

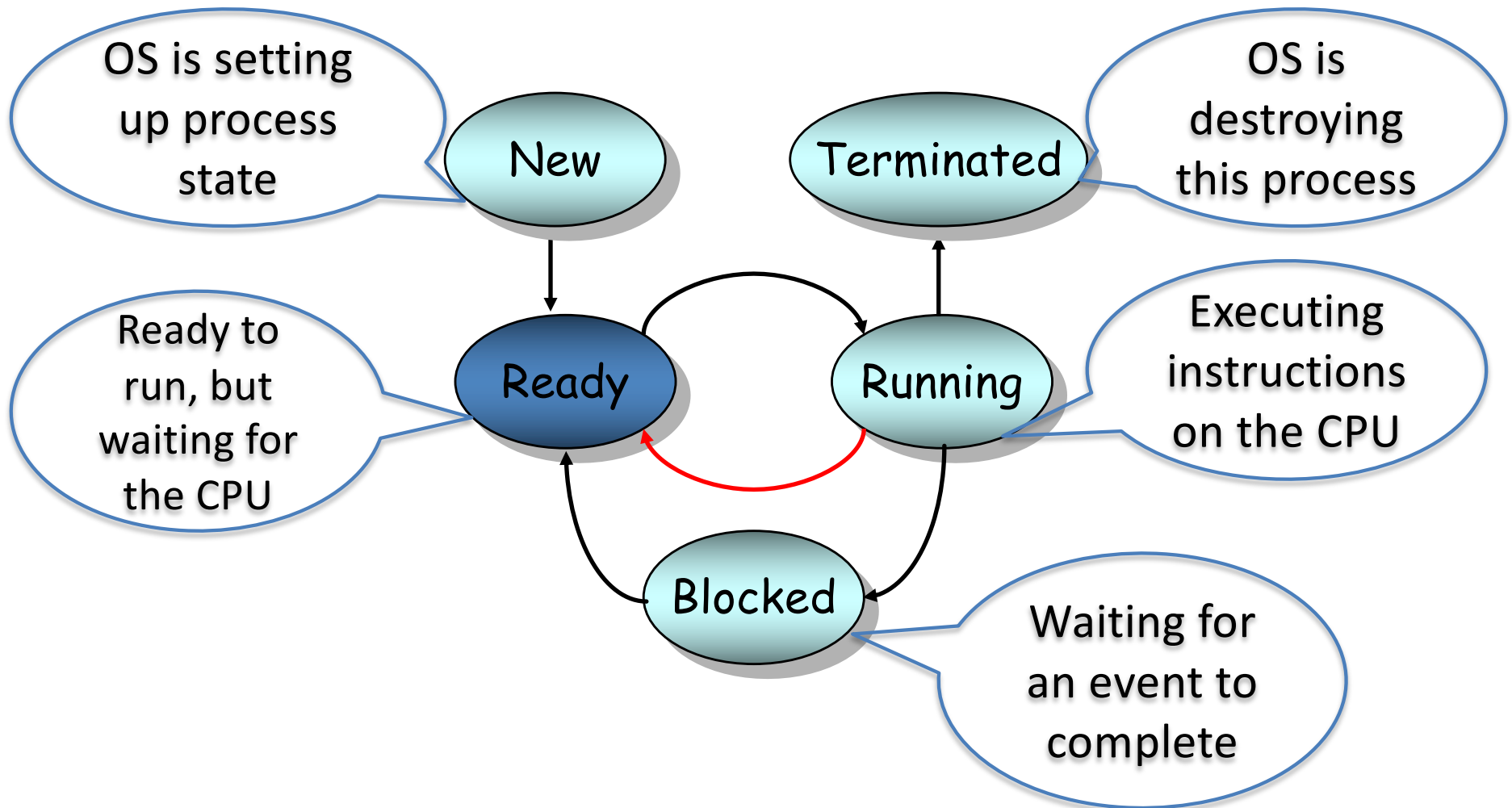
CPU Scheduling

Bringing It Together

- A process is a unit of execution
 - Defines an address space
 - An abstraction for protection
 - Basic unit of execution in the OS
 - Processes are represented as Process Control Blocks in the OS
 - At any time, a process is either New, Ready, Blocked, Running, or Terminated
 - Processes are created and managed through system calls
 - system calls exist for other things, too
 - Supports the OS in its roles as referee and illusionist
- We're in a multiprogramming environment
- We also want to support timesharing and interactivity

