# CPU Scheduling (1)

Dr. Jun Zheng
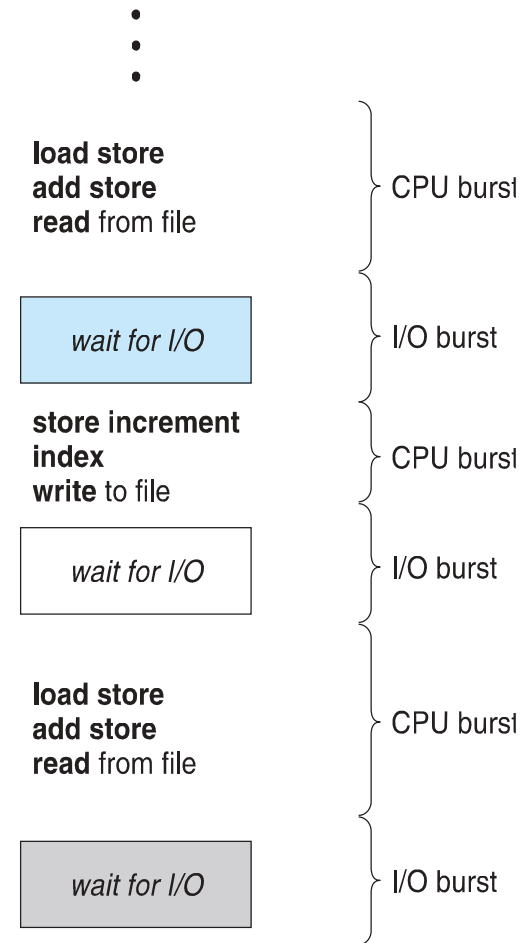
CSE325 Principles of Operating Systems

