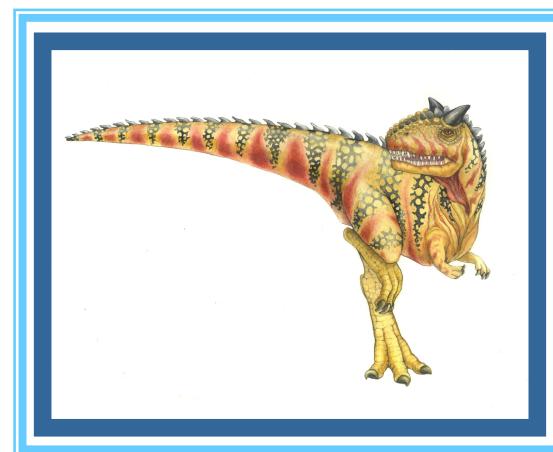
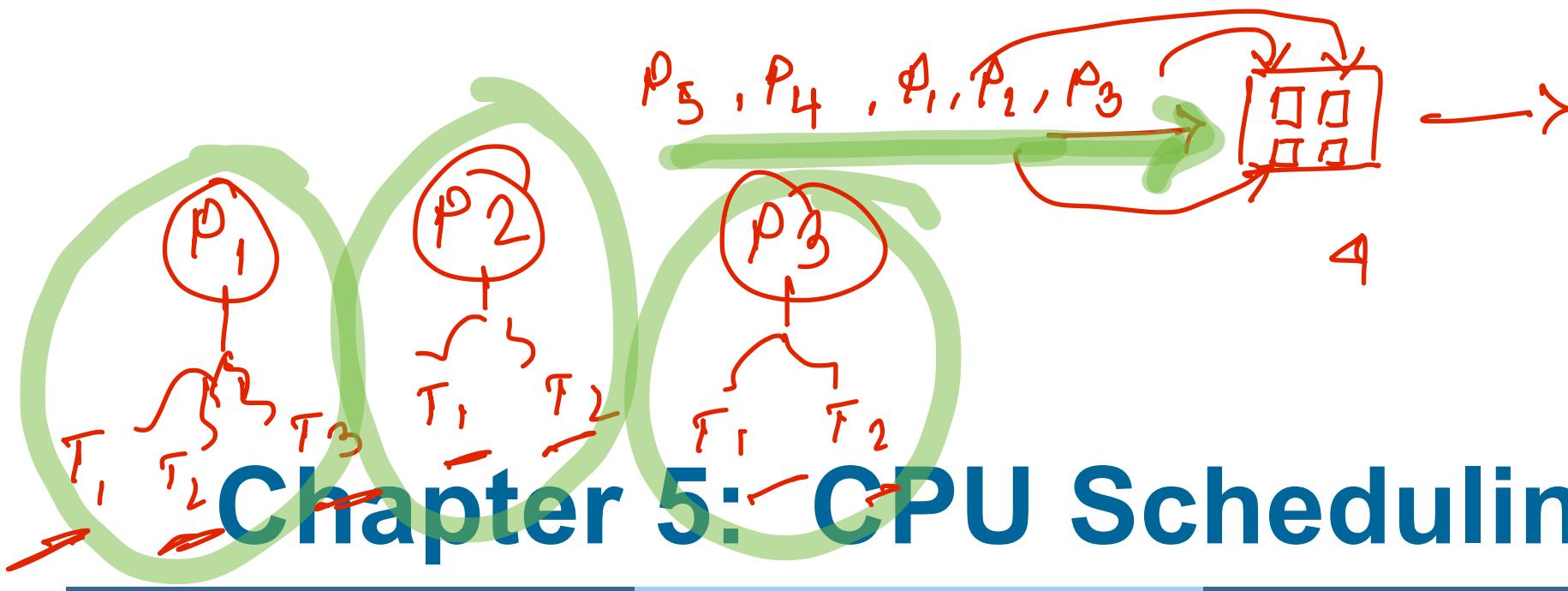


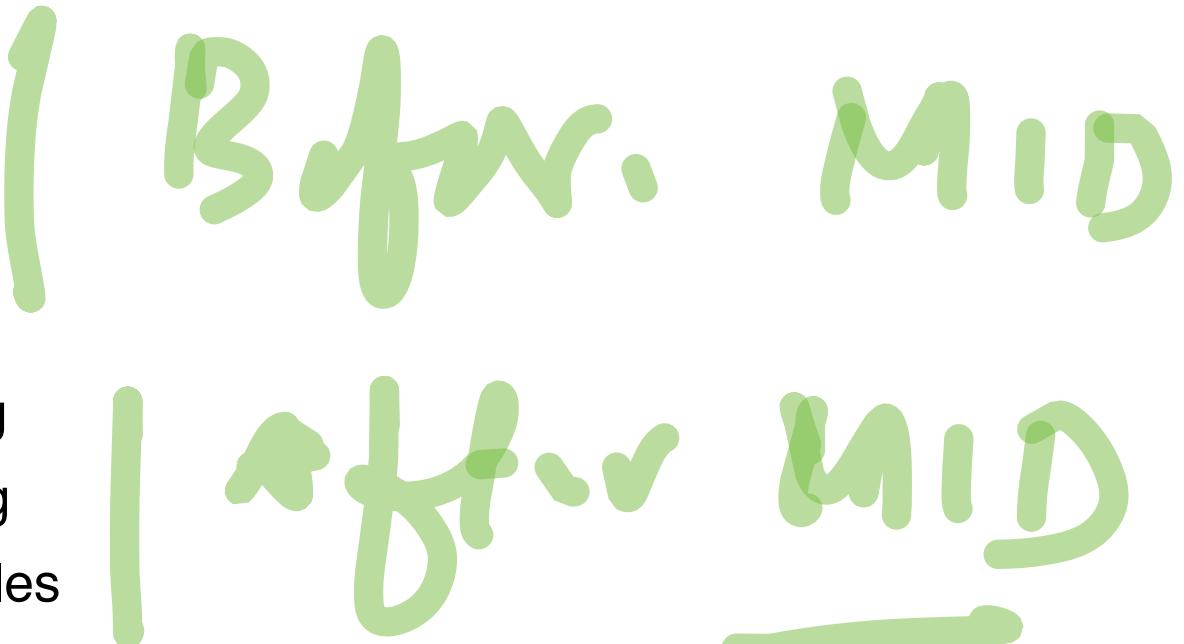
Chapter 5: CPU Scheduling





Chapter 5: CPU Scheduling

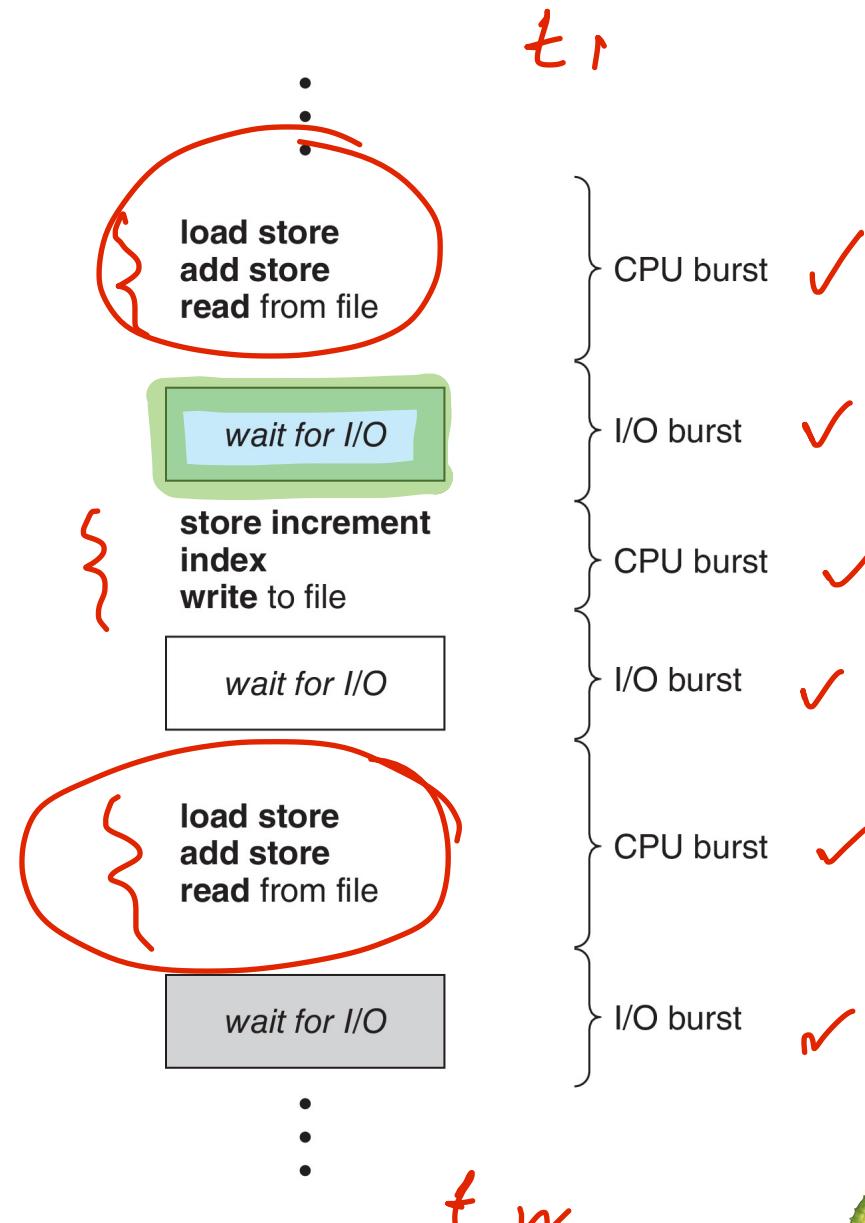
- Basic Concepts ✓
- Scheduling Criteria ✓
- Scheduling Algorithms ✓
- Thread Scheduling ✓
- Multi-Processor Scheduling
- Real-Time CPU Scheduling
- Operating Systems Examples
- Algorithm Evaluation





