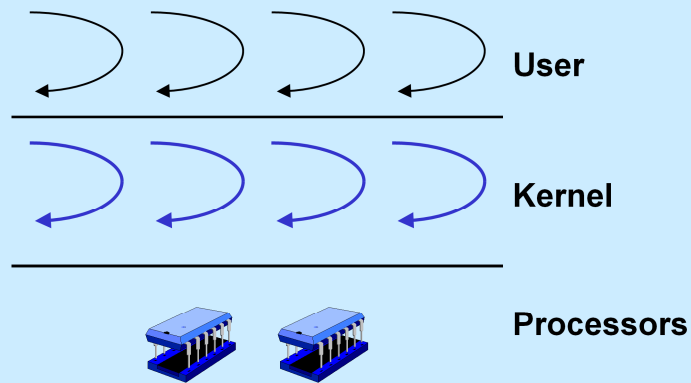


CS 33

Threads and the Operating System

One-Level Threads Model

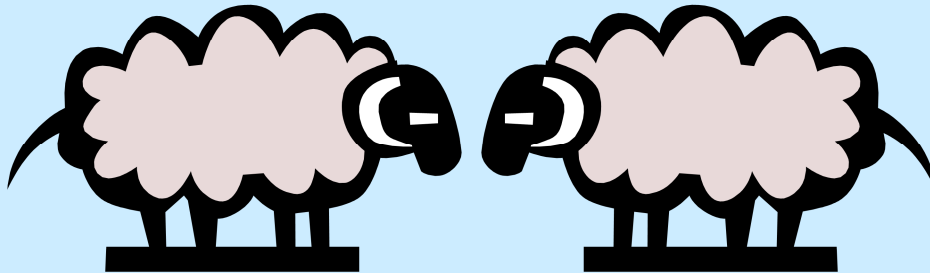


In most systems there are actually two components of the execution context: the user context and the kernel context. The former is for use when a processor is executing user code; the latter is for use when the processor is executing kernel code (on behalf of the chore). How these contexts are manipulated is one of the more crucial aspects of a threads implementation.

The conceptually simplest approach is what is known as the one-level model: each thread consists of both contexts. Thus a thread is scheduled to a processor and the processor can switch back and forth between the two types of contexts. A single scheduler in the kernel can handle all the multiplexing duties. The threading implementation in Windows is (mostly) done this way.

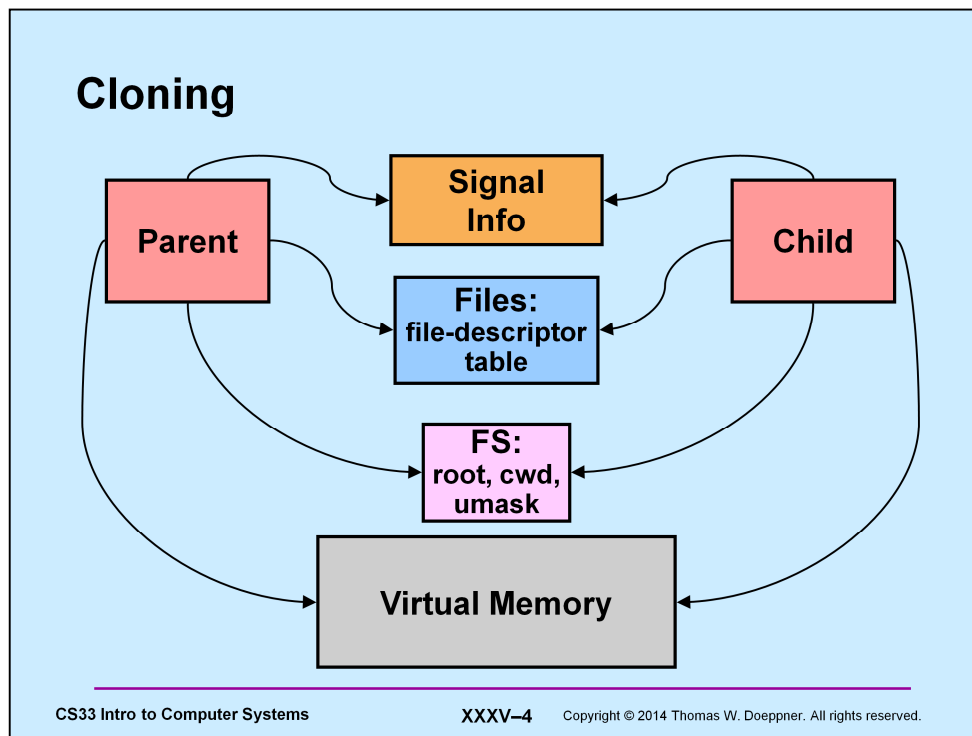
Variable-Weight Processes in Linux

- Variant of one-level model
- Portions of parent process selectively *copied* into or *shared* with child process
- Children created using *clone* system call



