
 - » int WW: Number of waiting writers; initially = 0
 - » Condition okToRead = NIL
 - » Condition okToWrite = NIL



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.9

Recall: Code for a Reader

```
Reader() {
    // First check self into system
    acquire(&lock);

    while ((AW + WW) > 0) { // Is it safe to read?
        WR++;              // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--;              // No longer waiting
    }

    AR++;                  // Now we are active!
    release(&lock);

    // Perform actual read-only access
    AccessDatabase(ReadOnly);

    // Now, check out of system
    acquire(&lock);
    AR--;                  // No longer active
    if (AR == 0 && WW > 0) // No other active readers
        cond_signal(&okToWrite); // Wake up one writer
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.10

Recall: Code for a Writer

```
Writer() {
    // First check self into system
    acquire(&lock);

    while ((AW + AR) > 0) { // Is it safe to write?
        WW++;              // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--;              // No longer waiting
    }

    AW++;                  // Now we are active!
    release(&lock);

    // Perform actual read/write access
    AccessDatabase(ReadWrite);

    // Now, check out of system
    acquire(&lock);
    AW--;                  // No longer active
    if (WW > 0) {           // Give priority to writers
        cond_signal(&okToWrite); // Wake up one writer
    } else if (WR > 0) {    // Otherwise, wake reader
        cond_broadcast(&okToRead); // Wake all readers
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.11

Simulation of Readers/Writers Solution

- Use an example to simulate the solution
- Consider the following sequence of operators:
 - R1, R2, W1, R3
- Initially: AR = 0, WR = 0, AW = 0, WW = 0

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.12

Simulation of Readers/Writers Solution

- R1 comes along (no waiting threads)
- AR = 0, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.13

Simulation of Readers/Writers Solution

- R1 comes along (no waiting threads)
- AR = 0, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.14

Simulation of Readers/Writers Solution

- R1 comes along (no waiting threads)
- AR = 1, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.15

Simulation of Readers/Writers Solution

- R1 comes along (no waiting threads)
- AR = 1, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.16

Simulation of Readers/Writers Solution

- R1 accessing dbase (no other threads)
- AR = 1, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.17

Simulation of Readers/Writers Solution

- R2 comes along (R1 accessing dbase)
- AR = 1, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.18

Simulation of Readers/Writers Solution

- R2 comes along (R1 accessing dbase)
- AR = 1, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.19

Simulation of Readers/Writers Solution

- R2 comes along (R1 accessing dbase)
- AR = 2, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.20

Simulation of Readers/Writers Solution

- R2 comes along (R1 accessing dbase)
- AR = 2, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

AccessDBase(ReadOnly);

Simulation of Readers/Writers Solution

- R1 and R2 accessing dbase
- AR = 2, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
```

Assume readers take a while to access database
Situation: Locks released, only AR is non-zero

Simulation of Readers/Writers Solution

- W1 comes along (R1 and R2 are still accessing dbase)
- AR = 2, WR = 0, AW = 0, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond signal(&okToWrite);
    } else if (WR > 0) {
        cond broadcast(&okToRead);
    }
    release(&lock);
}
```

AccessDBase(ReadWrite);

Simulation of Readers/Writers Solution

- W1 comes along (R1 and R2 are still accessing dbase)
- AR = 2, WR = 0, AW = 0, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond signal(&okToWrite);
    } else if (WR > 0) {
        cond broadcast(&okToRead);
    }
    release(&lock);
}
```

AccessDBase(ReadWrite);

Simulation of Readers/Writers Solution

- W1 comes along (R1 and R2 are still accessing dbase)
- AR = 2, WR = 0, AW = 0, WW = 1

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.25

Simulation of Readers/Writers Solution

- R3 comes along (R1 and R2 accessing dbase, W1 waiting)
- AR = 2, WR = 0, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0) {
        cond_signal(&okToWrite);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.26

Simulation of Readers/Writers Solution

- R3 comes along (R1 and R2 accessing dbase, W1 waiting)
- AR = 2, WR = 0, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0) {
        cond_signal(&okToWrite);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.27

Simulation of Readers/Writers Solution

- R3 comes along (R1 and R2 accessing dbase, W1 waiting)
- AR = 2, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    lock.release();

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0) {
        cond_signal(&okToWrite);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.28

Simulation of Readers/Writers Solution

- R3 comes along (R1, R2 accessing dbase, W1 waiting)
- AR = 2, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.29

Simulation of Readers/Writers Solution

- R1 and R2 accessing dbase, W1 and R3 waiting
- AR = 2, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
```

Status:

- R1 and R2 still reading
- W1 and R3 waiting on okToWrite and okToRead, respectively

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.30

Simulation of Readers/Writers Solution

- R2 finishes (R1 accessing dbase, W1 and R3 waiting)
- AR = 2, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.31

Simulation of Readers/Writers Solution

- R2 finishes (R1 accessing dbase, W1 and R3 waiting)
- AR = 1, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.32

Simulation of Readers/Writers Solution

- R2 finishes (R1 accessing dbase, W1 and R3 waiting)
- AR = 1, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.33

Simulation of Readers/Writers Solution

- R2 finishes (R1 accessing dbase, W1 and R3 waiting)
- AR = 1, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.34

Simulation of Readers/Writers Solution

- R1 finishes (W1 and R3 waiting)
- AR = 1, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.35

Simulation of Readers/Writers Solution

- R1 finishes (W1, R3 waiting)
- AR = 0, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.36

Simulation of Readers/Writers Solution

- R1 finishes (W1, R3 waiting)
- AR = 0, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.37

Simulation of Readers/Writers Solution

- R1 signals a writer (W1 and R3 waiting)
- AR = 0, WR = 1, AW = 0, WW = 1

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.38

Simulation of Readers/Writers Solution

- W1 gets signal (R3 still waiting)
- AR = 0, WR = 1, AW = 0, WW = 1

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else-if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.39

Simulation of Readers/Writers Solution

- W1 gets signal (R3 still waiting)
- AR = 0, WR = 1, AW = 0, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else-if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.40

Simulation of Readers/Writers Solution

- W1 gets signal (R3 still waiting)
- AR = 0, WR = 1, AW = 1, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.41

Simulation of Readers/Writers Solution

- W1 accessing dbase (R3 still waiting)
- AR = 0, WR = 1, AW = 1, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.42

Simulation of Readers/Writers Solution

- W1 finishes (R3 still waiting)
- AR = 0, WR = 1, AW = 1, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.43

Simulation of Readers/Writers Solution

- W1 finishes (R3 still waiting)
- AR = 0, WR = 1, AW = 0, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.44

Simulation of Readers/Writers Solution

- W1 finishes (R3 still waiting)
- AR = 0, WR = 1, AW = 0, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.45

Simulation of Readers/Writers Solution

- W1 signaling readers (R3 still waiting)
- AR = 0, WR = 1, AW = 0, WW = 0

```
Writer() {
    acquire(&lock);
    while ((AW + AR) > 0) { // Is it safe to write?
        WW++; // No. Active users exist
        cond_wait(&okToWrite, &lock); // Sleep on cond var
        WW--; // No longer waiting
    }
    AW++;
    release(&lock);

    AccessDBase(ReadWrite);

    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okToWrite);
    } else if (WR > 0) {
        cond_broadcast(&okToRead);
    }
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.46

Simulation of Readers/Writers Solution

- R3 gets signal (no waiting threads)
- AR = 0, WR = 1, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.47

Simulation of Readers/Writers Solution

- R3 gets signal (no waiting threads)
- AR = 0, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okToWrite);
    release(&lock);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.48

Simulation of Readers/Writers Solution

- R3 accessing dbase (no waiting threads)
- AR = 1, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okToWrite);
    release(&lock);
}
```

AccessDBase(ReadOnly);

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.49

Simulation of Readers/Writers Solution

- R3 finishes (no waiting threads)
- AR = 1, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okToWrite);
    release(&lock);
}
```

AccessDBase(ReadOnly);

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.50

Simulation of Readers/Writers Solution

- R3 finishes (no waiting threads)
- AR = 0, WR = 0, AW = 0, WW = 0

```
Reader() {
    acquire(&lock);
    while ((AW + WW) > 0) { // Is it safe to read?
        WR++; // No. Writers exist
        cond_wait(&okToRead, &lock); // Sleep on cond var
        WR--; // No longer waiting
    }
    AR++; // Now we are active!
    release(&lock);

    AccessDBase(ReadOnly);

    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okToWrite);
    release(&lock);
}
```

AccessDBase(ReadOnly);

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.51

Questions

- Can readers starve? Consider Reader() entry code:


```
while ((AW + WW) > 0) { // Is it safe to read?
    WR++; // No. Writers exist
    cond_wait(&okToRead, &lock); // Sleep on cond var
    WR--; // No longer waiting
}
AR++; // Now we are active!
```
- What if we erase the condition check in Reader exit?