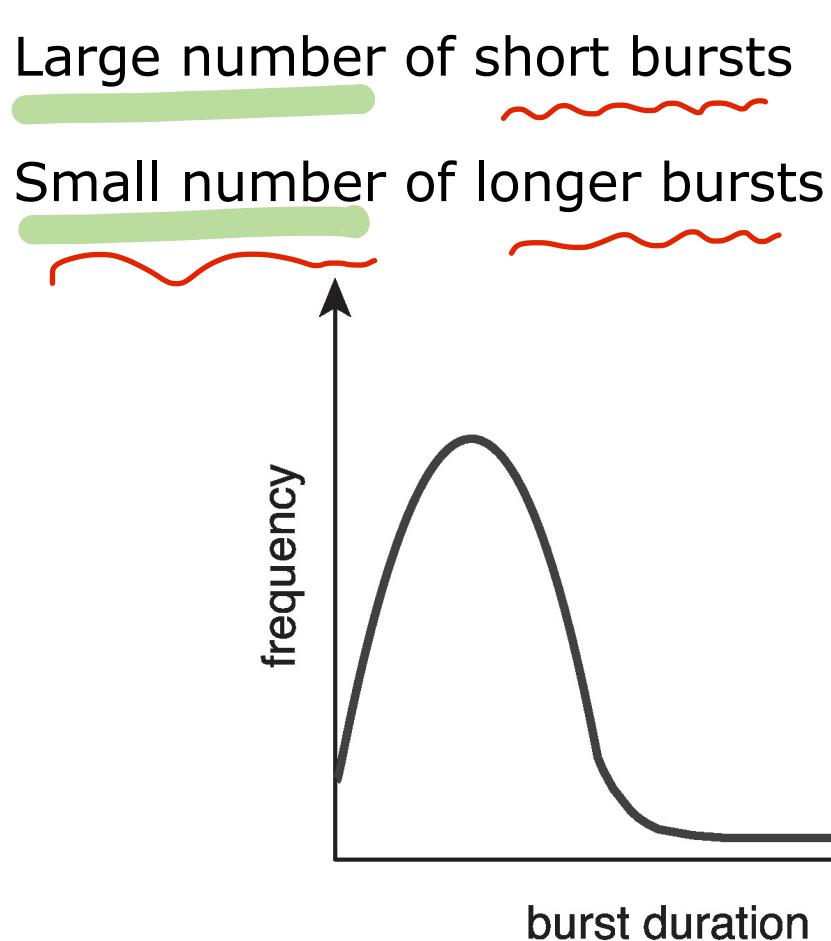
Basic Concepts

- Maximum CPU utilization obtained with multiprogramming
- CPU–I/O Burst Cycle – Process execution consists of a **cycle** of CPU execution and I/O wait
- **CPU burst** followed by **I/O burst**
- CPU burst distribution is of main concern

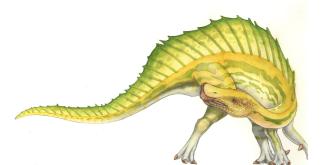




Histogram of CPU-burst Times



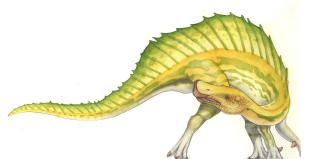
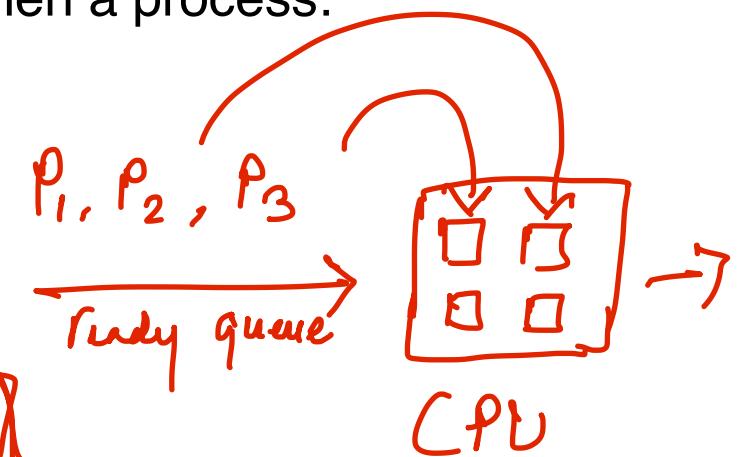
1 ms
10 ms } for example





CPU Scheduler

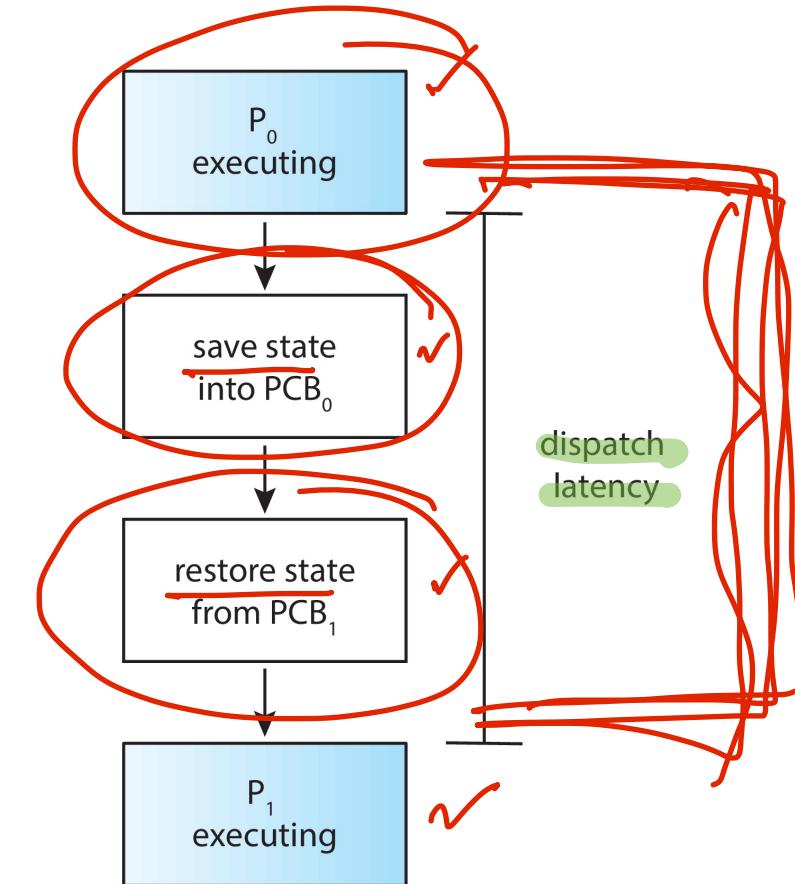
- The **CPU scheduler** selects from among the processes in ready queue, and allocates the a CPU core to one of them
 - Queue may be ordered in various ways
- CPU scheduling decisions may take place when a process:
 1. Switches from running to waiting state
 2. Switches from running to ready state
 3. Switches from waiting to ready
 4. Terminates
- Scheduling under 1 and 4 is **nonpreemptive**
- All other scheduling is **preemptive**
 - Consider access to shared data ✓
 - Consider preemption while in kernel mode ✓
 - Consider interrupts occurring during crucial OS activities ✓



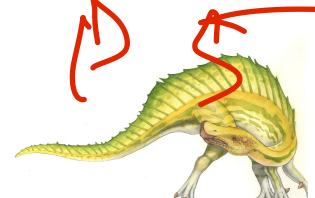


Dispatcher

- Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:
 - switching context ✓
 - switching to user mode ✓
 - jumping to the proper location in the user program to restart that program ✓
- **Dispatch latency** – time it takes for the dispatcher to stop one process and start another running



$P_0, P_1, P_2, P_3, P_4, P_5$





Scheduling Criteria

- ✓ ■ **CPU utilization** – keep the CPU as busy as possible
- ✓ ■ **Throughput** – # of processes that complete their execution per time unit
- ✓ ■ **Turnaround time** – amount of time to execute a particular process
- ✓ ■ **Waiting time** – amount of time a process has been waiting in the ready queue
- ✓ ■ **Response time** – amount of time it takes from when a request was submitted until the first response is produced, not output (for time-sharing environment)





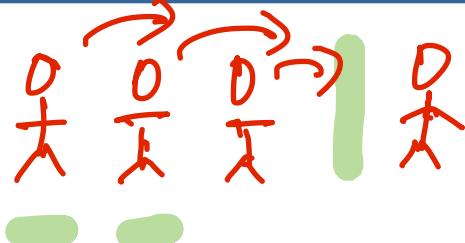
Scheduling Algorithm Optimization Criteria

- Max CPU utilization ✓
- Max throughput ✓
- Min turnaround time ✓
- Min waiting time ✓
- Min response time ✓





First-Come, First-Served (FCFS) Scheduling



Process

P_1

P_2

P_3

Burst Time

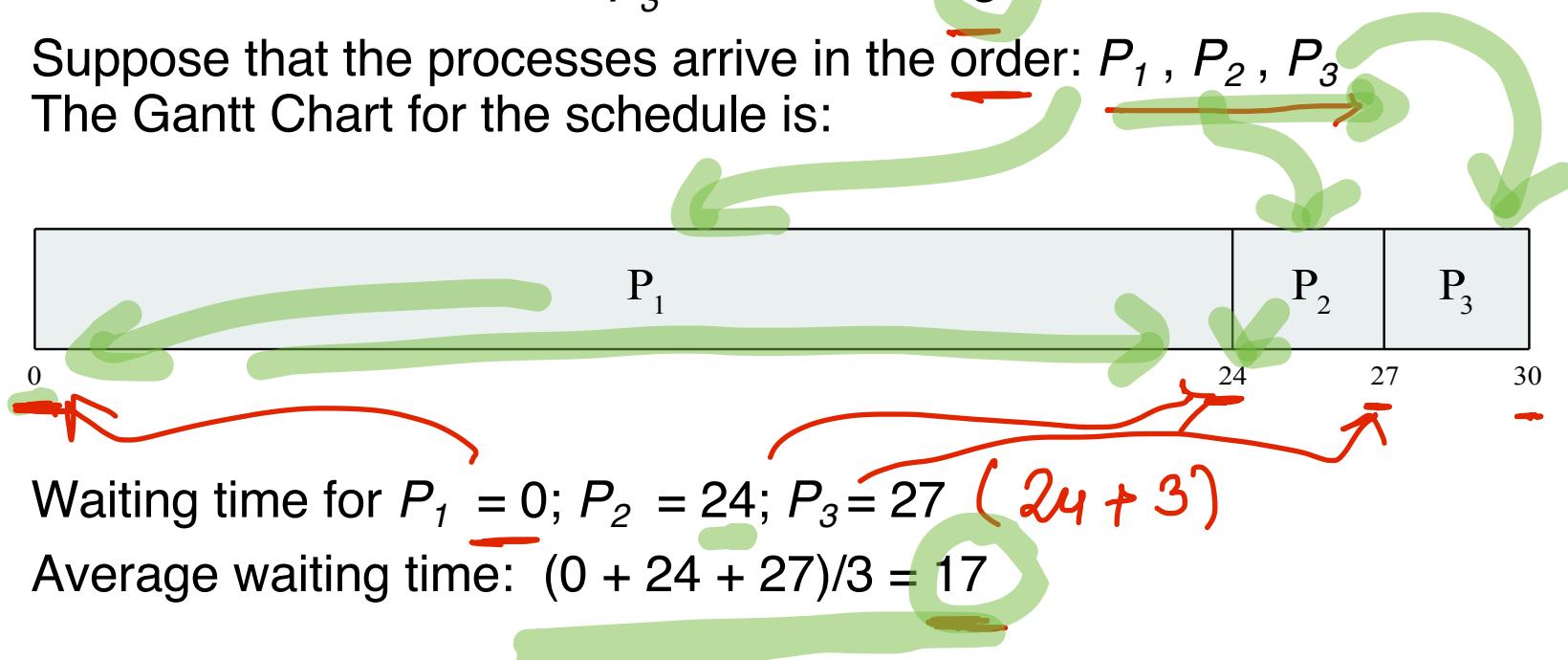
24

3

3

FIFO ✓
Queuing ✓

- Suppose that the processes arrive in the order: P_1, P_2, P_3 .
The Gantt Chart for the schedule is:



- Waiting time for $P_1 = 0$; $P_2 = 24$; $P_3 = 27$ ($24 + 3$)
- Average waiting time: $(0 + 24 + 27)/3 = 17$

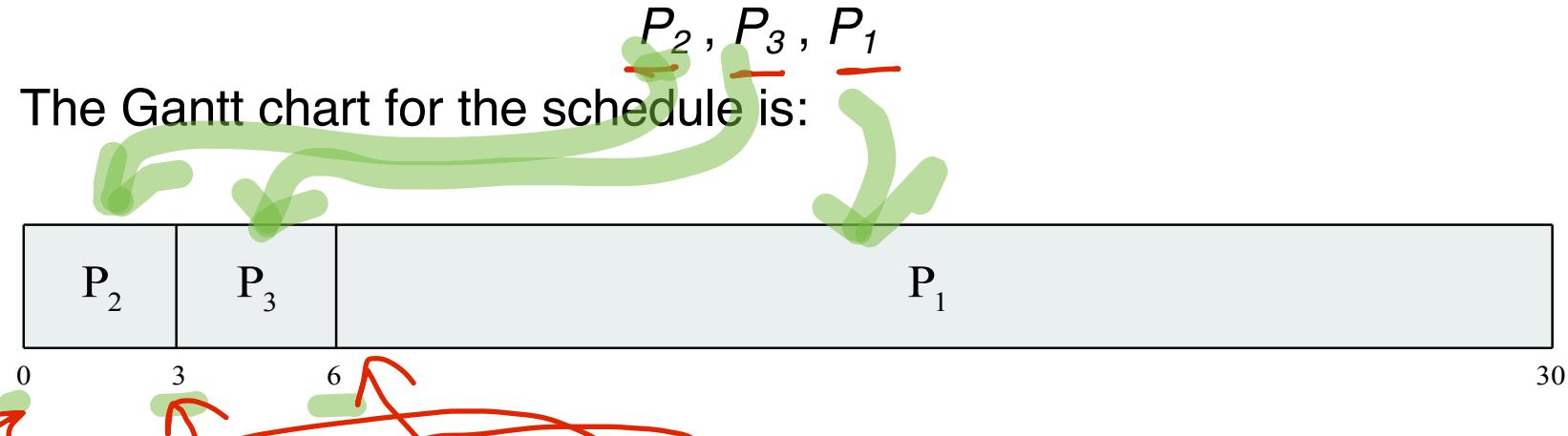




FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order:

- The Gantt chart for the schedule is:



- Waiting time for P₁ = 6; P₂ = 0; P₃ = 3
- Average waiting time: $(6 + 0 + 3)/3 = 3$
- Much better than previous case
- Convoy effect** - short process behind long process
 - Consider one CPU-bound and many I/O-bound processes

Best Case Scenario





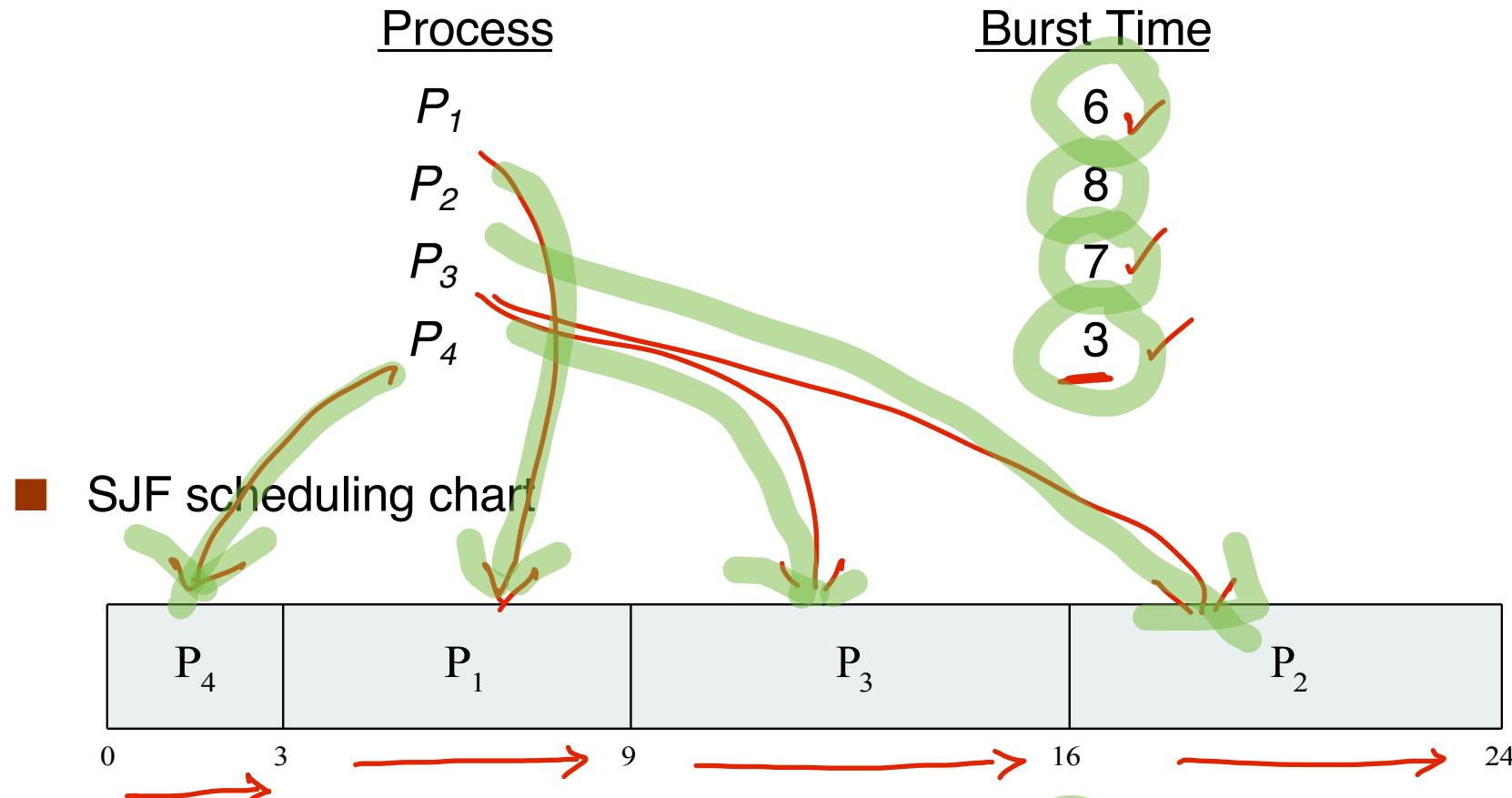
Shortest-Job-First (SJF) Scheduling

- Associate with each process the length of its next CPU burst
 - Use these lengths to schedule the process with the shortest time
- SJF is optimal – gives minimum average waiting time for a given set of processes
 - The difficulty is knowing the length of the next CPU request
 - Could ask the user