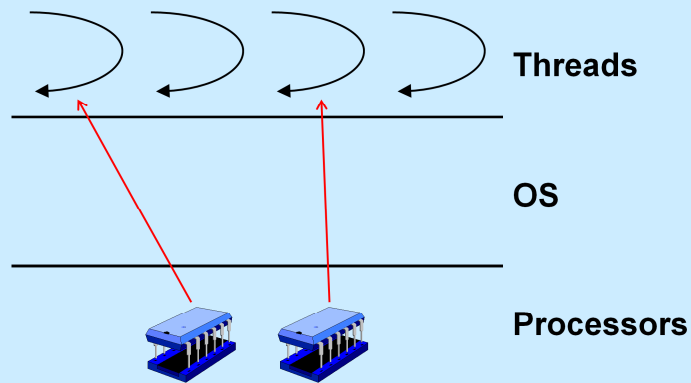
Unlike most other Unix systems, which make a distinction between processes and threads, allowing multithreaded processes, Linux maintains the one-thread-per-process approach. However, so that we can have multiple threads sharing an address space, Linux supports the *clone* system call, a variant of *fork*, via which a new process can be created that shares resources (in particular, its address space) with the parent. The result is a variant of the one-level model.

This approach is not unique to Linux. It was used in SGI's IRIX and was first discussed in early '89, when it was known as variable-weight processes. (See "Variable-Weight Processes with Flexible Shared Resources," by Z. Aral, J. Bloom, T. Doeppner, I. Gertner, A. Langerman, G. Schaffer, *Proceedings of Winter 1989 USENIX Association Meeting*.)



As implemented in Linux, a process may be created with the *clone* system call (in addition to using the *fork* system call). One can specify, for each of the resources shown in the slide, whether a copy is made for the child or the child shares the resource with the parent. Only two cases are generally used: everything is copied (equivalent to *fork*) or everything is shared (creating what we ordinarily call a thread, though the “thread” has a separate process ID).

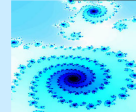
Scheduling



The operating system is responsible for multiplexing the execution of threads on the available processors. The OS's *scheduler* is responsible for assigning threads to processors. Periodically, say every millisecond, each processor is interrupted and calls upon the OS to determine if another thread should run. If so, the current thread on the processor is preempted in favor of the next thread. Assuming all threads are treated equally, over a sufficient period of time each thread gets its fair share of available processor time. Thus, even though a system may have only one processor, all threads make progress and give the appearance of running simultaneously.

Real-Life Example

- Your iPhone is broken
 - mp3 player
- You're watching a phenomenal YouTube video
- You're computing a fractal model of the universe
- A 33 assignment is due
 - editor, compiler, debugger
- You've got to do everything on one computer
- Can your scheduler hack it?



```
void *Timeout_WatchDog(void *arg) {  
    Timeout_t *timeout = (Timeout_t *)arg;  
    pthread_barrier_wait(&timeout->client->started);  
    struct timespec to;  
    to.tv_sec = Timeout_wait_secs;  
    to.tv_nsec = 0;  
    struct timeval now;  
    while(1) {  
        if (nanosleep(&to, 0) == -1) {  
            perror("nanosleep");  
            exit(1);  
        }  
        if (timeout->active == 0) {  
            pthread_cancel(timeout->client->thread);  
            pthread_exit(0);  
        } else  
            timeout->active = 0;  
    }  
}
```

Scheduling

- **Aims**
 - provide timely response
 - provide quick response
 - use resources equitably

The slide lists three possible aims of an operating system's scheduler. Note that “timely response” and “quick response” are two very different things!

Timely Response

- “Hard” real time
 - chores *must* be completed on time
 - » controlling a nuclear power plant
 - » landing (softly) on a comet

Providing timely response is the realm of what's known as *hard real-time systems*, where there are various chores that simply must be completed on time (if not, disaster occurs). What's important in such systems is predictable behavior.

Fast Response

- “Soft” real time
 - the longer it takes, the less useful a chore’s result becomes
 - » responding to user input
 - » playing streaming audio or video

Fast response is what’s required of *soft real-time systems*, where timely response is important, but disaster doesn’t occur if time requirements are not met. However, it is important that the requirements are met “most of the time.”