```
AR--; // No longer active
if (AR == 0 && WW > 0) // No other active readers
    cond_signal(&okToWrite); // Wake up one writer
```
- Further, what if we turn the signal() into broadcast()


```
AR--; // No longer active
cond_broadcast(&okToWrite); // Wake up sleepers
```
- Finally, what if we use only one condition variable (call it "okContinue") instead of two separate ones?
 - Both readers and writers sleep on this variable
 - Must use broadcast() instead of signal()

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.52

Use of Single CV: okContinue

```

Reader() {
    // check into system
    acquire(&lock);
    while ((AW + WW) > 0) {
        WR++;
        cond_wait(&okContinue,&lock);
        WR--;
    }
    AR++;
    release(&lock);

    // read-only access
    AccessDbase(ReadOnly);

    // check out of system
    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okContinue);
    release(&lock);
}

Writer() {
    // check into system
    acquire(&lock);
    while ((AW + AR) > 0) {
        WW++;
        cond_wait(&okContinue,&lock);
        WW--;
    }
    AW++;
    release(&lock);

    // read/write access
    AccessDbase(ReadWrite);

    // check out of system
    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okContinue);
    } else if (WR > 0) {
        cond_broadcast(&okContinue);
    }
    release(&lock);
}

```

What if we turn okToWrite and okToRead into okContinue
(i.e. use only one condition variable instead of two)?

9/28/20

Lec 9.53

Use of Single CV: okContinue

```

Reader() {
    // check into system
    acquire(&lock);
    while ((AW + WW) > 0) {
        WR++;
        cond_wait(&okContinue,&lock);
        WR--;
    }
    AR++;
    release(&lock);

    // read-only access
    AccessDbase(ReadOnly);

    // check out of system
    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_signal(&okContinue);
    release(&lock);
}

Writer() {
    // check into system
    acquire(&lock);
    while ((AW + AR) > 0) {
        WW++;
        cond_wait(&okContinue,&lock);
        WW--;
    }
    AW++;
    release(&lock);

    // read/write access
    AccessDbase(ReadWrite);

    // check out of system
    acquire(&lock);
    AW--;
    if (WW > 0) {
        cond_signal(&okContinue);
    } else if (WR > 0) {
        cond_broadcast(&okContinue);
    }
    release(&lock);
}

```

Consider this scenario:

- R1 arrives
- W1, R2 arrive while R1 still reading → W1 and R2 wait for R1 to finish
- Assume R1's signal is delivered to R2 (not W1)

9/28/20

Lec 9.54

Use of Single CV: okContinue

```

Reader() {
    // check into system
    acquire(&lock);
    while ((AW + WW) > 0) {
        WR++;
        cond_wait(&okContinue,&lock);
        WR--;
    }
    AR++;
    release(&lock);

    // read-only access
    AccessDbase(ReadOnly);

    // check out of system
    acquire(&lock);
    AR--;
    if (AR == 0 && WW > 0)
        cond_broadcast(&okContinue);
    release(&lock);
}

Writer() {
    // check into system
    acquire(&lock);
    while ((AW + AR) > 0) {
        WW++;
        cond_wait(&okContinue,&lock);
        WW--;
    }
    AW++;
    release(&lock);

    // read/write access
    AccessDbase(ReadWrite);

    // check out of system
    acquire(&lock);
    AW--;
    if (WW > 0 || WR > 0) {
        cond_broadcast(&okContinue);
    }
    release(&lock);
}

```

Need to change to
broadcast()!

Must broadcast()
to sort things out!

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.55

Administrivia

- Midterm 1: Thursday (October 1st): 5-7PM
 - We understand that this partially conflicts with CS170, but those of you in CS170 can start that exam after 7PM (according to CS170 staff)
 - Video Proctored, No curve, Use of computer to answer questions
 - More details as we get closer to exam
- Midterm topics:
 - Everything from lecture up to today's lecture
 - Scheduling is not part of exam...
 - Homework and Project work is fair game
- Midterm Review: Tomorrow, 7-9pm
 - Details TBA

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.56

Can we construct Monitors from Semaphores?

- Locking aspect is easy: Just use a mutex
- Can we implement condition variables this way?

```
Wait(Semaphore *thesema) { semaP(thesema); }
Signal(Semaphore *thesema) { semaV(thesema); }
```
- Does this work better?

```
Wait(Lock *thelock, Semaphore *thesema) {
    release(thelock);
    semaP(thesema);
    acquire(thelock);
}
Signal(Semaphore *thesema) {
    semaV(thesema);
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.57

Construction of Monitors from Semaphores (con't)

- Problem with previous try:
 - P and V are commutative – result is the same no matter what order they occur
 - Condition variables are NOT commutative
- Does this fix the problem?

```
Wait(Lock *thelock, Semaphore *thesema) {
    release(thelock);
    semaP(thesema);
    acquire(thelock);
}
Signal(Semaphore *thesema) {
    if semaphore queue is not empty
        semaV(thesema);
}
```

 - Not legal to look at contents of semaphore queue
 - There is a race condition – signaler can slip in after lock release and before waiter executes semaphore.P()
- It is actually possible to do this correctly
 - Complex solution for Hoare scheduling in book
 - Can you come up with simpler Mesa-scheduled solution?

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.58

Mesa Monitor Conclusion

- Monitors represent the synchronization logic of the program
 - Wait if necessary
 - Signal when change something so any waiting threads can proceed
- Typical structure of monitor-based program:

```
lock
while (need to wait) {
    condvar.wait();
}
unlock

do something so no need to wait

lock
condvar.signal();
unlock
```

} Check and/or update state variables
Wait if necessary

} Check and/or update state variables

9/28/20

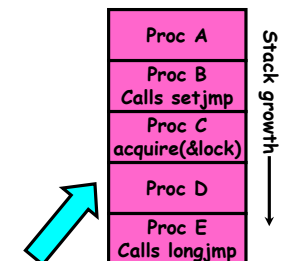
Kubiatowicz CS162 © UCB Fall 2020

Lec 9.59

C-Language Support for Synchronization

- C language: Pretty straightforward synchronization
 - Just make sure you know *all* the code paths out of a critical section

```
int Rtn() {
    acquire(&lock);
    ...
    if (exception) {
        release(&lock);
        return errReturnCode;
    }
    ...
    release(&lock);
    return OK;
}
```
 - Watch out for `setjmp/longjmp`!
 - » Can cause a non-local jump out of procedure
 - » In example, procedure E calls `longjmp`, popping stack back to procedure B
 - » If Procedure C had `lock.acquire`, problem!



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.60

Concurrency and Synchronization in C

- Harder with more locks

```
void Rtn() {
    lock1.acquire();
    if (error) {
        lock1.release();
        return;
    }
    ...
    lock2.acquire();
    ...
    if (error) {
        lock2.release();
        lock1.release();
        return;
    }
    ...
    lock2.release();
    lock1.release();
}
```

- Is goto a solution???

```
void Rtn() {
    lock1.acquire();
    if (error) {
        goto release_lock1_and_return;
    }
    ...
    lock2.acquire();
    ...
    if (error) {
        goto release_both_and_return;
    }
    ...
release_both_and_return:
    lock2.release();
release_lock1_and_return:
    lock1.release();
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.61

C++ Language Support for Synchronization

- Languages with exceptions like C++

- Languages that support exceptions are problematic (easy to make a non-local exit without releasing lock)

- Consider:

```
void Rtn() {
    lock.acquire();
    ...
    DoFoo();
    ...
    lock.release();
}
void DoFoo() {
    ...
    if (exception) throw errException;
    ...
}
```

- Notice that an exception in DoFoo() will exit without releasing the lock!

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.62

C++ Language Support for Synchronization (con't)

- Must catch all exceptions in critical sections

- Catch exceptions, release lock, and re-throw exception:

```
void Rtn() {
    lock.acquire();
    try {
        ...
        DoFoo();
        ...
    } catch (...) { // catch exception
        lock.release(); // release lock
        throw; // re-throw the exception
    }
    lock.release();
}
void DoFoo() {
    ...
    if (exception) throw errException;
    ...
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.63

Much better: C++ Lock Guards

```
#include <mutex>
int global_i = 0;
std::mutex global_mutex;

void safe_increment() {
    std::lock_guard<std::mutex> lock(global_mutex);
    ...
    global_i++;
    // Mutex released when 'lock' goes out of scope
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.64

Python with Keyword

- More versatile than we show here (can be used to close files, database connections, etc.)