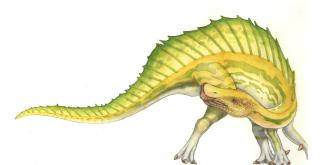


Example of SJF



Average waiting time = $(3 + 16 + 9 + 0) / 4 = 7$

P₁ P₂ P₃ P₄





Determining Length of Next CPU Burst

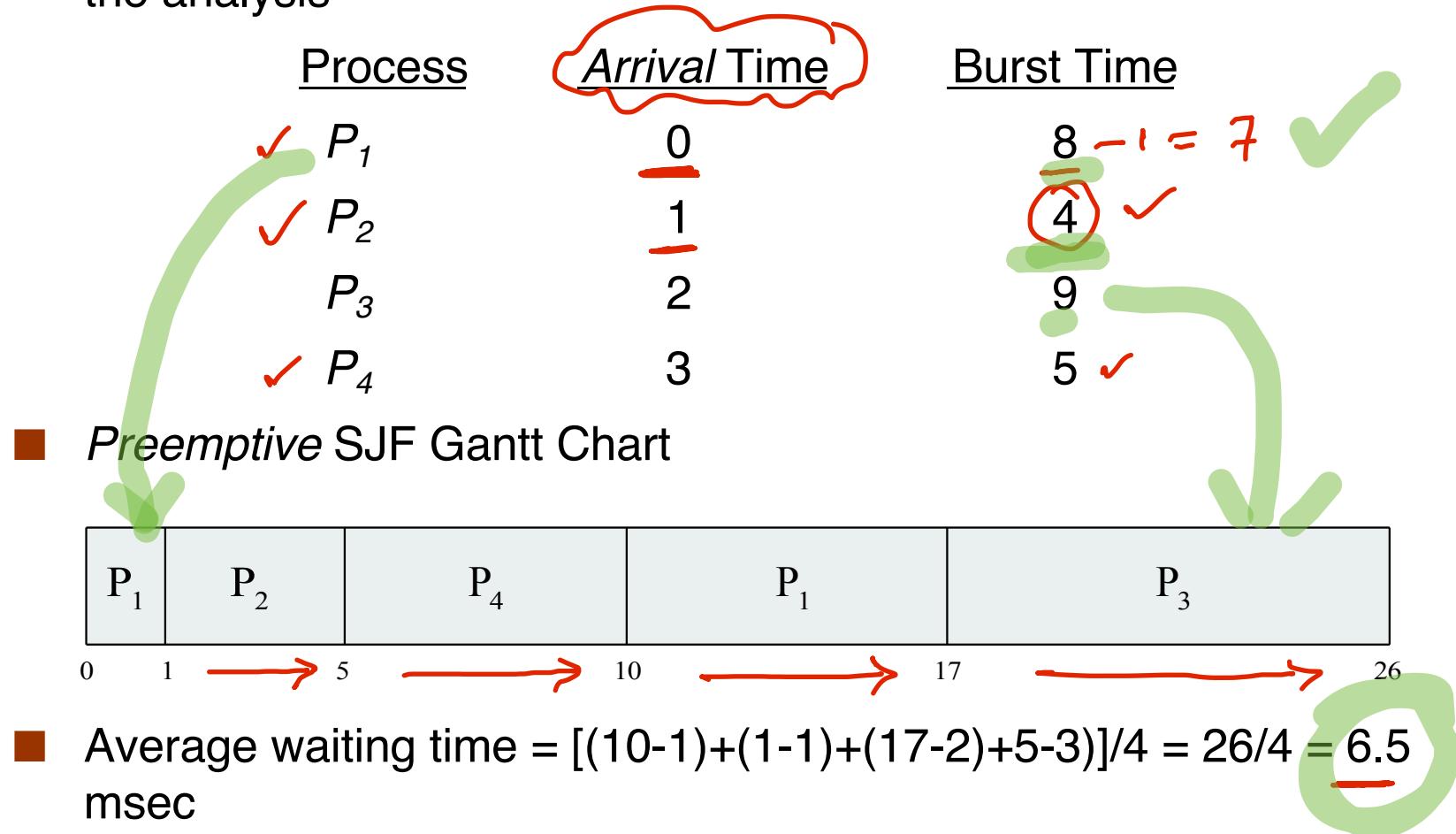
- Can only estimate the length – should be similar to the previous one
 - Then pick process with shortest predicted next CPU burst
- Can be done by using the length of previous CPU bursts, using exponential averaging
 1. t_n = actual length of n^{th} CPU burst
 2. τ_{n+1} = predicted value for the next CPU burst
 3. $\alpha, 0 \leq \alpha \leq 1$
 4. Define : $\tau_{n+1} = \alpha t_n + (1 - \alpha) \tau_n$.
- Commonly, α set to $\frac{1}{2}$
- Preemptive version called **shortest-remaining-time-first**





Example of Shortest-remaining-time-first

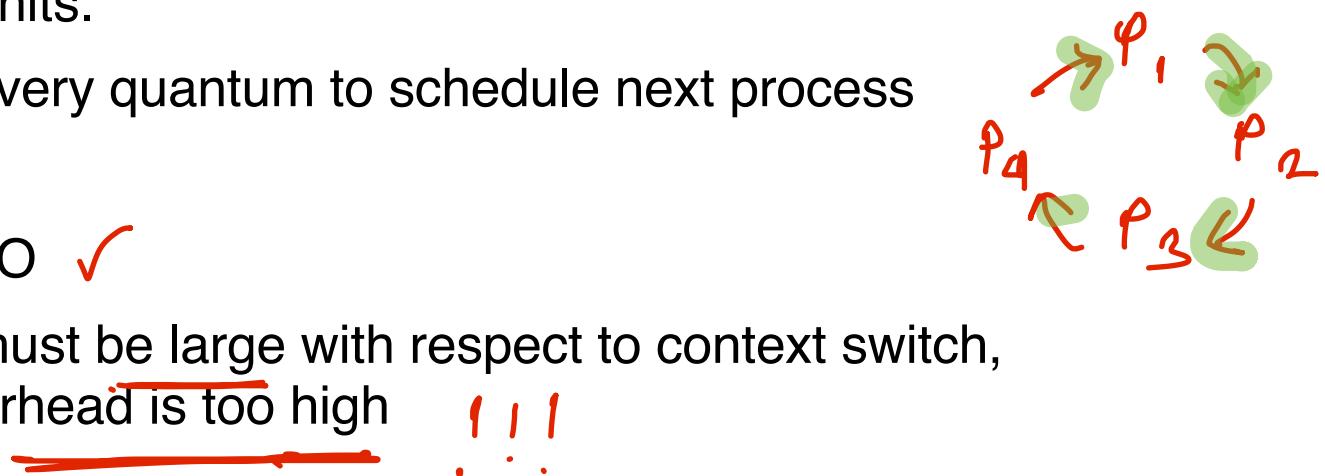
- Now we add the concepts of varying arrival times and preemption to the analysis





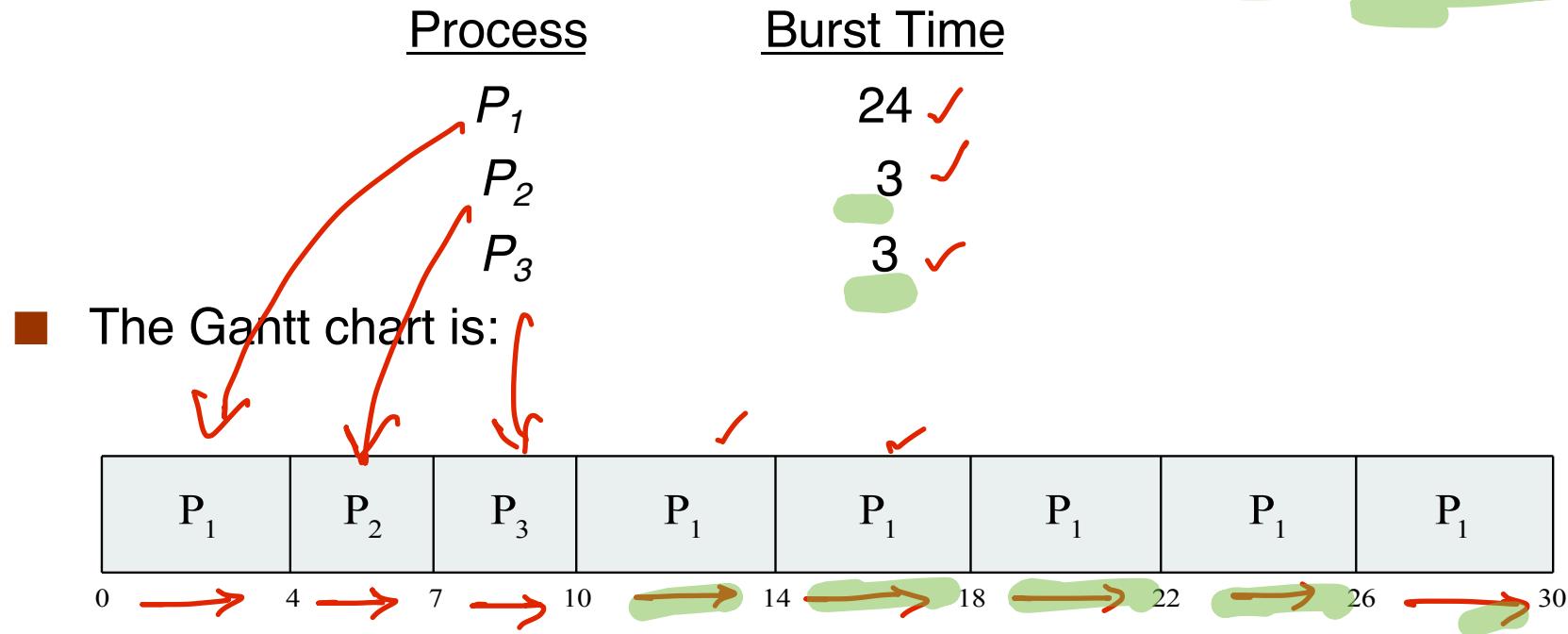
Round Robin (RR)

- Each process gets a small unit of CPU time (**time quantum** q), usually 10-100 milliseconds. After this time has elapsed, the process is preempted and added to the end of the ready queue.
- If there are n processes in the ready queue and the time quantum is q , then each process gets $1/n$ of the CPU time in chunks of at most q time units at once. No process waits more than $(n-1)q$ time units.
- Timer interrupts every quantum to schedule next process
- Performance
 - q large \Rightarrow FIFO ✓
 - q small \Rightarrow q must be large with respect to context switch, otherwise overhead is too high !!!

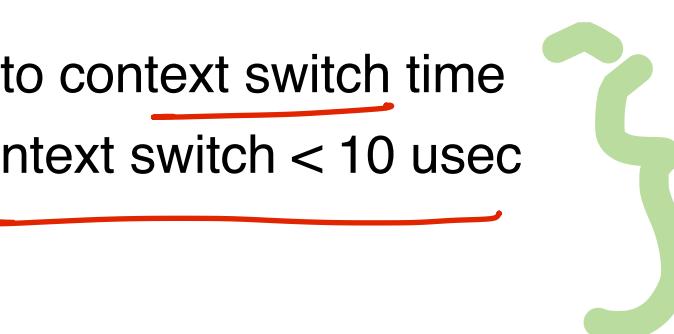




Example of RR with Time Quantum = 4



- Typically, higher average turnaround than SJF, but better response ✓
- q should be large compared to context switch time
- q usually 10ms to 100ms, context switch < 10 usec





Priority Scheduling

- A priority number (integer) is associated with each process
- The CPU is allocated to the process with the highest priority
(smallest integer = highest priority)