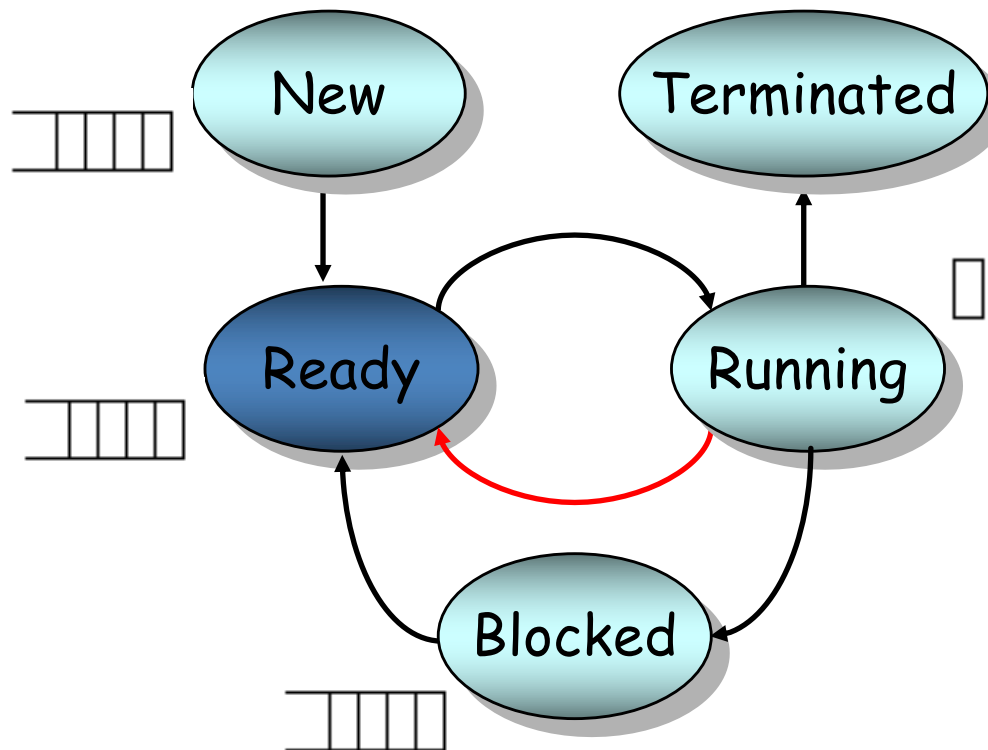
Recall: Process Life Cycle

Processes are always either *running*, *ready to run* or *blocked waiting for an event* to occur



Process Life Cycle

All of the processes the OS is currently managing reside in exactly one of these state queues.



So we know...

The OS uses *multiprogramming*, or *concurrency*, (one process on the CPU running and one or more doing I/O) to increase system utilization and throughput by overlapping I/O and CPU activities.

So now we need to make some decisions about how to manage this multiprogramming stuff.

Questions to Consider

- Long-Term Scheduling
 - How does the OS determine the degree of multiprogramming, i.e., the number of jobs residing at once in the primary memory?
- Short-Term Scheduling:
 - How does (or should) the OS select a process from the ready queue to execute?

Short-Term Scheduling

The *scheduler* schedules the next process on the CPU

- it is the *mechanism*
- it implements a *policy*

The scheduler (maybe) executes when:

- a process switches from running to blocking
- a process is created or terminated
- an interrupt occurs
 - I/O device finishes
 - Periodic timer

Scheduling Policies:

Two Big Categories

Scheduling algorithms are divided into categories based on their use of the timer interrupt.

- Non-preemptive
 - An executing process is not ever pre-empted
 - Scheduler runs when process blocks or terminates---not on hardware interrupts
- Preemptive
 - An executing process might be pre-empted from the CPU
 - Scheduler runs on interrupts, mostly timer, but also system calls and other hardware device interrupts
 - Most OSes today have this

Evaluating Scheduling Policies

Criteria

- **CPU Utilization:** percentage of time that the CPU is busy
- **Throughput:** number of processes completing in a unit of time
- **Turnaround time:** length of time to run a process from initialization to termination, including all the waiting time
- **Response time:** time between issuing a command and getting a result
- **Waiting time:** total amount of time that a process is in the ready queue.