```
lock = threading.Lock()
...
with lock: # Automatically calls acquire()
    some_var += 1
...
# release() called however we leave block
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.65

Java synchronized Keyword

- Every Java object has an associated lock:
 - Lock is acquired on entry and released on exit from a **synchronized** method
 - Lock is properly released if exception occurs inside a **synchronized** method
 - Mutex execution of synchronized methods (beware deadlock)

```
class Account {
    private int balance;
    // object constructor
    public Account (int initialBalance) {
        balance = initialBalance;
    }
    public synchronized int getBalance() {
        return balance;
    }
    public synchronized void deposit(int amount) {
        balance += amount;
    }
}
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.66

Java Support for Monitors

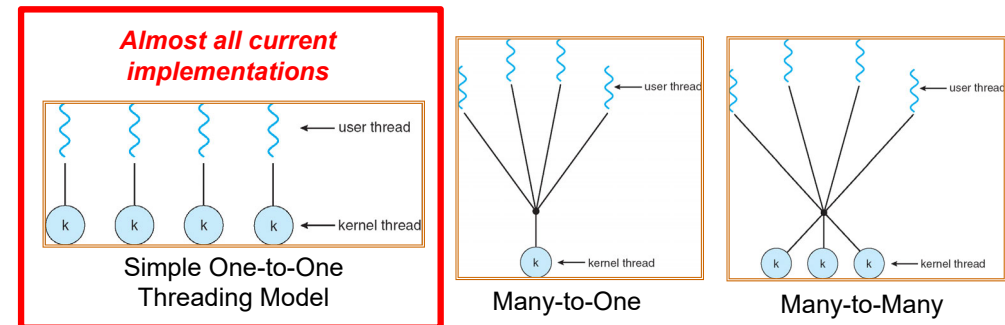
- Along with a lock, every object has a single condition variable associated with it
- To wait inside a synchronized method:
 - void wait();
 - void wait(long timeout);
- To signal while in a synchronized method:
 - void notify();
 - void notifyAll();

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.67

Recall: User/Kernel Threading Models



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.68

Recall: Thread State in the Kernel

- For every thread in a process, the kernel maintains:
 - The thread's TCB
 - A kernel stack used for syscalls/interrupts/traps
 - This kernel-state is sometimes called the “**kernel thread**”
 - The “kernel thread” is suspended (but ready to go) when thread is running in user-space
- Additionally, some threads just do work in the kernel
 - Still has TCB
 - Still has kernel stack
 - But not part of any process, and never executes in user mode

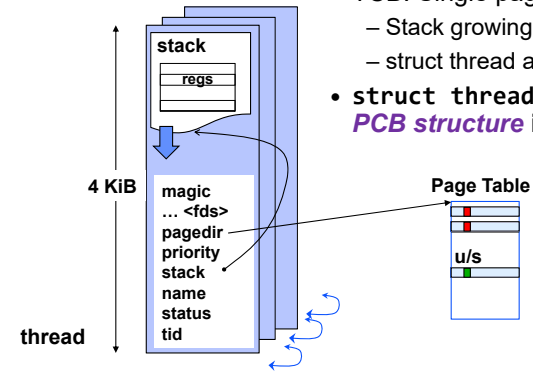
9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.69

In Pintos, Processes are Single-Threaded

- Pintos processes have only one thread
- TCB: Single page (4 KiB)
 - Stack growing from the top (high addresses)
 - struct thread at the bottom (low addresses)
- struct thread** defines the TCB structure *and* **PCB structure** in Pintos



Pintos: thread.c

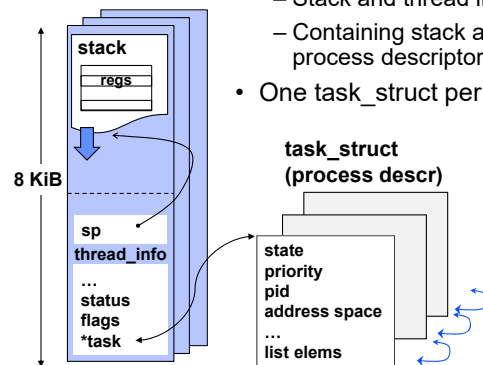
9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.70

(Aside): Linux “Task”

- Linux “Kernel Thread”: 2 pages (8 KiB)
 - Stack and thread information on opposite sides
 - Containing stack and thread information + process descriptor
- One task_struct per thread



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.71

Multithreaded Processes (not in Pintos)

- Traditional implementation strategy:
 - One PCB (process struct) per process
 - Each PCB contains (or stores pointers to) each thread's TCB
- Linux's strategy:
 - One task_struct per thread
 - Threads belonging to the same process happen to share some resources
 - Like address space, file descriptor table, etc.
- To what extent does this actually matter?**

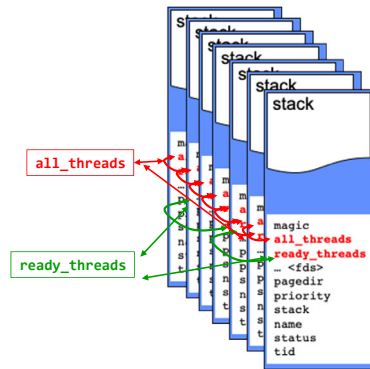
9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.72

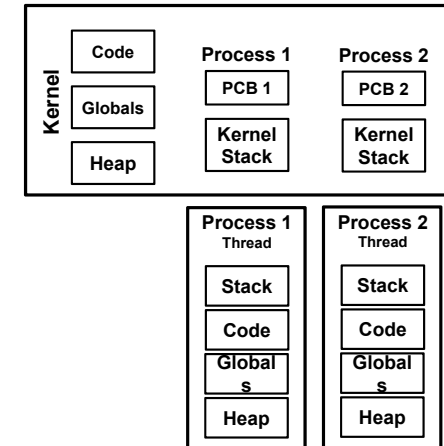
Aside: Polymorphic Linked Lists in C

- Many places in the kernel need to maintain a "list of X"
 - This is tricky in C, which has no polymorphism
 - Essentially adding an *interface* to a package
- In Linux and Pintos this is done by embedding a `list_elem` in the struct
 - Macros allow shift of view between object and list
 - You saw this in Homework 1

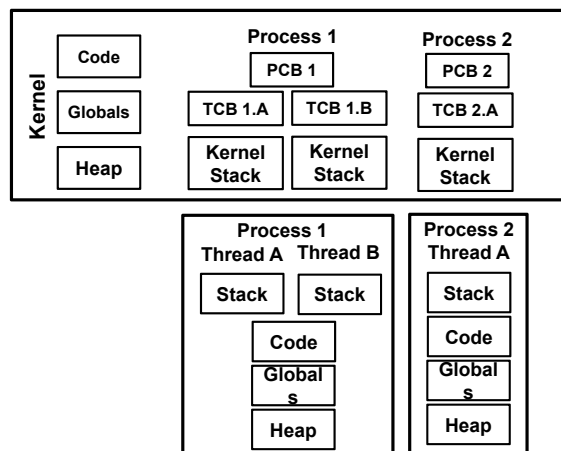


Pintos: `list.c`

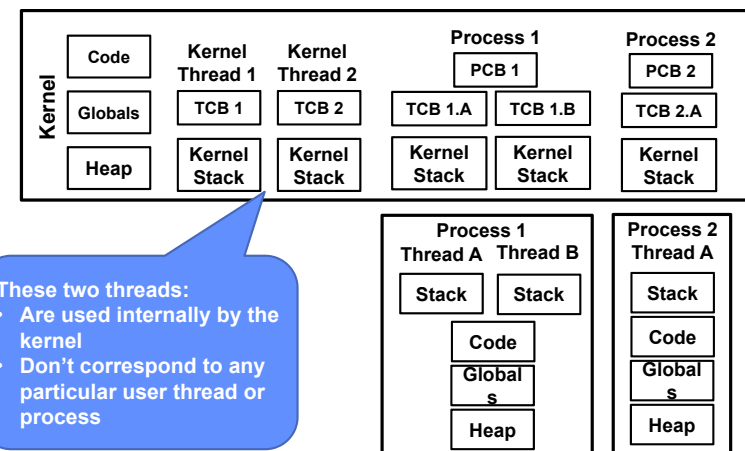
Kernel Structure So Far (1/3)



Kernel Structure So Far (2/3)



Kernel Structure So Far (3/3)



These two threads:

- Are used internally by the kernel
- Don't correspond to any particular user thread or process

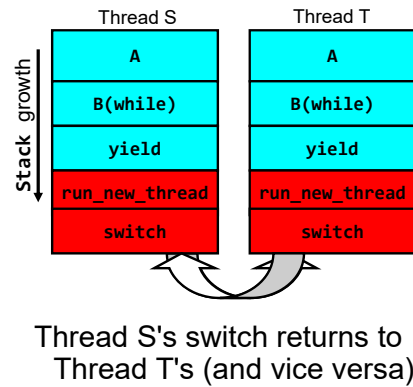
Recall: Multithreaded Stack Example

- Consider the following code blocks:

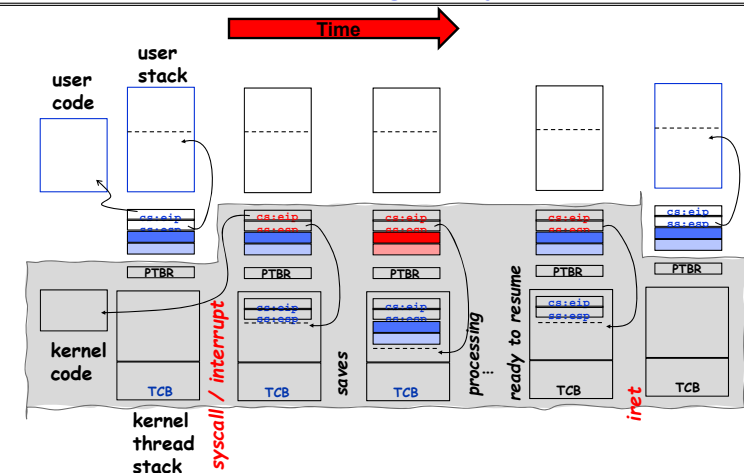
```