Sharing

- All active threads share processor time equally

Linux Scheduler Evolution

- **Original scheduler**
 - very simple
 - poor scaling
- **O(1) scheduler**
 - introduced in 2.5
 - less simple
 - better scaling
- **Completely fair scheduler (CFS)**
 - even better
 - simpler in concept
 - much less so in implementation

The original scheduler was introduced in the early 1990s. It's very simple and works reasonably well on lightly loaded uniprocessors. The O(1) scheduler, introduced in release 2.5, is considerably more sophisticated. CFS is based on stride scheduling.

Original Scheduler

- **Three per-process scheduling variables**
 - *policy*: which one
 - » we look at only the default policy
 - *priority*: time-slice parameter (“nice” value)
 - *counter*: records processor consumption

We start by discussing the original scheduler. Three variables are used for scheduling processes, as shown in the slide. The first two are inherited from a process's parent. Using (privileged) system calls, one can change the first two. The second, *priority*, can be worsened with non-privileged system calls, but can be improved only with a privileged system call.

Original Scheduler: Time Slicing



- Clock “ticks” HZ times per second
 - interrupt/tick
- Per-process *counter*
 - current process’s is decremented by one each tick
 - time slice over when counter reaches 0

Implementing time slicing requires the support of the clock-interrupt handler. Each process’s counter variable is initialized to the process’s *priority*. At each clock “tick”, *counter* is decremented by one. When it reaches zero, the process’s time slice is over and it must relinquish the processor. (How *counter* gets a positive value again is discussed in the next slide.) In current versions of Linux, “HZ” is 1000.

Original Scheduler: Throughput

- Scheduling cycle
 - length, in “ticks,” is sum of priorities
 - each process gets *priority* ticks/cycle
 - » *counter* set to *priority*
 - » cycle over when *counters* for runnable processes are all 0
 - sleeping processes get “boost” at wakeup
 - » at beginning of each cycle, for each process:
$$\text{counter} = \text{counter}/2 + \text{priority}$$

The default scheduling policy is a throughput policy, which means that all processes are guaranteed to make progress, though some (those with better priority) make faster progress (get a higher percentage of processor cycles) than others. It's easiest to understand if we assume a fixed number of processes, all of which use the default policy and all of which remain runnable. Scheduling is based on cycles, the length of which (measured in ticks) is the sum of the (runnable) processes' priorities. In each cycle, each process thus gets *priority* ticks of processor time, for that process's value of *priority*. This is done by setting each process's *counter* to *priority* when the cycle starts.

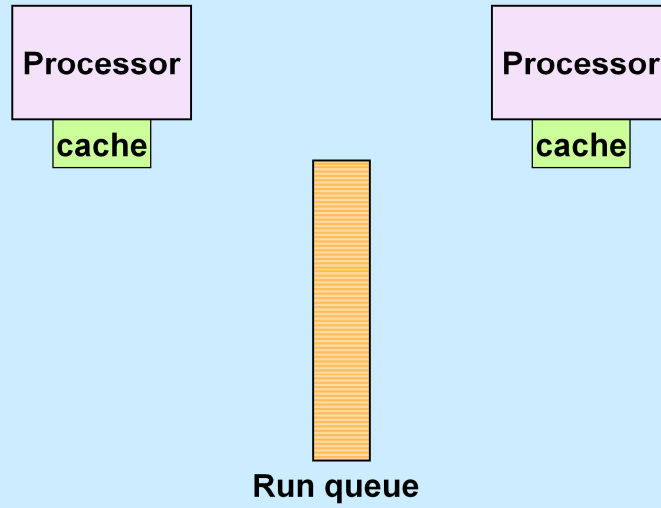
Sleeping processes are given a “boost” when they wake up. The rationale is that we want to favor interactive and I/O intensive requests (the latter don't use much processor time and thus let's quickly get them back to waiting for an I/O operation to complete). To implement this, at the end of each cycle, all processes (not just all runnable processes) have their *counters* set as shown in the slide. (Thus the maximum value of *counter* is twice the process's priority.)

Original Scheduler: Who's Next?

- Run-queue searched beginning to end
 - new arrivals go to front
- Next running process is first process with highest “goodness”
 - *counter* for default processes

When a process gives up the processor, it executes (in the kernel) a routine called *schedule* which determines the next process to run, as shown in the slide. Note that all runnable processes are examined so as to find the “best” one. This is not a good strategy if there are a lot of them; on a PC there will be few, but on a busy server there could be many.