Scheduling Policies

- Ideal CPU scheduler
 - Maximizes CPU utilization and throughput
 - Minimizes turnaround time, waiting time, and response time
 - Requires knowledge of process behavior we typically don't have
- Real CPU schedulers implement particular policy
 - Minimize response time
 - Minimize variance of average response time
 - predictability may be more important than a low average with a high variance.
 - Maximize throughput
 - minimize overhead (OS overhead, context switching)
 - efficient use of system resources (CPU, I/O devices)
 - Minimize waiting time
 - be *fair* by ensuring each process waits similar amounts of time
 - often increases average response time

Scheduling Policies

Simplifying Assumptions

- One process per user
- One thread per process
 - More on this topic next time
- Processes are independent

Researchers developed these algorithms in the 70s when these assumptions were more realistic, and it is still an open problem how to efficiently relax these assumptions.

First-Come-First-Served (FCFS)

- Also known as First-In-First-Out (FIFO)
- Scheduler executes jobs in arrival order
- Jobs run until they either complete or block on I/O
- In early FCFS schedulers, the job did not relinquish the CPU even when it was doing I/O.

Round Robin

- Add a timer and use a pre-emptive policy
- Run each process for its time slice (*scheduling quantum*)
- After each time slice, move the running process to the back of the queue
- Selecting a time slice:
 - Too large: waiting time suffers, degenerates to FCFS if processes are never preempted
 - Too small: throughput suffers because too much time is spent context switching
 - Balance the two by selecting a time slice where context switching is roughly 1% of the time slice.

Most time sharing systems use some variation of this policy.

A typical time slice today is between 10-100 milliseconds, with a context switch time of 0.1 to 1 millisecond.

Scheduling Policies:

Context Switches

- Changing from one process to another is known as a *context switch*
- Steps:
 - Process is running
 - Process blocks or is interrupted
 - Mode switch to kernel
 - It's time for the scheduler to run!
 - Scheduler saves executing process state (to its PCB)
 - Scheduler chooses new process to run
 - Scheduler loads its state (from its PCB)
 - Process is running
- Main source of overhead for a scheduling policy

SJF : Shortest Job First

Shortest Job First (SJF):

- Schedule the job that has the least amount of work (measured in CPU time) to do until its next I/O request or termination.
 - I/O bound jobs get priority over CPU bound jobs.
- Works for preemptive and non-preemptive schedulers.

Using the Past to Predict the Future

- If a process is I/O bound in the past, it is also likely to be I/O bound in the future (programs turn out not to be random).
- To exploit this behavior, the scheduler can favor jobs (schedule them sooner) when they use very little CPU time (absolutely or relatively), thus approximating SJF.
- This policy is **adaptive** because it relies on past behavior and changes in behavior result in changes to scheduling decisions.

Approximating SJF: Multilevel Feedback Queues

- Multiple queues with different priorities.
- OS uses Round Robin scheduling at each priority level, running the jobs in the highest priority queue first.
- Once those finish, OS runs jobs out of the next highest priority queue, etc.
 - Can lead to starvation if highest priority jobs run forever
- Round robin time slice increases exponentially at lower priorities.

Multilevel Feedback Queues: Example

	Priority	Time Slice				
<table><tr><td></td><td>G</td><td>F</td><td>A</td></tr></table>		G	F	A	1	1
	G	F	A			
<table><tr><td></td><td></td><td>E</td></tr></table>			E	2	2	
		E				
<table><tr><td></td><td>D</td><td>B</td></tr></table>		D	B	3	4	
	D	B				
<table><tr><td></td><td>C</td></tr></table>		C	4	8		
	C					

Approximating SJF:

Multilevel Feedback Queues

Adjust priorities as follows (details can vary):

1. Job starts in the highest priority queue
2. If job's time slice expires, drop its priority one level, but no further than the lowest level
3. If job's time slice does not expire (the context switch comes from an I/O request instead), then increase its priority one level, up to the top priority level.

==> In practice, CPU bounds drop like a rock in priority and I/O bound jobs stay at high priority

Improving Fairness

Since SJF is optimal, but unfair, any increase in fairness by giving long jobs a fraction of the CPU when shorter jobs are available will degrade average waiting time. Possible solutions:

- Give each queue a fraction of the CPU time. This solution is only fair if there is an even distribution of jobs among queues.
- Adjust the priority of jobs as they do not get serviced (Unix originally did this.) This ad hoc solution avoids starvation but average waiting time suffers when the system is overloaded because all the jobs end up with a high priority.