- Preemptive

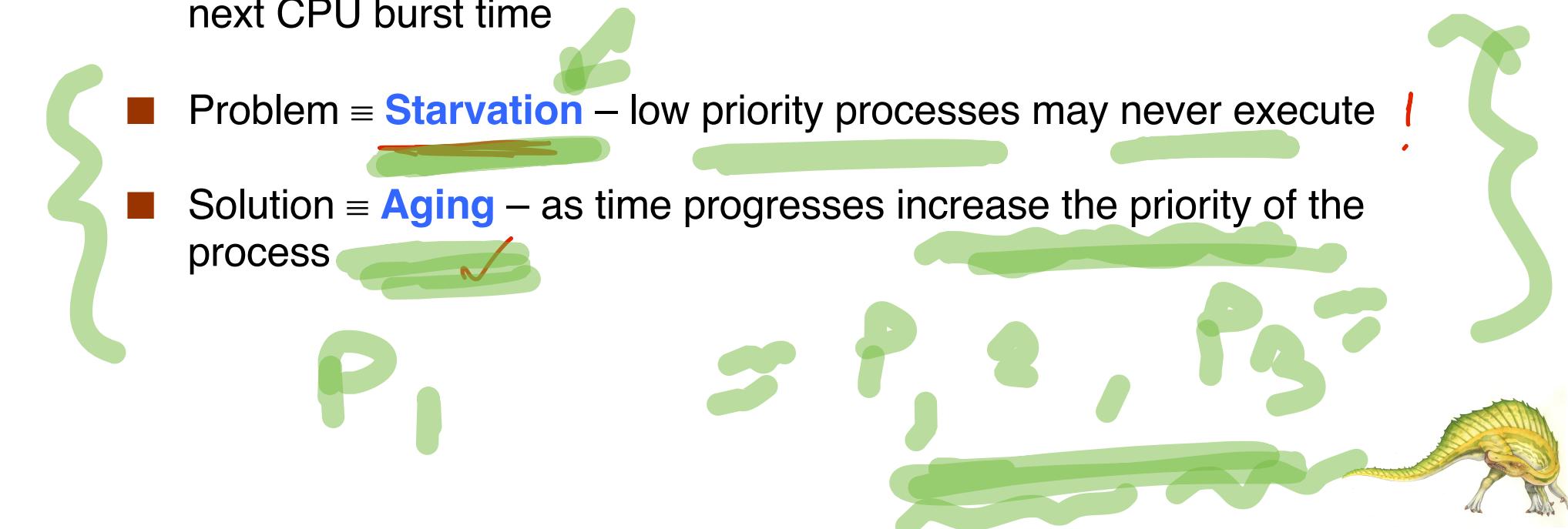


- Nonpreemptive

- SJF is priority scheduling where priority is the inverse of predicted next CPU burst time

- Problem = **Starvation** – low priority processes may never execute

- Solution = **Aging** – as time progresses increase the priority of the process

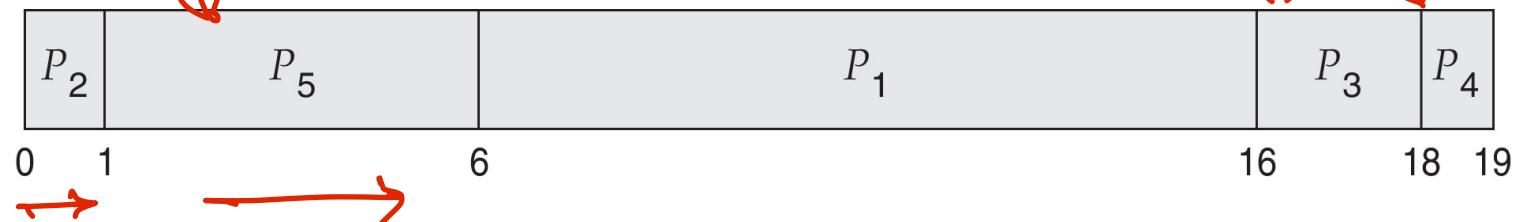




Example of Priority Scheduling

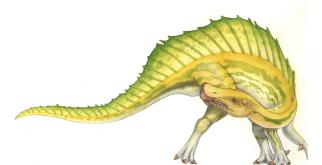
<u>Process</u>	<u>Burst Time</u>	<u>Priority</u>
P_1	10	3
P_2	1	1
P_3	2	4
P_4	1	5
P_5	5	2

■ Priority scheduling Gantt Chart



■ Average waiting time = 8.2 msec

8.2



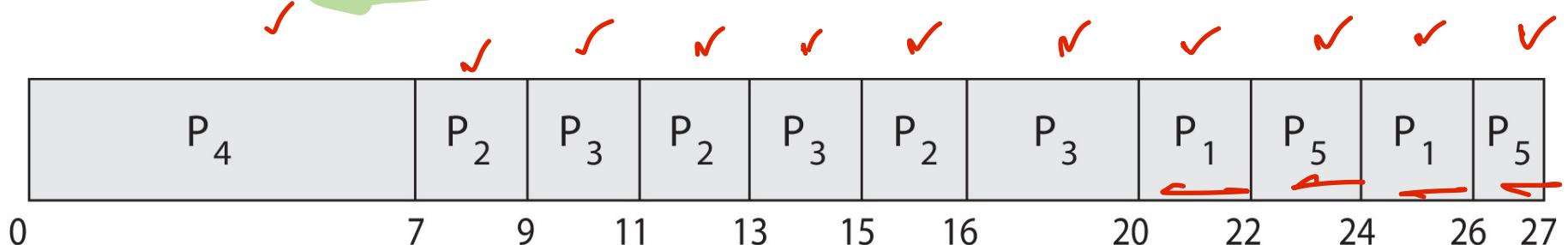


Priority Scheduling w/ Round-Robin

Process	Burst Time	Priority
P_1	4	3
P_2	5	2
P_3	8	2
P_4	7	1
P_5	3	3

- Run the process with the highest priority. Processes with the same priority run round-robin

- Gantt Chart with 2 ms time quantum



$$P_1 = 20 + 2 = 22$$

$$P_2 = 7 + 2 + 2 = 11$$

$$P_3 = 9 + 2 + 1 = 12$$

$$P_4 = 0$$

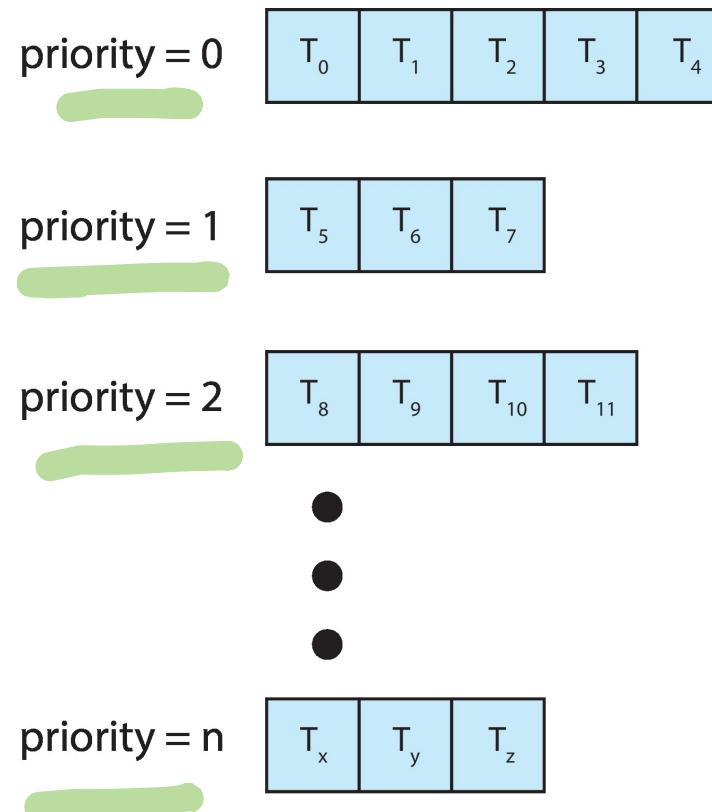
$$P_5 = 22 + 2 = 24$$

$$\text{AV. W. T.} = \frac{22 + 11 + 12 + 24 + 0}{5} = 13.8$$



Multilevel Queue

- With priority scheduling, have separate queues for each priority.
- Schedule the process in the highest-priority queue!



Red handwritten annotations:

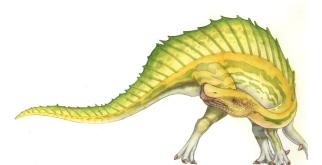
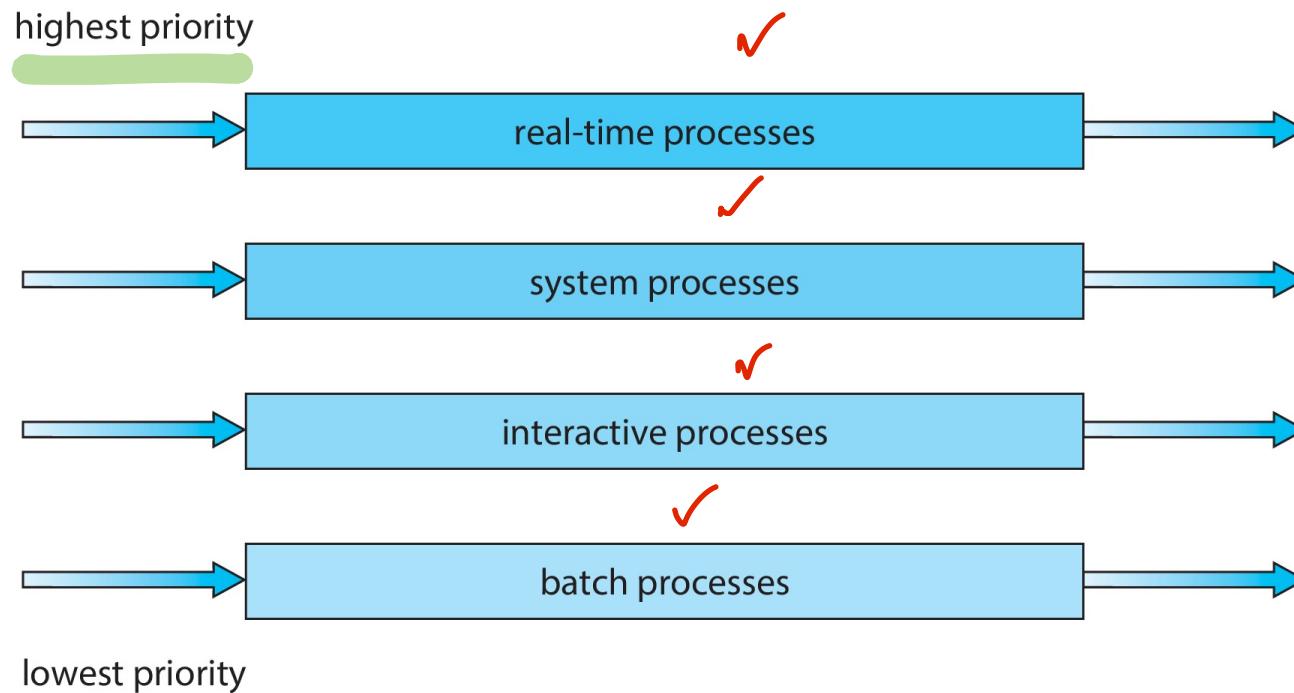
- A horizontal red line with the letters "A", "F", "J", "L", "W" written along it.
- The letters "M", "I", "D" written below the third level.
- A red circle drawn over the second level.