proc A() {
    B();
}
proc B() {
    while(TRUE) {
        yield();
        run_new_thread();
    }
}
    
```

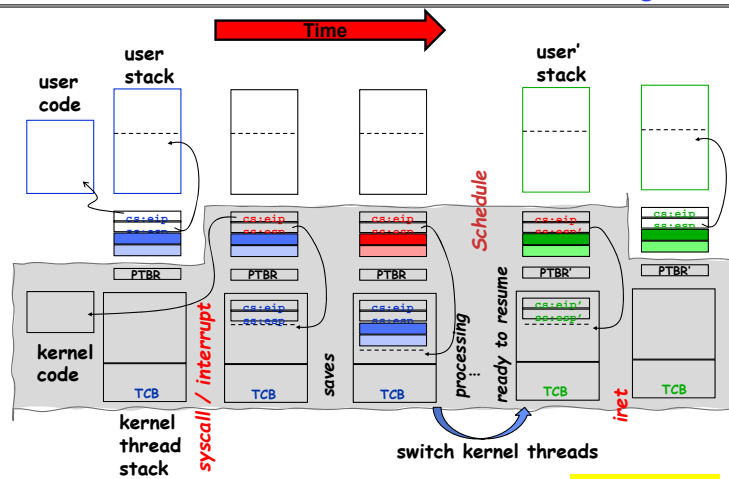
- Suppose we have 2 threads:
 - Threads S and T



Recall: Kernel Crossing on Syscall or Interrupt

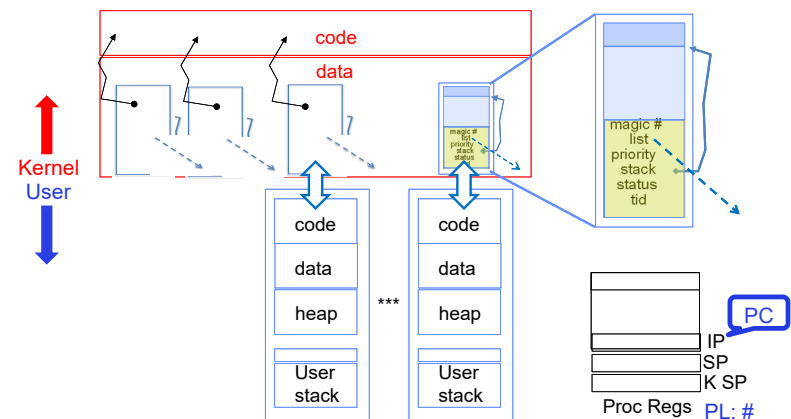


Recall: Context Switch – Scheduling



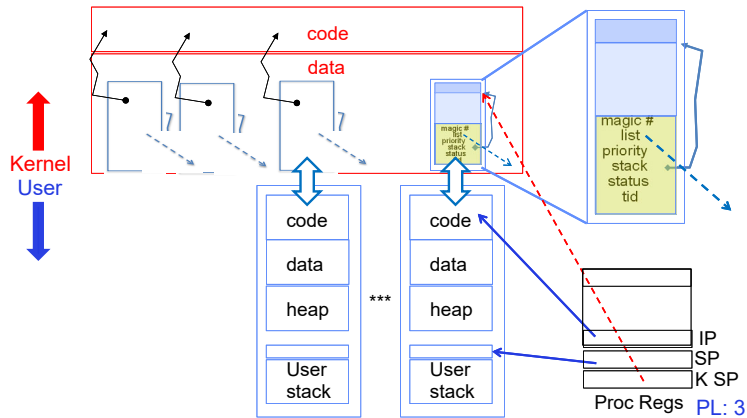
Pintos: switch.S

MT Kernel 1T Process ala Pintos/x86



- Each user process/thread associated with a kernel thread, described by a 4KB page object containing TCB and kernel stack for the kernel thread

In User thread, w/ Kernel thread waiting



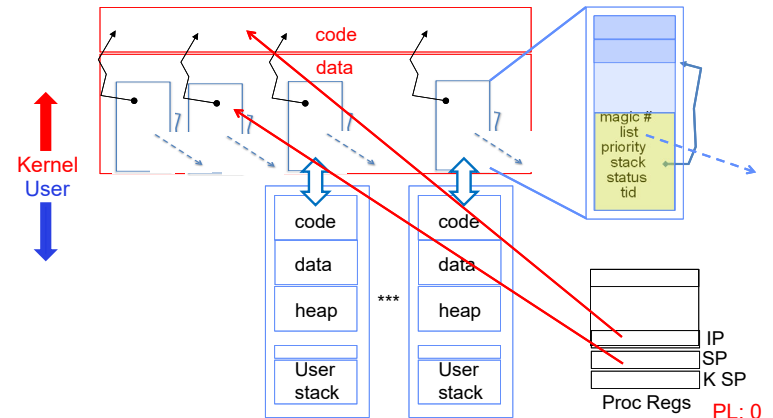
- x86 CPU holds interrupt SP in register
- During user thread execution, associated kernel thread is “standing by”

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.81

In Kernel Thread: No User Component



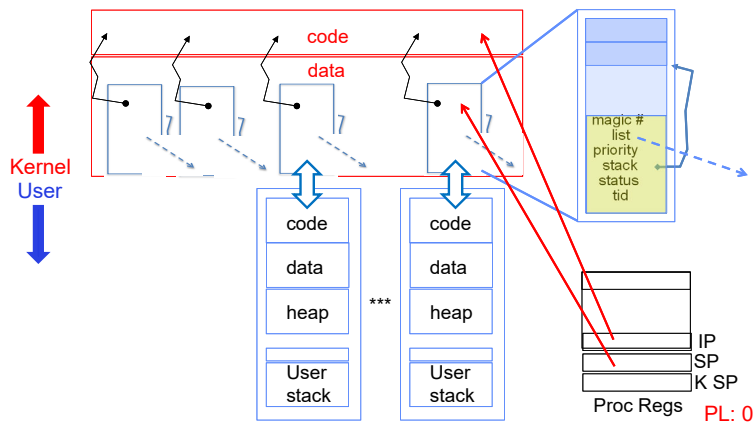
- Kernel threads execute with small stack in thread structure
- Pure kernel threads have no corresponding user-mode thread

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.82

User → Kernel (exceptions, syscalls)



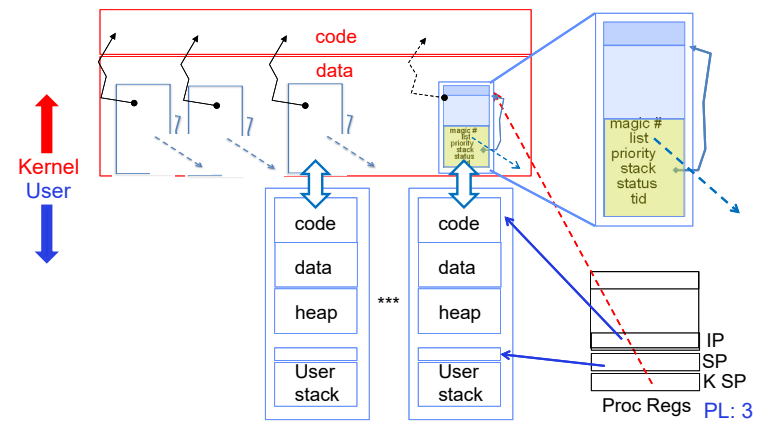
- Mechanism to resume k-thread goes through interrupt vector

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.83

Kernel → User



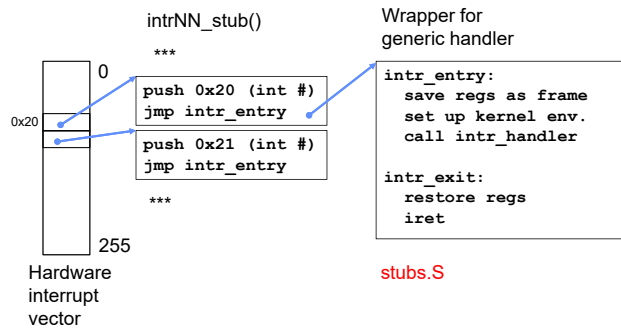
- Interrupt return (iret) restores user stack, IP, and PL

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.84

Pintos Interrupt Processing

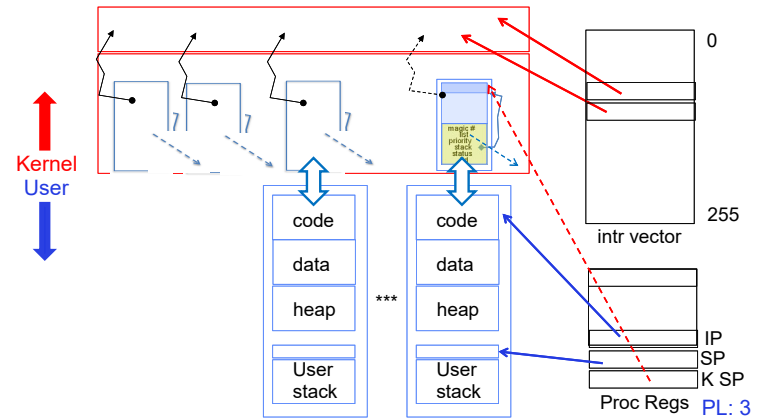


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.85

User → Kernel via interrupt vector



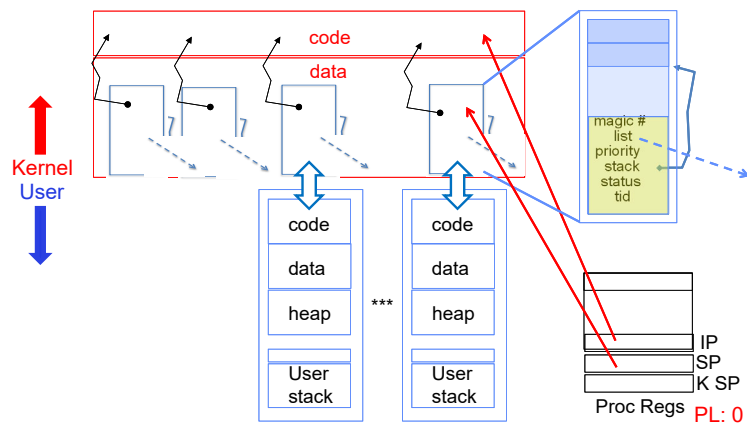