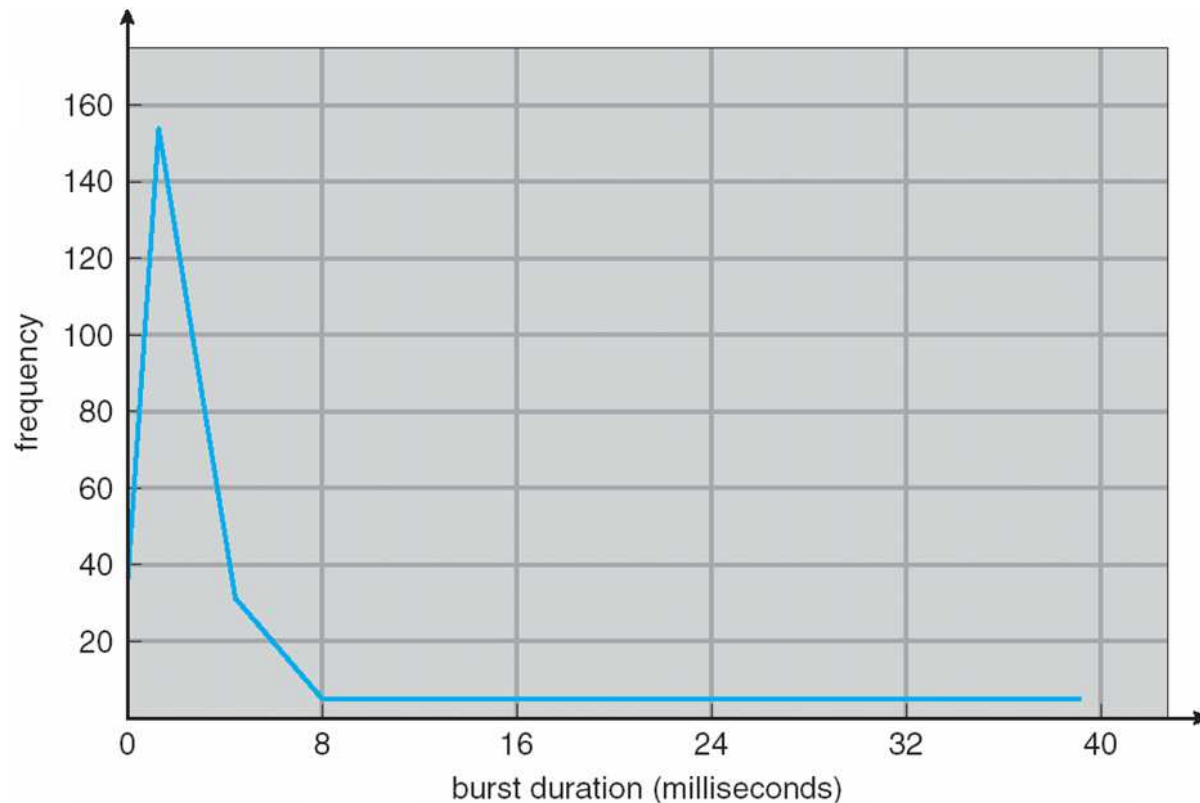9/9/2019

**NEW MEXICO TECH**

SCIENCE · ENGINEERING · RESEARCH · UNIVERSITY

# 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

| | |
|---|---|
| load store<br>add store<br>**read** from file | CPU burst |
| *wait for I/O* | I/O burst |
| store increment<br>index<br>**write** to file | CPU burst |
| *wait for I/O* | I/O burst |
| load store<br>add store<br>**read** from file | CPU burst |
| *wait for I/O* | I/O burst |

# Histogram of CPU-burst Times

# CPU Scheduler

❑ **Short-term scheduler** selects from among the processes in ready queue, and allocates the CPU 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 (e.g. I/O request)
  2. Switches from running to ready state (e.g. an interrupt)
  3. Switches from waiting to ready (e.g. completion of I/O)
  4. Terminates

❑ Scheduling under 1 and 4 is **nonpreemptive**

❑ All other scheduling is **preemptive**

# 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

# 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

# Gantt Chart

Illustrates how jobs are scheduled over time
on CPU

Example:

| A | B | C |
|---|---|---|

0                                    10    12                    16

Time  →

# First- Come, First-Served (FCFS) Scheduling

| Process | Burst Time |
|---------|------------|
| $P_1$   | 24         |
| $P_2$   | 3          |
| $P_3$   | 3          |

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

| $P_1$ | $P_2$ | $P_3$ |
|-------|-------|-------|

0　　　　　　　　　　　　　　　　　　24　　　27　　　30

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

# FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order:

$$P_2, P_3, P_1$$

❑ The Gantt chart for the schedule is:

| P$_2$ | P$_3$ | P$_1$ |
|:---:|:---:|:---:|

0    3    6                                            30

❑ Waiting time for $P_1 = 6$; $P_2 = 0$; $P_3 = 3$

❑ Average waiting time:   $(6 + 0 + 3)/3 = 3$

❑ Much better than previous case

❑ **Convoy effect** - short process behind long process

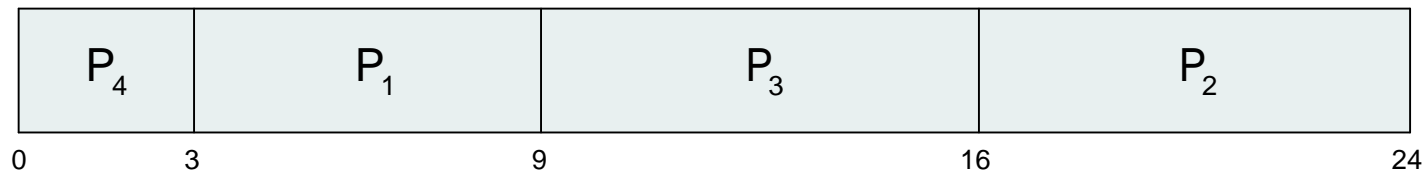  ❑   Consider one CPU-bound and many I/O-bound processes

# 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

| Process | Burst Time |
|---------|-----------|
| $P_1$ | 6 |
| $P_2$ | 8 |
| $P_3$ | 7 |
| $P_4$ | 3 |

- SJF scheduling chart

| $P_4$ | $P_1$ | $P_3$ | $P_2$ |
|:-----:|:-----:|:-----:|:-----:|
| 0      3 | 9 | 16 | 24 |

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

NEW MEXICO TECH
SCIENCE · ENGINEERING · RESEARCH · UNIVERSITY

# 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
    1. $t_n$ = actual length of $n^{th}$ CPU burst
    2. $\tau_{n+1}$ = predicted value for the next CPU burst
    3. $\alpha, 0 \leq \alpha \leq 1$
    4. Define : $\tau_{n+1} = \alpha\, t_n + (1-\alpha)\tau_n$.
- ❑ Can be done by using the length of previous CPU bursts, using exponential averaging
- ❑ Commonly, $\alpha$ set to ½
- ❑ Preemptive version called **shortest-remaining-time-first**

NEW MEXICO TECH
SCIENCE · ENGINEERING · RESEARCH · UNIVERSITY