Diagram



Original Scheduler: Problems

- **$O(n)$ execution**
- **Poor interactive performance with heavy loads**
- **Contention for run-queue lock**
- **Processor affinity**
 - cache “footprint”

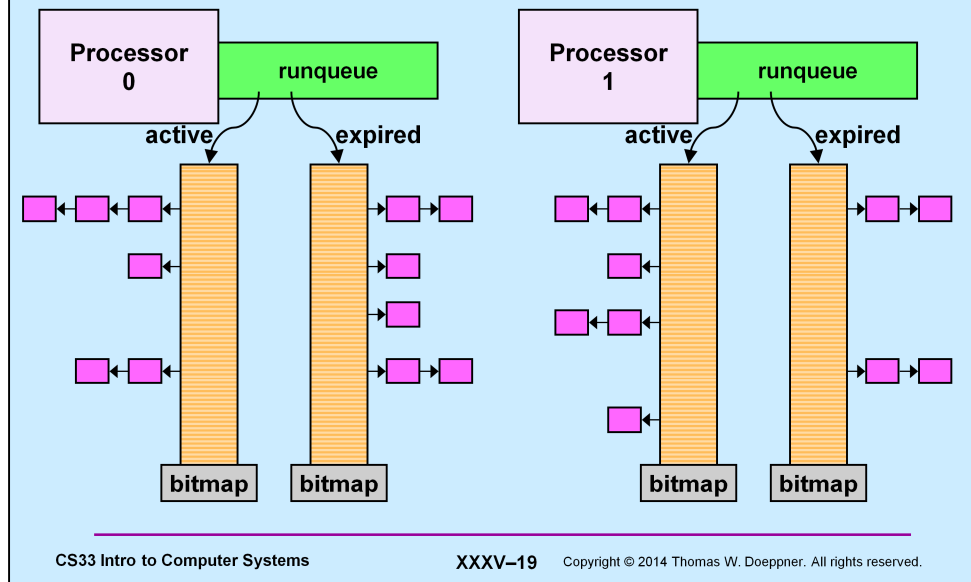
The main problems with the original scheduler are summarized in the slide. As we've seen, the scheduler must periodically examine all processes, as well as examining all runnable processes for each scheduling decision. When a system is running with a heavy load, interactive processes will get a much smaller percentage of processor cycles than they do under a lighter load. Since there's one run queue feeding all processors, all must contend for its lock. Finally, though some attempt is made at dealing with processor-affinity issues, processes still tend to move around among available processors (and thus losing any advantage of their “cache footprint”).

O(1) Scheduler

- **All concerns of original scheduler plus:**
 - **efficient, scalable execution**
 - **identify and favor interactive processes**
 - **good MP performance**
 - » **minimal lock overhead**
 - » **processor affinity**

The O(1) scheduler deals with all the concerns dealt with by the original one along with some additional concerns, as shown in the slide.

O(1) Scheduler: Data Structures



The O(1) scheduler associates a pair of queues with each processor via a per-processor *runqueue*. Each queue is actually an array of lists of processes, one list for each possible priority value. There are 140 priorities, running from 0 to 139; “good” priorities have low numbers; “bad” priorities have high numbers. Real-time priorities run from 0 to 99; normal priorities run from 100 to 139. Associated with each queue is a 140-element bit map indicating which priority values have a non-empty list of processes.

O(1) Scheduler: Queues

- **Two queues per processor**
 - **active: processes with remaining time slice**
 - **expired: processes with no more time slice**
 - **each queue is an array of lists of processes of the same priority**
 - » **bitmap indicates which priorities have processes**
 - **processors scheduled from private queues**
 - » **infrequent lock contention**
 - » **good affinity**

When processes become runnable they are assigned time slices (based on their priority) and put on some processor's *active* queue. When a processor needs a process to run, it chooses the highest priority process on its active queue (this can be done quickly (and in time bounded by a constant) by scanning the queue's bit vector to find the first non-empty priority level, then selecting the first process from the list at that priority). When a process completes its time slice, it goes back to either the active or the expired queue of its processor (as explained shortly). If it blocks for some reason and later wakes up, it will generally go back to the active queue of the processor it last was on. Thus processes tend to stay on the same processor (providing good use of the cache footprint). Since processors rarely access other processors' queues, there is very little lock contention.