- Interrupt transfers control through the Interrupt Vector (IDT in x86)
- `iret` restores user stack and priority level (PL)

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.86

Switch to Kernel Thread for Process

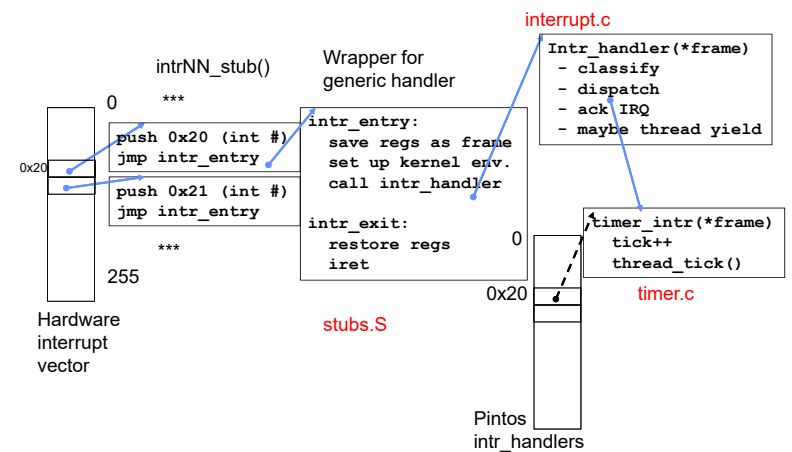


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.87

Pintos Interrupt Processing



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.88

Timer may trigger thread switch

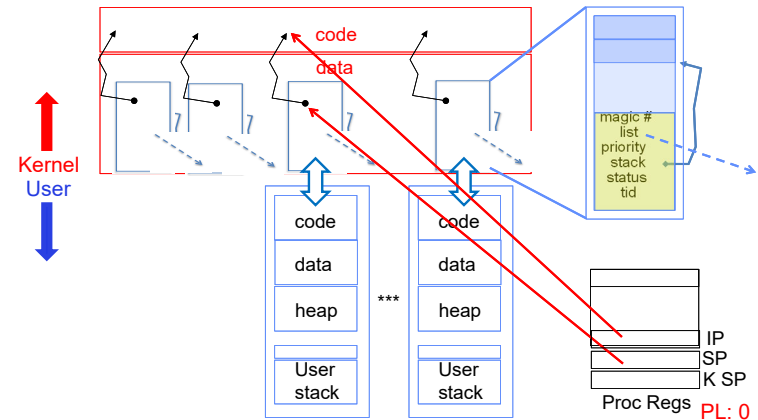
- `thread_tick`
 - Updates thread counters
 - If quanta exhausted, sets yield flag
- `thread_yield`
 - On path to rtn from interrupt
 - Sets current thread back to READY
 - Pushes it back on `ready_list`
 - Calls `schedule` to select next thread to run upon `iret`
- **Schedule (Next Lecture!)**
 - Selects next thread to run
 - Calls `switch_threads` to change regs to point to stack for thread to resume
 - Sets its status to RUNNING
 - If user thread, activates the process
 - Returns back to `intr_handler`

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.89

Thread Switch (switch.S)



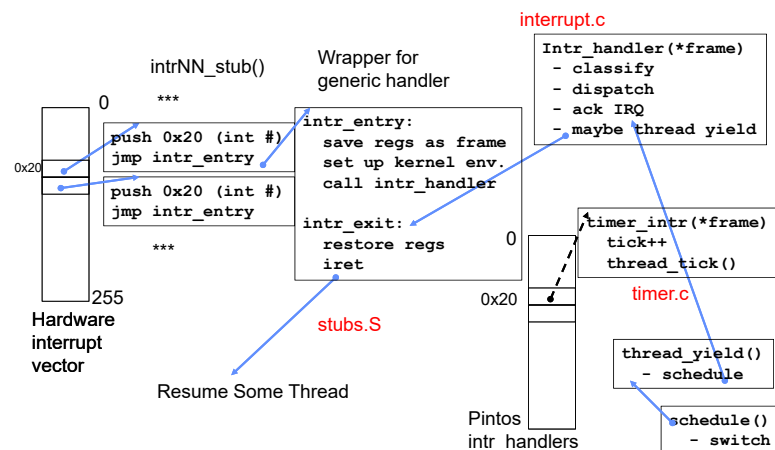
- `switch_threads`: save regs on current small stack, change SP, return from destination threads call to `switch_threads`

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.90

Pintos Return from Processing

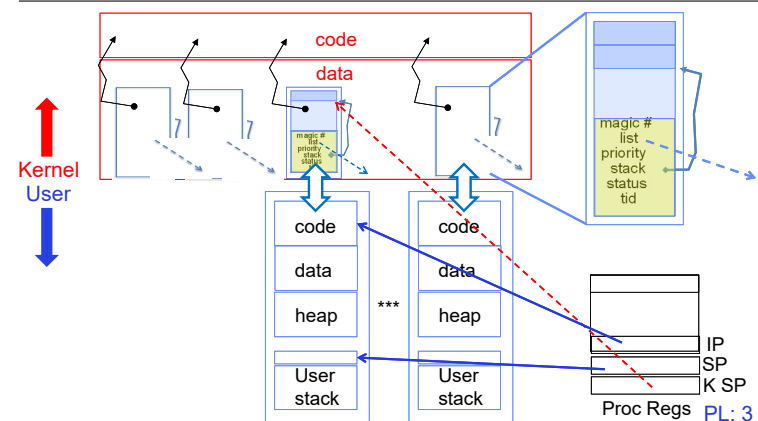


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.91

Kernel → Different User Thread



- `iret` restores user stack and priority level (PL)

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.92

Famous Quote WRT Scheduling: Dennis Richie

Dennis Richie,
Unix V6, slp.c:

```

2230 /*
2231  * If the new process paused because it was
2232  * swapped out, set the stack level to the last call
2233  * to savu(u_ssav). This means that the return
2234  * which is executed immediately after the call to aretu
2235  * actually returns from the last routine which did
2236  * the savu.
2237  *
2238  * You are not expected to understand this.
2239  */
    
```

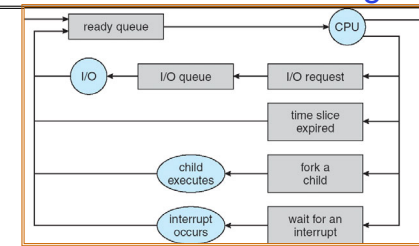
"If the new process paused because it was swapped out, set the stack level to the last call to savu(u_ssav). This means that the return which is executed immediately after the call to aretu actually returns from the last routine which did the savu."

"You are not expected to understand this."

Source: Dennis Ritchie, Unix V6 slp.c (context-switching code) as per The Unix Heritage Society(tuhs.org); gif by Eddie Koehler.

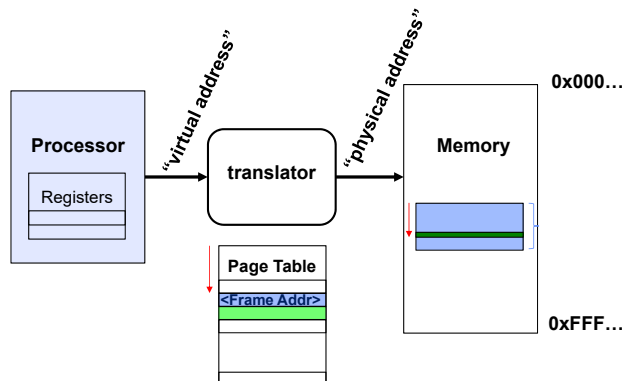
Included by Ali R. Butt in CS3204 from Virginia Tech

Recall: Scheduling



- Question: How is the OS to decide which of several tasks to take off a queue?
- **Scheduling**: deciding which threads are given access to resources from moment to moment
 - Often, we think in terms of CPU time, but could also think about access to resources like network BW or disk access
- Next time: we dive into scheduling!

Recall: Address Space



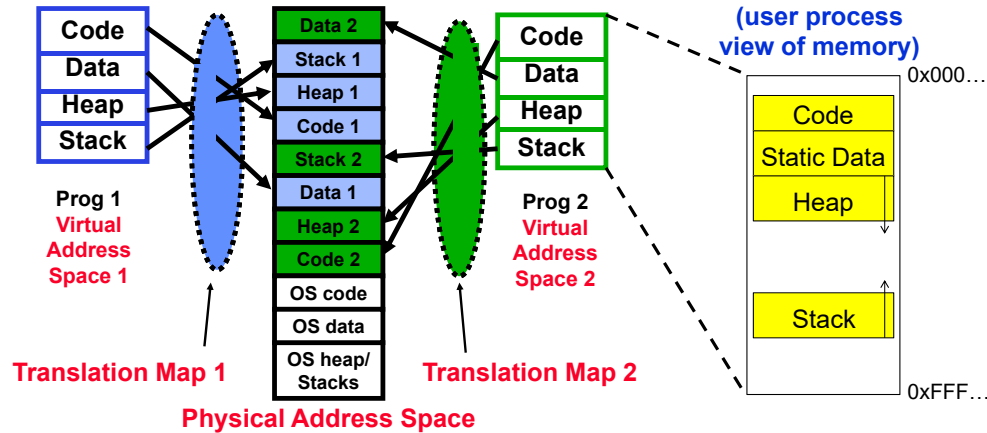
- Program operates in an address space that is distinct from the physical memory space of the machine

Understanding "Address Space"

- Page table is the primary mechanism
- Privilege Level determines which regions can be accessed
 - Which entries can be used
- System (PL=0) can access all, User (PL=3) only part
- Each process has its own address space
- The "System" part of all of them is the same

All system threads share the same system address space and same memory

Page Table Mapping (Rough Idea)

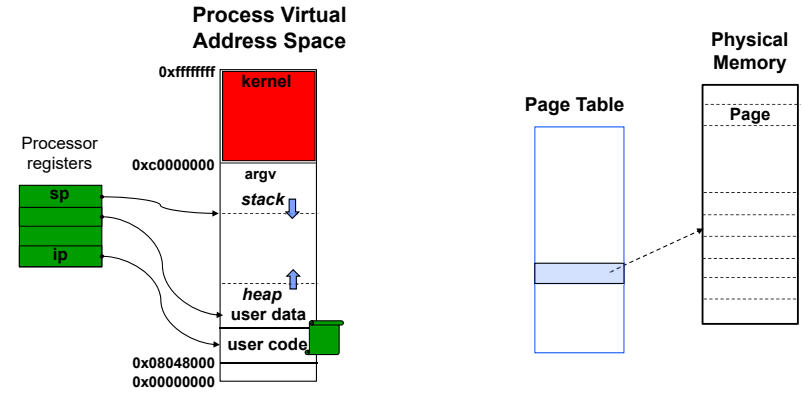


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.97

User Process View of Memory

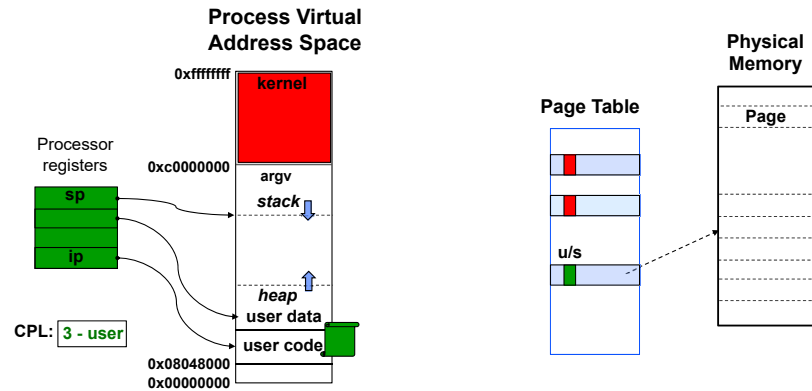


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.98

Processor Mode (Privilege Level)

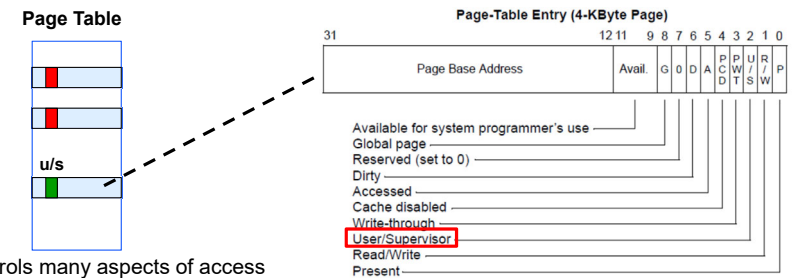


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.99