Multilevel Queue

- Prioritization based upon process type





Multilevel Feedback Queue

Starvation



- A process can move between the various queues; aging can be implemented this way
- Multilevel-feedback-queue scheduler defined by the following parameters:
 - number of queues ✓
 - ✓ ● scheduling algorithms for each queue
 - method used to determine when to upgrade a process → Promote
 - method used to determine when to demote a process
 - method used to determine which queue a process will enter when that process needs service





Example of Multilevel Feedback Queue

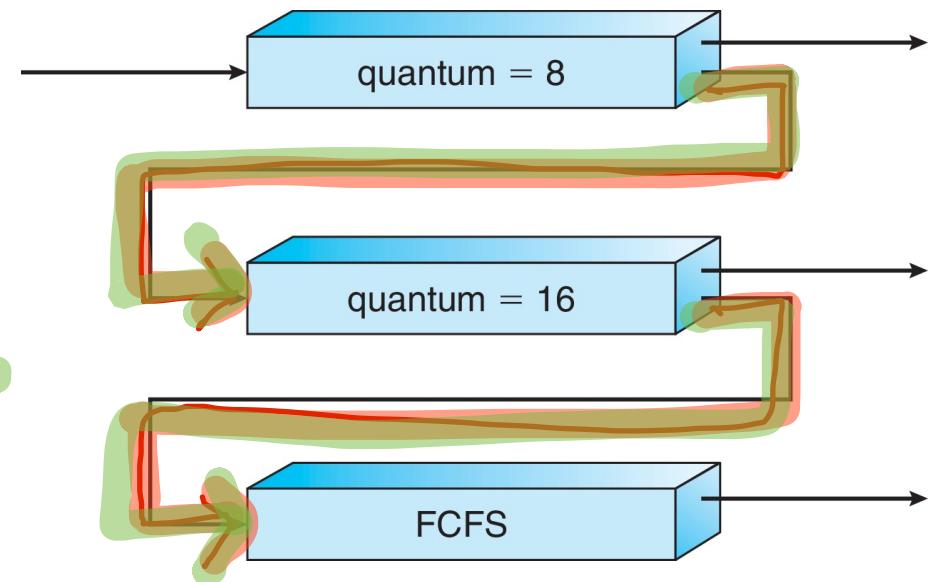
■ Three queues:

- Q_0 – RR with time quantum 8 milliseconds
- Q_1 – RR time quantum 16 milliseconds
- Q_2 – FCFS

For example,

■ Scheduling

- A new job enters queue Q_0 which is served FCFS
 - When it gains CPU, job receives 8 milliseconds
 - If it does not finish in 8 milliseconds, job is moved to queue Q_1
- At Q_1 job is again served FCFS and receives 16 additional milliseconds
 - If it still does not complete, it is preempted and moved to queue Q_2





Thread Scheduling

- Distinction between user-level and kernel-level threads
 - When threads supported, threads scheduled, not processes
 - Many-to-one and many-to-many models, thread library schedules user-level threads to run on LWP (*Light-weight Process*)
 - Known as **process-contention scope (PCS)** since scheduling competition is within the process
 - Typically done via priority set by programmer · in applications
 - Kernel thread scheduled onto available CPU is **system-contention scope (SCS)** – competition among all threads in system
- 4
4





Pthread Scheduling

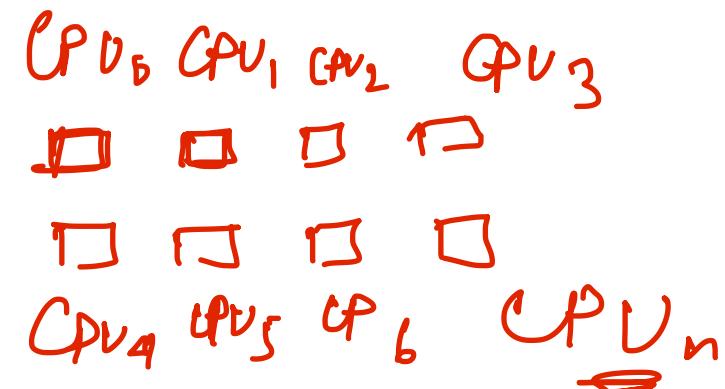
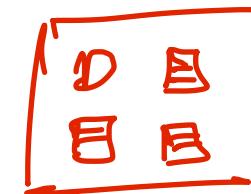
- API allows specifying either PCS or SCS during thread creation
 - PTHREAD_SCOPE_PROCESS schedules threads using PCS scheduling
 - PTHREAD_SCOPE_SYSTEM schedules threads using SCS scheduling
- Can be limited by OS – Linux and macOS only allow PTHREAD_SCOPE_SYSTEM 





Multiple-Processor Scheduling

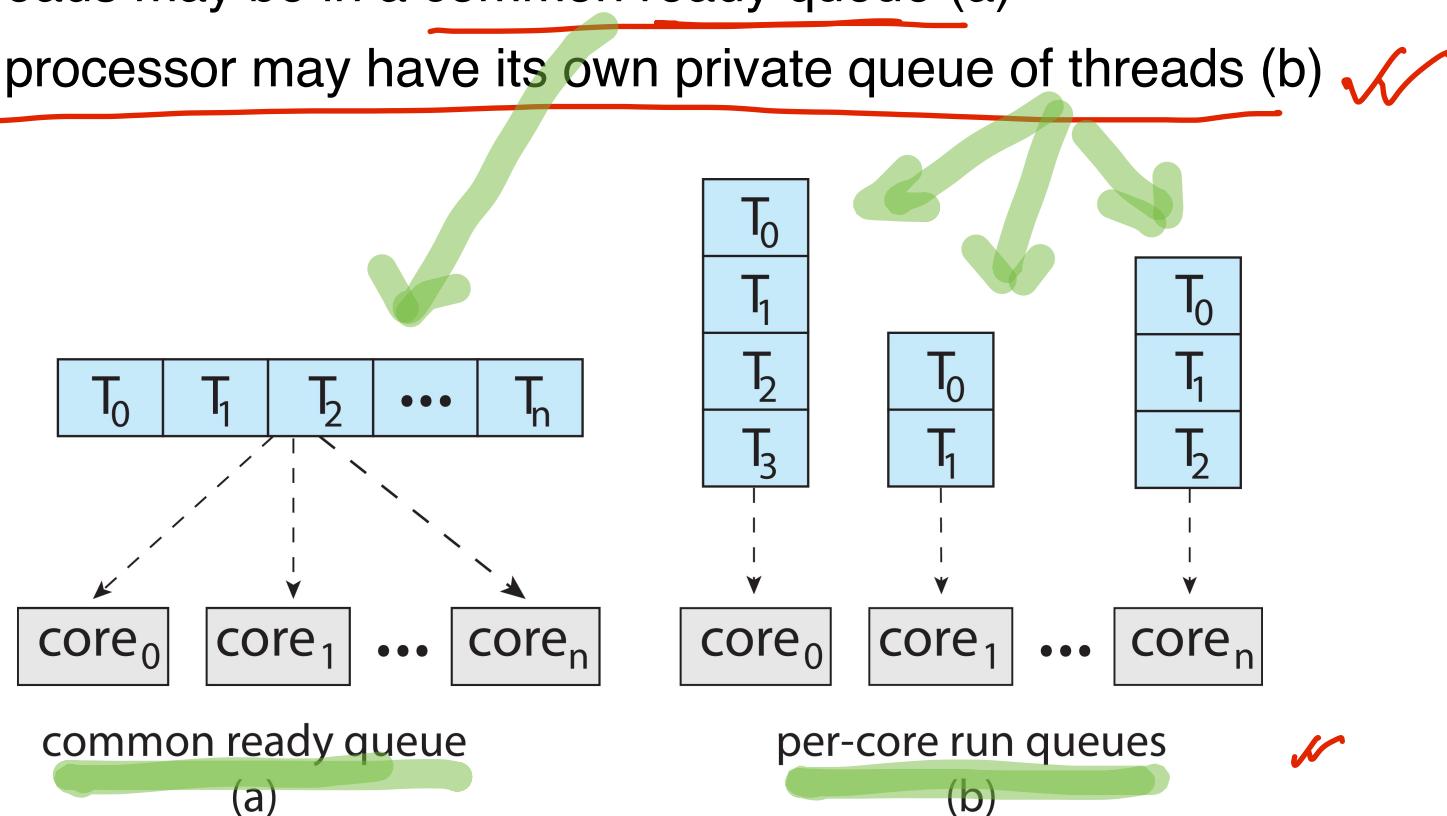
- CPU scheduling more complex when multiple CPUs are available
- Multiprocess may be any one of the following architectures: ✓
 - Multicore CPUs ✓
 - Multithreaded cores ✓
 - NUMA systems ✓
 - Heterogeneous multiprocessing ✓