O(1) Scheduler: Actions

- **Process switch**
 - pick best priority from active queue
 - » if empty, switch active and expired
 - new process's time slice is function of its priority
- **Wake up**
 - priority is boosted or dropped depending on sleep time
 - higher priority processes get longer time quanta
- **Time-slice expiration**
 - interactive processes (those waking up) rejoin active queue

When a process completes its time slice, it is inserted into either its processor's active queue or its processor's expired queue, depending on its priority. The intent is that interactive and real-time processes get another time slice on the processor, while other processes have to wait a bit on the inactive queue. When there are no processes remaining in the active queue, the two queues are switched. Of course, if there are interactive processes, the active queue might never empty out. So, if the processes in the expired queue have been waiting too long (how long this is depends on the number of runnable processes on the queue), interactive processes completing their time slices go to the expired queue rather than the active queue. Runnable real-time processes never go to the expired queue: they are always in the active queue (and always have priority over non-real-time processes). Thus two queues are employed as a means to guarantee that, in the absence of real-time processes, all processes get some processor time. Lower priority processes will remain on the active queue until all higher-priority processes have moved to the expired queue.

The net effect is similar to the original scheduler: in the absence of real-time processes, processes get the processor in proportion to their priority. However, interactive processes (those that have recently woken up) get extra time slices.

O(1) Scheduler: Load Balancing

- **Processors with empty queues steal from busiest processor**
 - checked every millisecond
- **Processors with relatively small queues also steal from busiest processor**
 - checked every 250 milliseconds

Since processors schedule strictly from their own private queues, load balancing is an issue (it wasn't with the old scheduler, since there was only one global queue serving all processors). Each processor checks its queues for emptiness every millisecond. If empty, it calls a load balancing routine to find the processor with the largest queues and then transfers processes from that processor to the idle one until they are no longer imbalanced (they are considered balanced if they are no more than 25% different in size). Similarly, each processor checks the other processors' queues every 250 milliseconds. If an imbalance is found (and it's not just a momentary imbalance but has been that way since the last time the processor's queues were examined), then load balancing is done.

Shared Servers

- **You and four friends each contribute \$1000 towards a server**
 - you, rightfully, feel you own 20% of it
- **Your friends are into threads, you're not**
 - they run 5-threaded programs
 - you run a 1-threaded program
- **Their programs each get 5/21 of the processor**
- **Your programs get 1/21 of the processor**
 - you should have paid more attention in class

Lottery Scheduling

- **25 lottery tickets are distributed equally to you and your four friends**
 - you give 5 tickets to your one thread
 - they give one ticket each to their threads
- **A lottery is held for every scheduling decision**
 - your thread is 5 times more likely to win than the others

Proportional-Share Scheduling

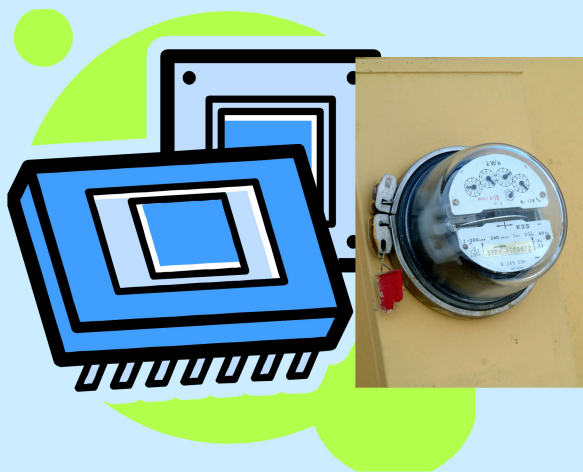
- **Lottery scheduling**
- **Stride scheduling**
 - 1995 paper by Waldspurger and Weihl
- **Completely fair scheduling (CFS)**
 - added to Linux in 2007

The Stride scheduling paper is Stride Scheduling: Deterministic Proportional-Share Resource Management, Technical Memorandum MIT/LCS/TM-528, MIT Laboratory for Computer Science, June 1995.

There really isn't any documentation for CFS. It's making its first appearance in Linux 2.6.23. Some information may be found at <http://kerneltrap.org/node/8059>.

CFS and stride scheduling are pretty much the same thing.

Metered Processors



To measure the usage of a processor, let's assume the existence of a meter.

Algorithm

- Each thread has a meter, which runs only when the thread is running on the processor
- At every clock tick
 - give processor to thread that's had the least processor time as shown on its meter
 - in case of tie, thread with lowest ID wins

Assuming all threads are equal, all started at the same time, and all run forever, the intent is to share the processor equitably. Note that as the time between clock ticks approaches zero, each thread gets $1/n$ of total processor time, where n is the number of threads.