Aside: x86 (32-bit) Page Table Entry



- Controls many aspects of access
- Later – discuss page table organization
 - For 32 (64?) bit VAS, how large? vs size of memory?
 - Used sparsely

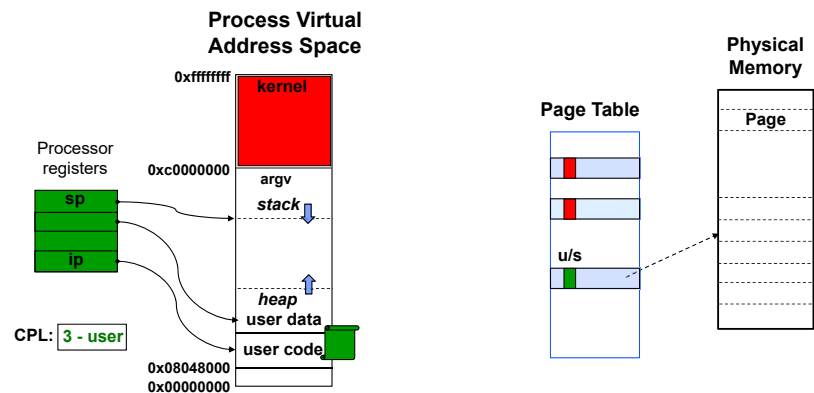
Pintos: page_dir.c

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.100

User → Kernel

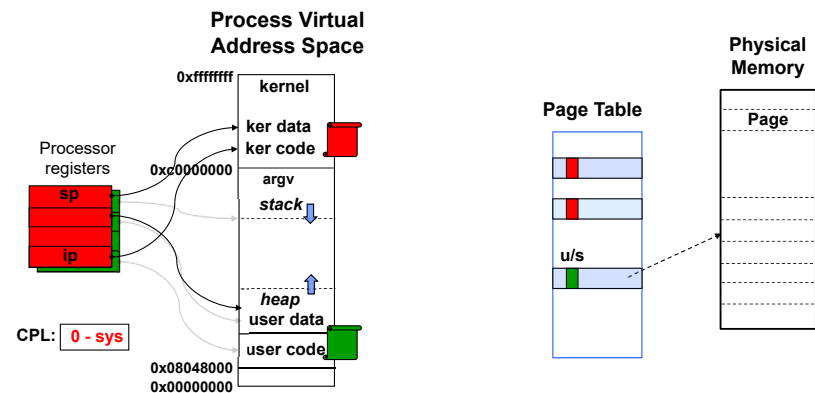


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.101

User → Kernel



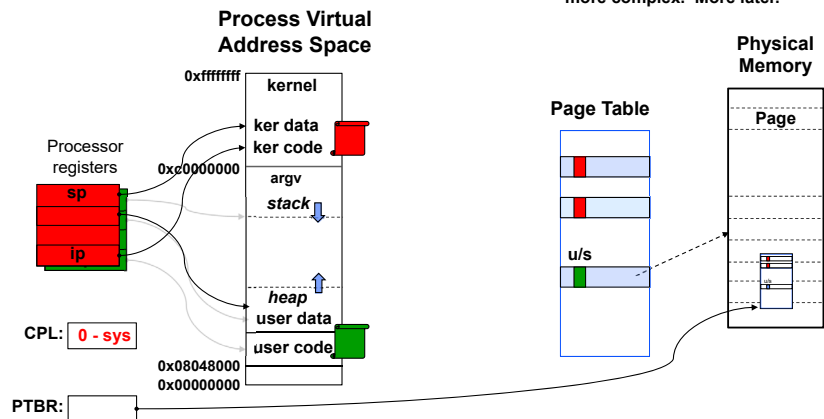
9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.102

Page Table Resides in Memory*

* In the simplest case. Actually more complex. More later.



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.103

Kernel Portion of Address Space

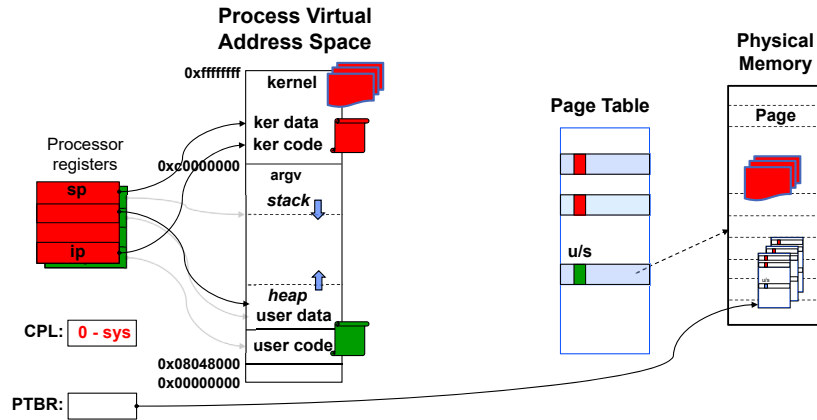
- Kernel memory is mapped into address space of *every* process
- Contains the kernel code
 - Loaded when the machine booted
- Explicitly mapped to physical memory
 - OS creates the page table
- Used to contain all kernel data structures
 - Lists of processes/threads
 - Page tables
 - Open file descriptions, sockets, ttys, ...
- Kernel stack for each thread

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.104

1 Kernel Code, Many Kernel Stacks

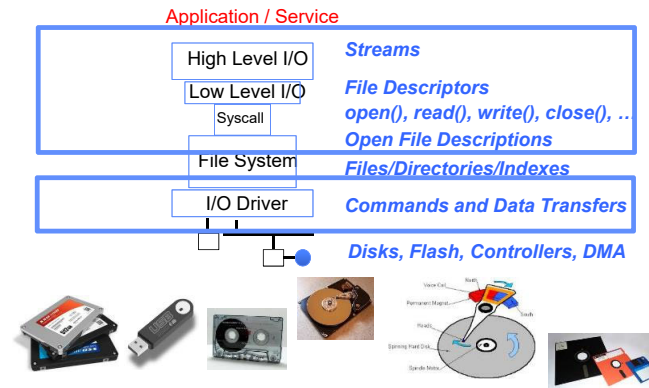


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.105

Recall: I/O and Storage Layers



What we've covered so far...

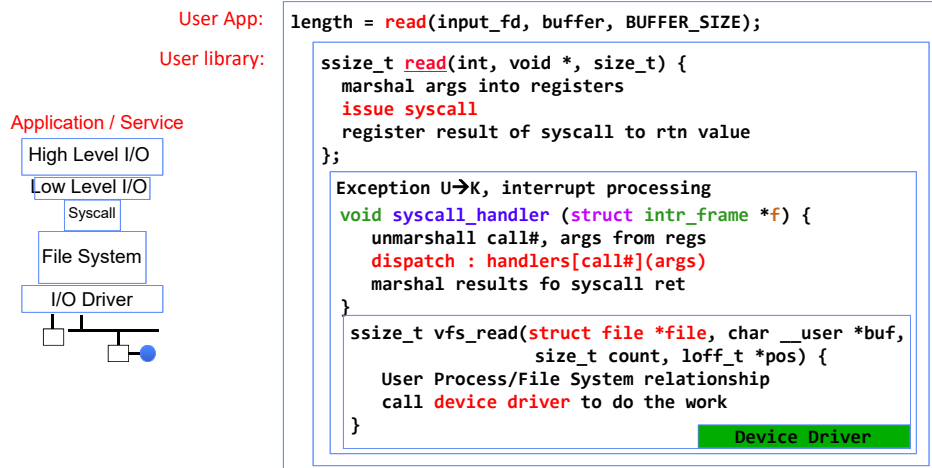
What we'll peek at next

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.106

Layers of I/O...

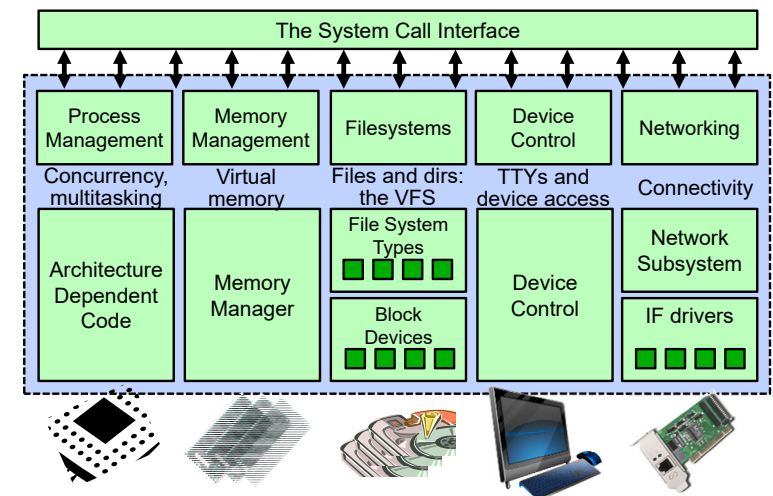


9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.107

Many different types of I/O



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.108

Recall: Internal OS File Description

- Internal Data Structure describing everything about the file

- Where it resides
- Its status
- How to access it

- Pointer: **struct file *file**

- Everything accessed with file descriptor has one of these

- **Struct file_operations *f_op:**
Describes how this particular device implements its operations

- For disks: points to file operations
- For pipes: points to pipe operations
- For sockets: points to socket operations