Multiple-Processor Scheduling

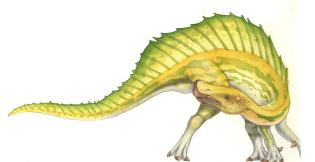
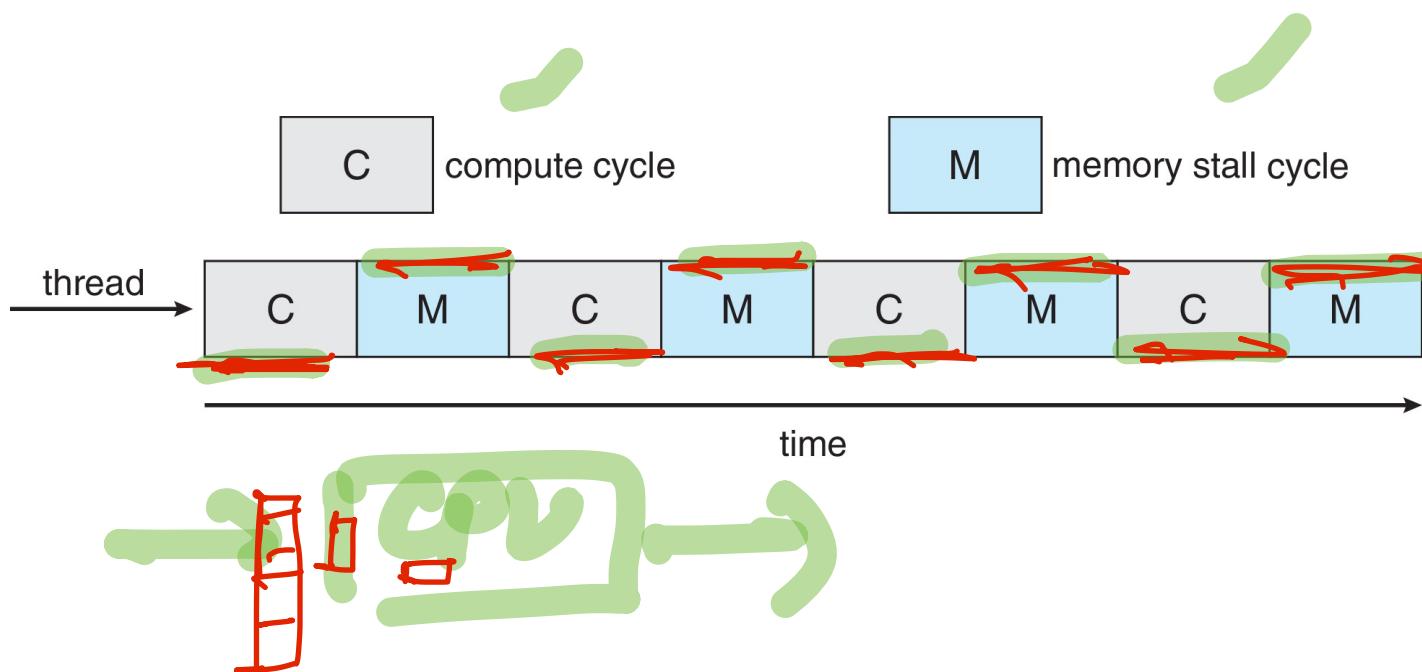
- Symmetric multiprocessing (SMP) is where each processor is self scheduling.
- All threads may be in a common ready queue (a)
- Each processor may have its own private queue of threads (b)





Multicore Processors

- Recent trend to place multiple processor cores on same physical chip
- Faster and consumes less power
- Multiple threads per core also growing
 - Takes advantage of memory stall to make progress on another thread while memory retrieve happens

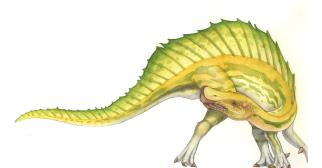
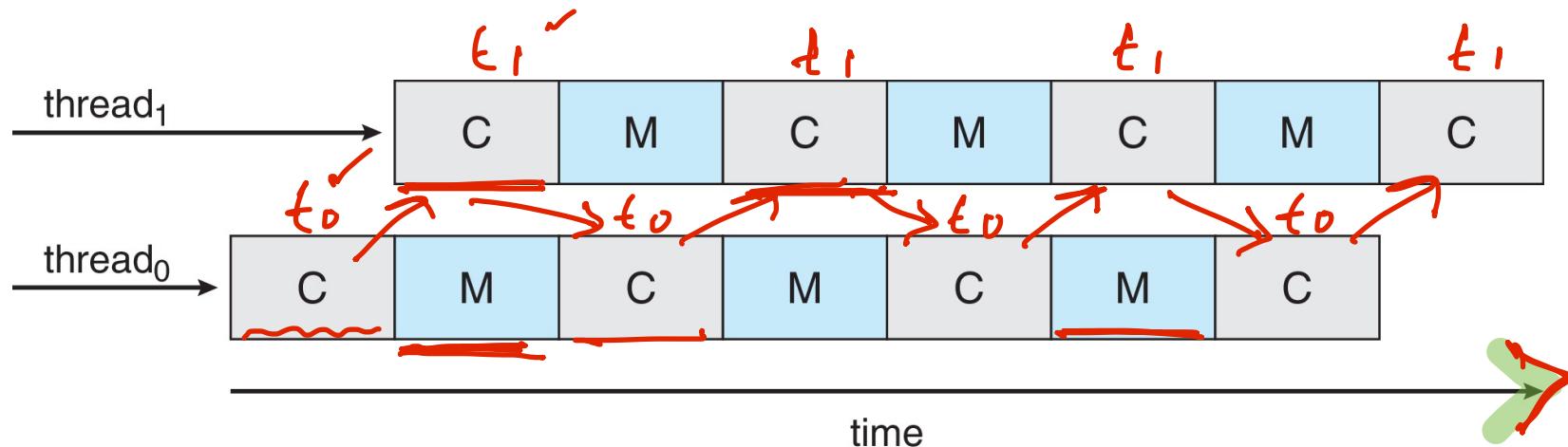




Multithreaded Multicore System

Each core has > 1 hardware threads.

If one thread has a memory stall, switch to another thread!





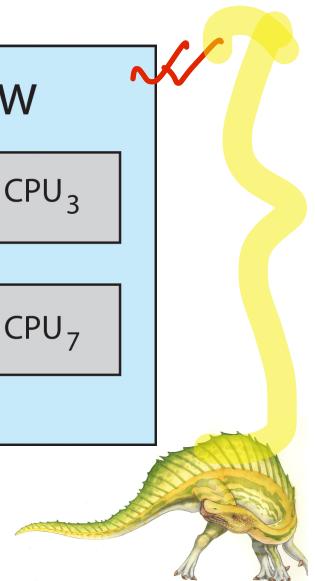
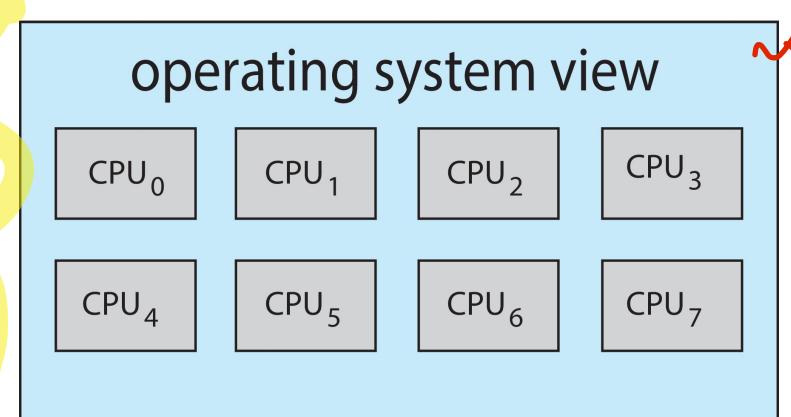
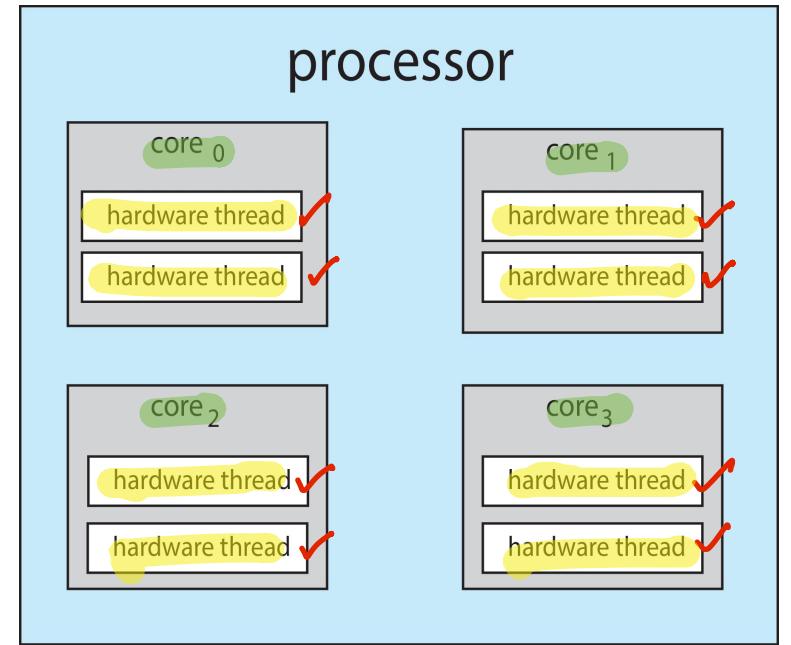
Multithreaded Multicore System

- Chip-multithreading (CMT) assigns each core multiple hardware threads. (Intel refers to this as hyperthreading.)



- On a quad-core system with 2 hardware threads per core, the operating system sees 8 logical processors.

4 × 2 = 8





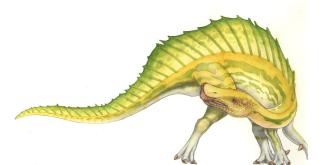
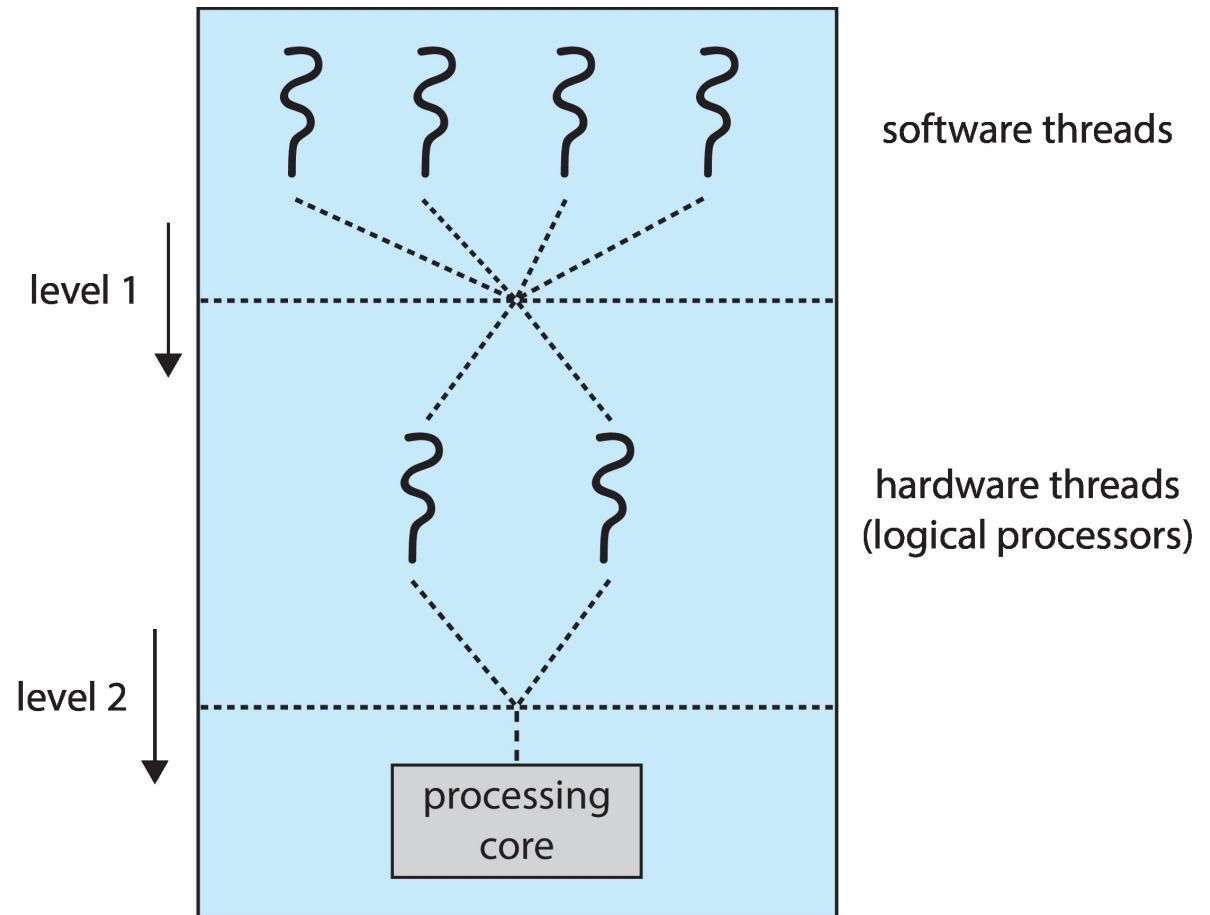
Multithreaded Multicore System