Boot Sequence

Booting an OS kernel

- CPU loads boot program from ROM (e.g. UEFI (BIOS) in PCs)
- Boot program:
 - Examines/checks machine configuration (number of CPUs, how much memory, number & type of hardware devices, etc.)
 - Builds a configuration structure describing the hardware
 - Loads the OS kernel and gives it the configuration structure

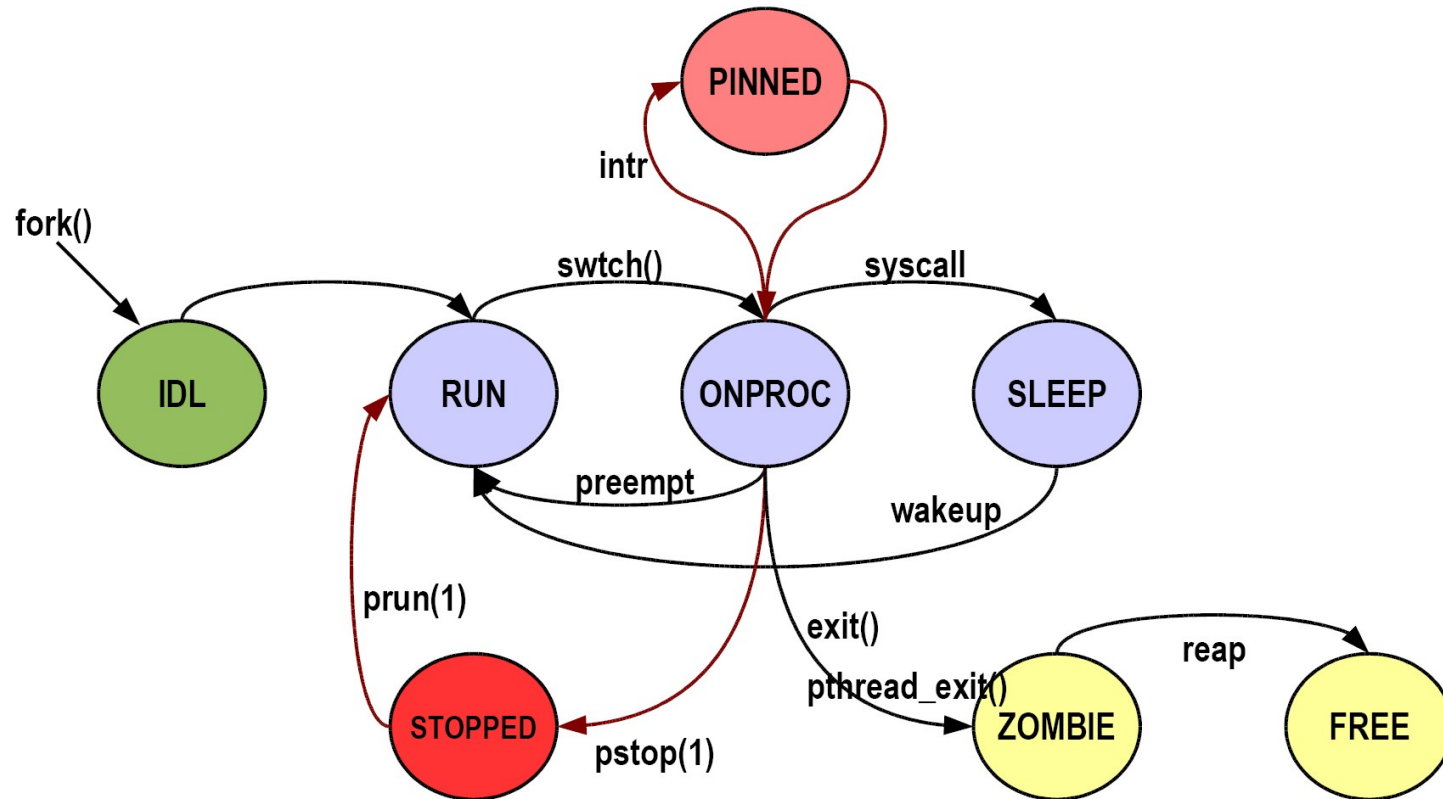
Booting an OS Kernel (Cont'd)

- Operating system initialization:
 - Initialize kernel data structures
 - Initialize the state of all hardware devices
 - Creates a number of processes to start operation (e.g. getty in UNIX, the Windowing system in Windows, etc.)
- After basic processes have started, the OS runs user programs, if available, otherwise enters the idle loop
 - In the idle loop:
 - OS performs some system management & profiling
 - OS halts the processor and enter in low-power mode (notebooks)
 - OS wakes up on interrupts from hardware devices

A Dose of Reality: Scheduling in Solaris



State Transitions



Close to our scheduling diagram, but more complicated