```
746 struct file {
747     union {
748         struct list_node fu_llist;
749         struct rcu_head fu_rcuhead;
750     } f_u;
751     struct path f_path;
752     #define f_dentry f_path.dentry
753     struct inode f_inode; /* caci
754     const struct file_operations *f_op;
755     /*
756     * Private f_op.links, f_flags.
757     * These must be taken from TQ context.
758     */
759     spinlock_t f_lock;
760     atomic_long_t f_count;
761     unsigned int f_flags;
762     fmode_t f_mode;
763     struct mutex f_pos_lock;
764     loff_t f_pos;
765     struct fown_struct f_owner;
766     const struct cred *f_cred;
767     struct file_ra_state f_ra;
768     u64 f_version;
769     #ifdef CONFIG_SECURITY
770     void *f_security;
771     #endif
772     /* needed for tty driver, and maybe others */
773     void *private_data;
774     #ifdef CONFIG_EPOLL
775     /* Used by fs/eventpoll.c to link all the hooks
776     struct list_head f_ep_links;
777     struct list_head f_tfile_llist;
778     #endif
779     #ifdef CONFIG_EPOLL
780     struct address_space *f_mapping;
781     /* __attribute__((aligned(4))) */ /* test something weird
782 }
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.109

File_operations: Why everything can look like a file

- Associated with particular hardware device or environment (i.e. file system)
- Registers / Unregisters itself with the kernel
- Handler functions for each of the file operations

```
struct file_operations {
    struct module *owner;
    loff_t (*llseek) (struct file *, loff_t, int);
    ssize_t (*read) (struct file *, char __user *, size_t, loff_t *);
    ssize_t (*write) (struct file *, const char __user *, size_t, loff_t *);
    ssize_t (*aio_read) (struct kiocb *, const struct iovec *, unsigned long, loff_t);
    ssize_t (*aio_write) (struct kiocb *, const struct iovec *, unsigned long, loff_t);
    int (*readdir) (struct file *, void *, filldir_t);
    unsigned int (*poll) (struct file *, struct poll_table_struct *);
    int (*ioctl) (struct inode *, struct file *, unsigned int, unsigned long);
    int (*mmap) (struct file *, struct vm_area_struct *);
    int (*open) (struct inode *, struct file *);
    int (*flush) (struct file *, filp_t);
    int (*release) (struct inode *, struct file *);
    int (*fsync) (struct file *, struct dentry *, int datasync);
    int (*fasync) (int, struct file *, int);
    int (*flock) (struct file *, int, struct file_lock *);
    ...
};
```

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.110

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
```

```
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count)))
        return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

- Read up to "count" bytes from "file" starting from "pos" into "buf".
- Return error or number of bytes read.

Linux: fs/read_write.c

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.111

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
```

```
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count)))
        return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

- Make sure we are allowed to read this file

Linux: fs/read_write.c

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.112

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count))) return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

Check if file has read methods

Linux: fs/read_write.c

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count))) return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

• Check whether we can write to buf (e.g., buf is in the user space range)
• unlikely(): hint to branch prediction this condition is unlikely

Linux: fs/read_write.c

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count))) return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

Check whether we read from a valid range in the file.

Linux: fs/read_write.c

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count))) return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

If driver provide a read function (f_op->read) use it; otherwise use do_sync_read()

Linux: fs/read_write.c

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count))) return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

Linux: fs/read_write.c

Notify the parent of this file that the file was read
(see <http://www.fieldses.org/~bfields/kernel/vfs.txt>)

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.117

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count))) return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

Linux: fs/read_write.c

Update the number of bytes read by "current" task (for scheduling purposes)

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.118

File System: From Syscall to Driver

```
ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
{
    ssize_t ret;
    if (!(file->f_mode & FMODE_READ)) return -EBADF;
    if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
        return -EINVAL;
    if (unlikely(!access_ok(VERIFY_WRITE, buf, count))) return -EFAULT;
    ret = rw_verify_area(READ, file, pos, count);
    if (ret >= 0) {
        count = ret;
        if (file->f_op->read)
            ret = file->f_op->read(file, buf, count, pos);
        else
            ret = do_sync_read(file, buf, count, pos);
        if (ret > 0) {
            fsnotify_access(file->f_path.dentry);
            add_rchar(current, ret);
        }
        inc_syscr(current);
    }
    return ret;
}
```

Linux: fs/read_write.c

Update the number of read syscalls by "current" task (for scheduling purposes)

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.119

Device Drivers

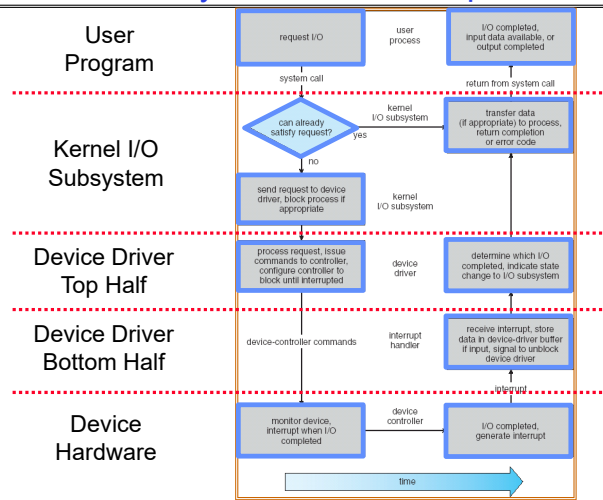
- **Device Driver:** Device-specific code in the kernel that interacts directly with the device hardware
 - Supports a standard, internal interface
 - Same kernel I/O system can interact easily with different device drivers
 - Special device-specific configuration supported with the `ioctl()` system call
- Device Drivers typically divided into two pieces:
 - Top half: accessed in call path from system calls
 - » implements a set of **standard, cross-device calls** like `open()`, `close()`, `read()`, `write()`, `ioctl()`, `strategy()`
 - » This is the kernel's interface to the device driver
 - » Top half will *start* I/O to device, may put thread to sleep until finished
 - Bottom half: run as interrupt routine
 - » Gets input or transfers next block of output
 - » May wake sleeping threads if I/O now complete

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.120

Life Cycle of An I/O Request



9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.121

Conclusion

- **Monitors:** A lock plus one or more condition variables
 - Always acquire lock before accessing shared data
 - Use condition variables to wait inside critical section
 - » Three Operations: **Wait()**, **Signal()**, and **Broadcast()**
- Monitors represent the logic of the program
 - Wait if necessary
 - Signal when change something so any waiting threads can proceed
 - Monitors supported natively in a number of languages
- Readers/Writers Monitor example
 - Shows how monitors allow sophisticated controlled entry to protected code
- Kernel Thread: Stack+State for independent execution in kernel
 - Every user-level thread paired one-to-one with kernel thread
 - Kernel thread associated with user thread is "suspended" (ready to go) when user-level thread is running
- Device Driver: Device-specific code in kernel that interacts directly with device hardware
 - Supports a standard, internal interface
 - Same kernel I/O system can interact easily with different device drivers

9/28/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 9.122