- Two levels of scheduling:

1. The operating system deciding which software thread to run on a logical CPU



2. How each core decides which hardware thread to run on the physical core.





Multiple-Processor Scheduling – Load Balancing

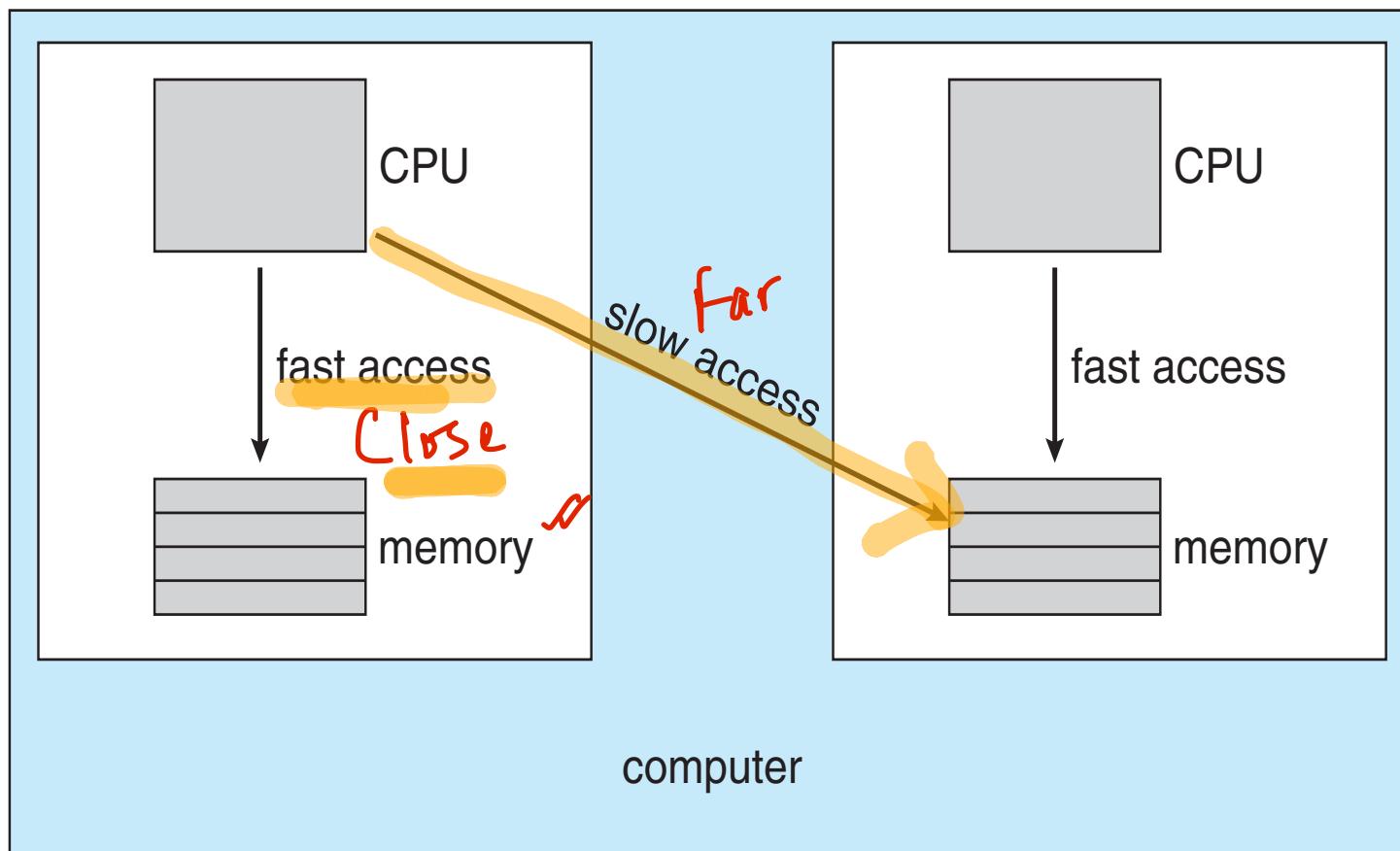
- If SMP, need to keep all CPUs loaded for efficiency
- **Load balancing** attempts to keep workload evenly distributed
- **Push migration** – periodic task checks load on each processor, and if found pushes task from overloaded CPU to other CPUs
- **Pull migration** – idle processors pulls waiting task from busy processor





NUMA and CPU Scheduling

If the operating system is **NUMA-aware**, it will assign memory closes to the CPU the thread is running on.

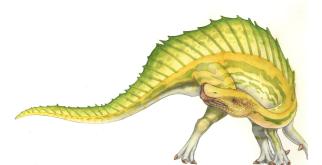




Real-Time CPU Scheduling

- Can present obvious challenges
- ✓ ■ **Soft real-time systems** – Critical real-time tasks have the highest priority, but no guarantee as to when tasks will be scheduled
- ✓ ■ **Hard real-time systems** – task must be serviced by its deadline

1





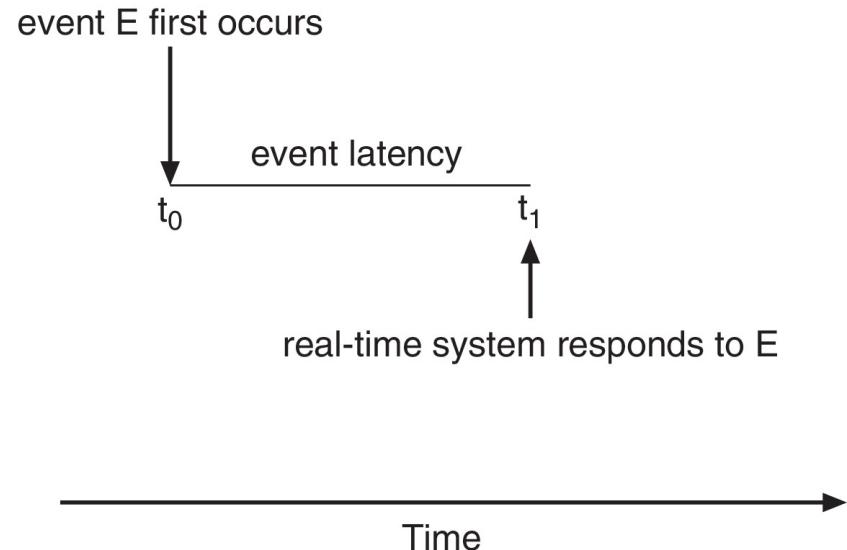
Real-Time CPU Scheduling

- Event latency – the amount of time that elapses from when an event occurs to when it is serviced.
- Two types of latencies affect performance

- ✓ 1. **Interrupt latency** – time from arrival of interrupt to start of routine that services interrupt
- ✓ 2. **Dispatch latency** – time for schedule to take current process off CPU and switch to another



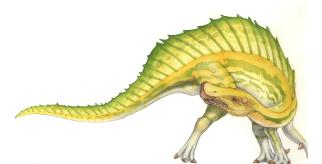
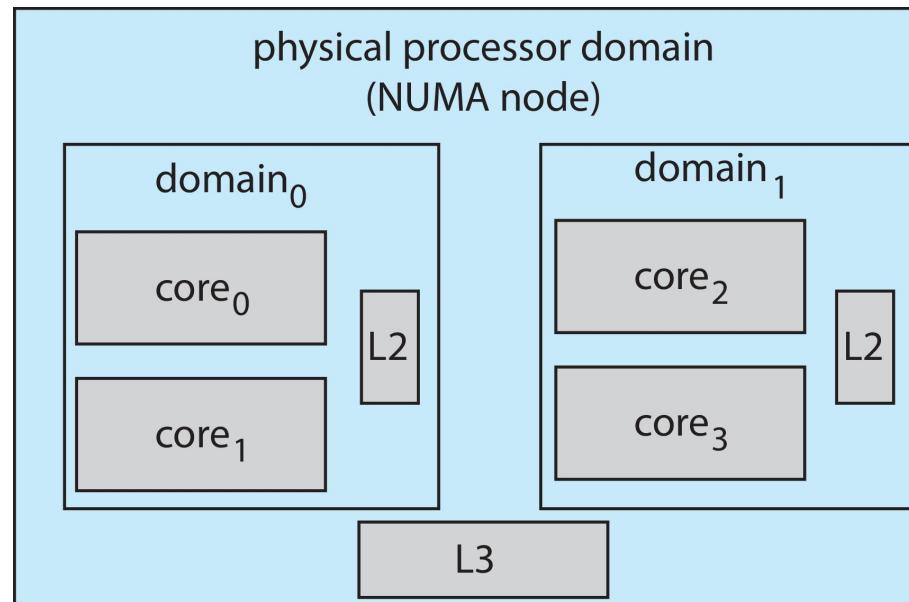
(Process switching time)





Linux Scheduling (Cont.)

- Linux supports load balancing, but is also NUMA-aware.
- **Scheduling domain** is a set of CPU cores that can be balanced against one another.
- Domains are organized by what they share (i.e. cache memory.) Goal is to keep threads from migrating between domains.



End of Chapter 5