Issues

- Some threads may be more important than others
- What if new threads enter system?
- What if threads block for I/O or synchronization?

We'll talk about the first problem in the next few slides. The other two problems are taken up in CS 167.

Metered Processors (Mafia Variation)



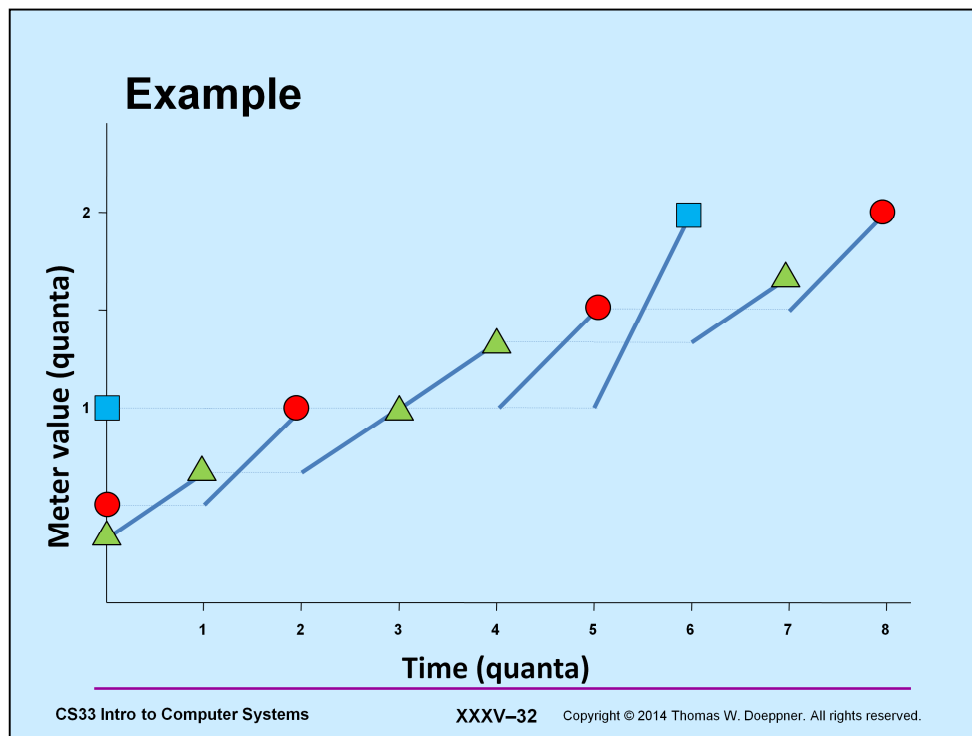
Let's now assume that meters can be "fixed" so that they run more slowly than they should. Thus a thread with a fixed meter gets charged for less processor time than it has actually used.

Details ...

- **Each thread pays a bribe**
 - the greater the bribe, the slower its meter runs
 - to simplify bribing, you buy “tickets”
 - » one ticket is required to get a fair meter
 - » two tickets get a meter running at half speed
 - » three tickets get a meter running at 1/3 speed
 - » etc.

New Algorithm

- Each thread has a (*possibly crooked*) meter, which runs only when the thread is running on the processor
- At every clock tick
 - give processor to thread that's had the least processor time as shown on its meter
 - in case of tie, thread with lowest ID wins



The slide illustrates the execution of three threads using stride scheduling. Thread 1 (labeled with a triangle) has paid a bribe of three tickets. Thread 2 (labeled with a circle) has paid a bribe of two tickets, and thread three (labeled with a square) had paid only one ticket. The thicker lines indicate when a thread is running. Their slopes are proportional to the meter rates (and inversely proportional to the bribe).

You'll Soon Finish CS 33 ...

- You might
 - celebrate
 - take another systems course
 - » 32
 - » 138
 - » 167



- become a 33 TA



Systems Courses Next Semester

- **CS 32**
 - you've mastered low-level systems programming
 - now do things at a higher level
 - learn software-engineering techniques using Java, XML, etc.
- **CS 138**
 - you now know how things work on one computer
 - what if you've got lots of computers?
 - some may have crashed, others may have been taken over by your worst (and brightest) enemy
- **CS 167/169**
 - still mystified about what the OS does?
 - write your own!

The End

Well, not quite ...

Database is due on 12/16.

**We'll do our best to get everything else back
next week.**

Happy coding and happy holidays!