

# Operating System

*Dr. GuoJun LIU*

Harbin Institute of Technology

<http://guojunos.hit.edu.cn>

# Chapter 09

## *Scheduling*

调度

# Learning Objectives

---

- **Explain the differences among long-, medium-, and short-term scheduling**
- **Assess the performance of different scheduling policies**

# Outline

---

## ■ Types of Processor Scheduling

- Long-Term Scheduling
- Medium-Term Scheduling
- Short-Term Scheduling

## ■ Scheduling Algorithms

- Short-Term Scheduling Criteria
- The Use of Priorities
- Alternative Scheduling Policies

# Processor Scheduling

---

- **Aim is to assign processes to be executed by the processor in a way that meets system objectives**
- **system objectives**
  - response time
  - throughput
  - processor efficiency

# Types of Scheduling

---

## ■ Long-term scheduling

- The decision to add to the **pool** of processes to be executed

## ■ Medium-term scheduling

- The decision to add to the number of processes that are **partially** or **fully** in **main memory**

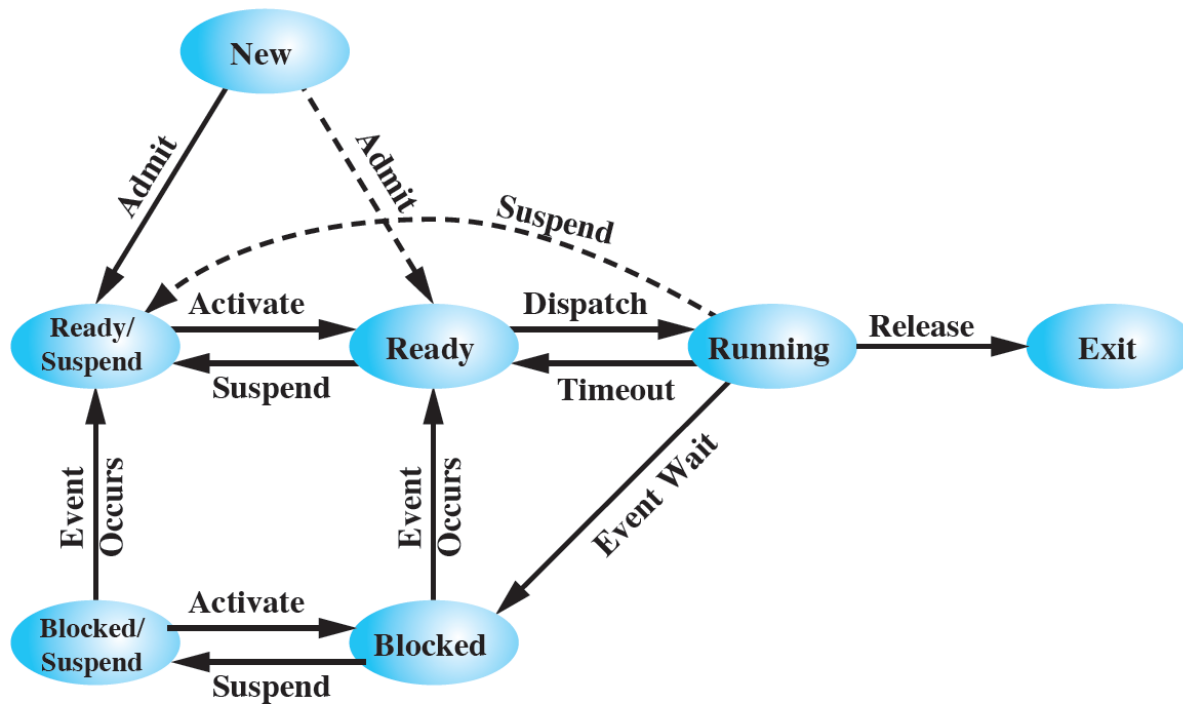
## ■ Short-term scheduling

- The decision as to which available process will be executed by the processor

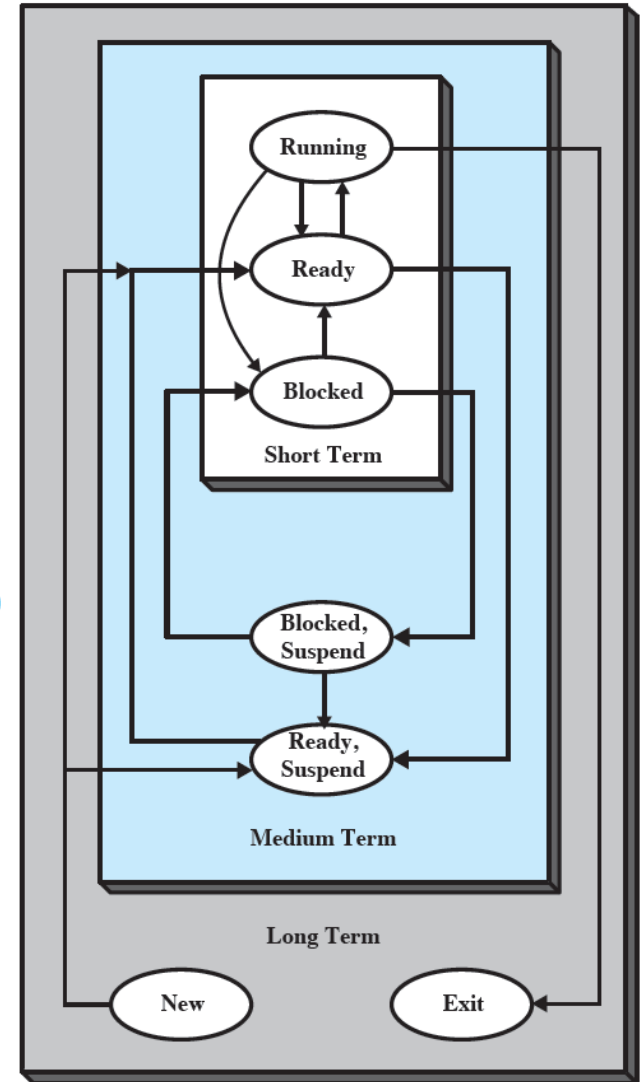
## ■ I/O scheduling

- The decision as to which process's pending I/O request shall be handled by an available I/O device

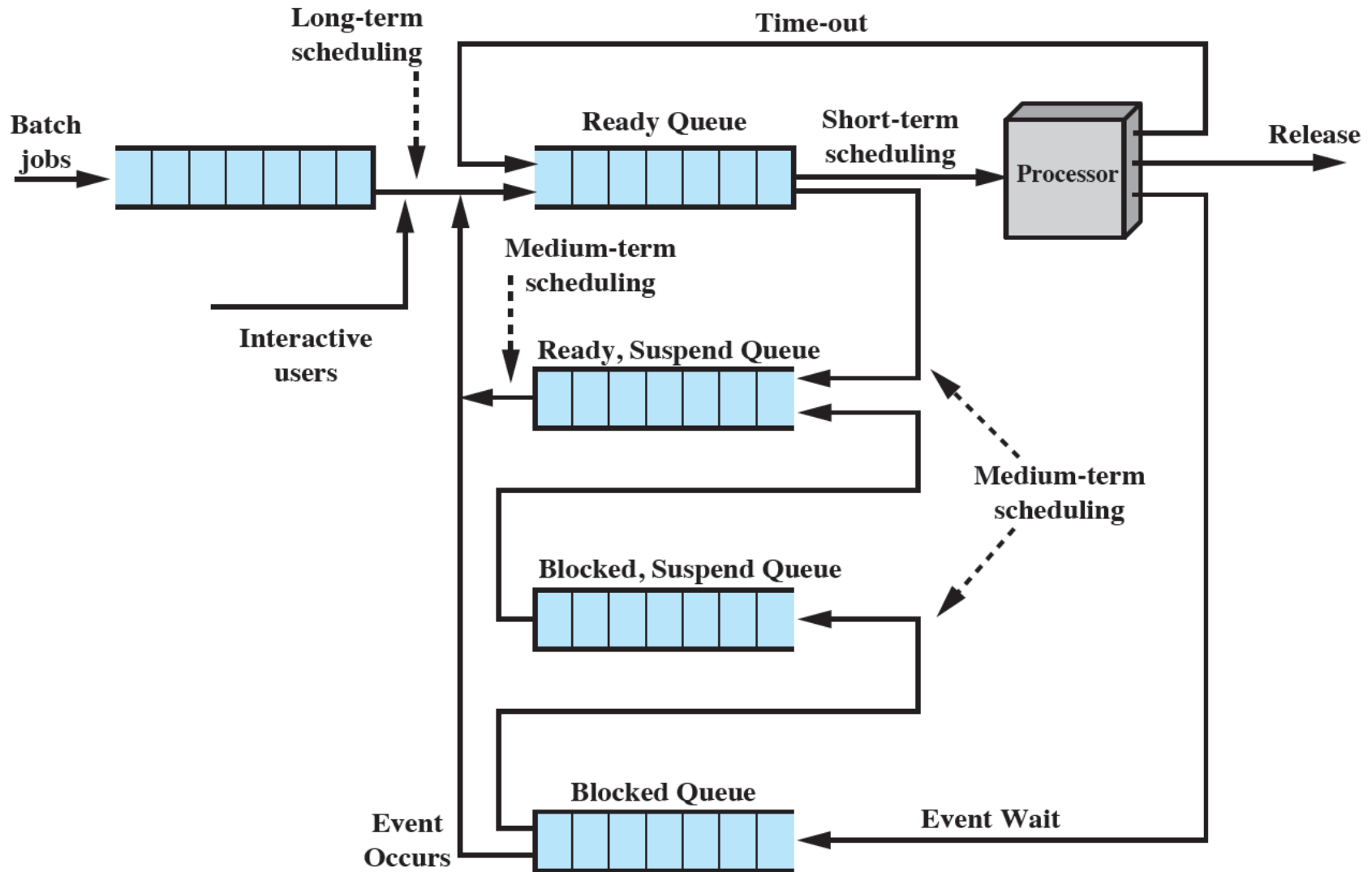
# Nesting of Scheduling Functions



(b) With Two Suspend States



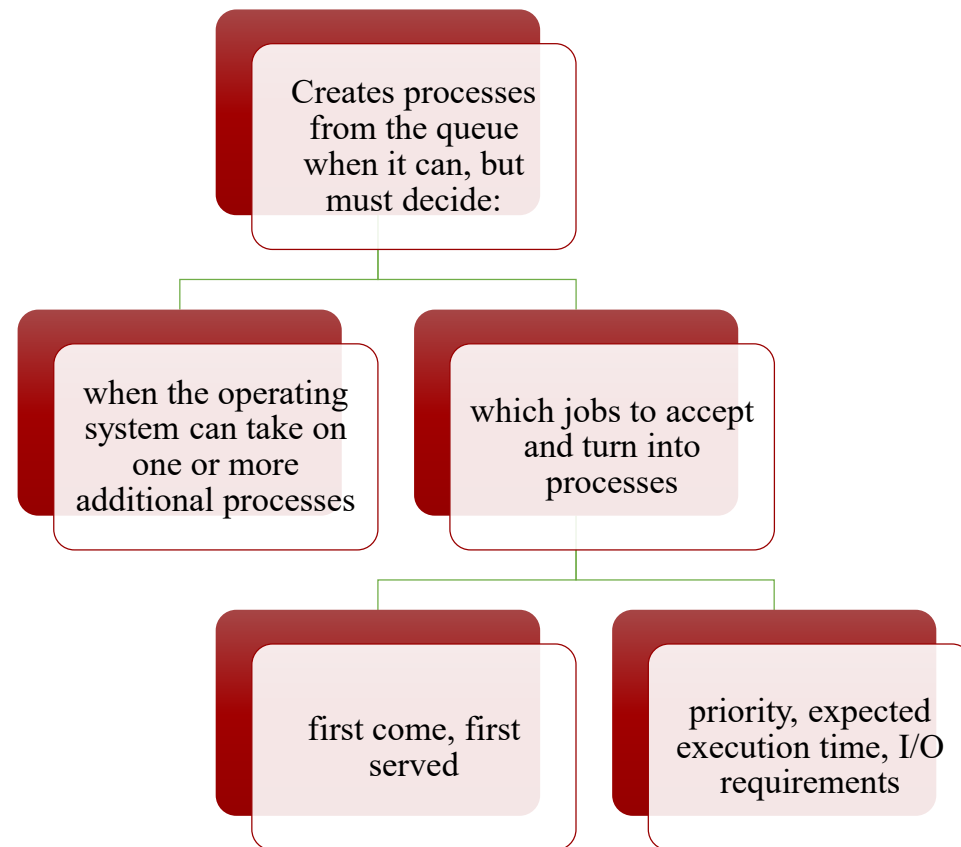
# Queuing Diagram for Scheduling





# Long-Term Scheduler

- **Determines which programs are admitted to the system for processing**
- **Controls the degree of multiprogramming**
  - the more processes that are created, the smaller the percentage of time that each process can be executed
  - may limit to provide satisfactory service to the current set of processes



# Medium-Term Scheduling

---

- Part of the swapping function
- Swapping-in decisions are based on the need to manage the degree of multiprogramming
  - considers the memory requirements of the swapped-out processes

# Short-Term Scheduling

---

- Known as the **dispatcher**, Executes most frequently
- Makes the fine-grained decision of which process to execute next
- Invoked when an event occurs that may lead to the blocking of the current process or that may provide an opportunity to preempt a currently running process in favor of another
  - Clock interrupts
  - I/O interrupts
  - Operating system calls
  - Signals (e.g., semaphores)

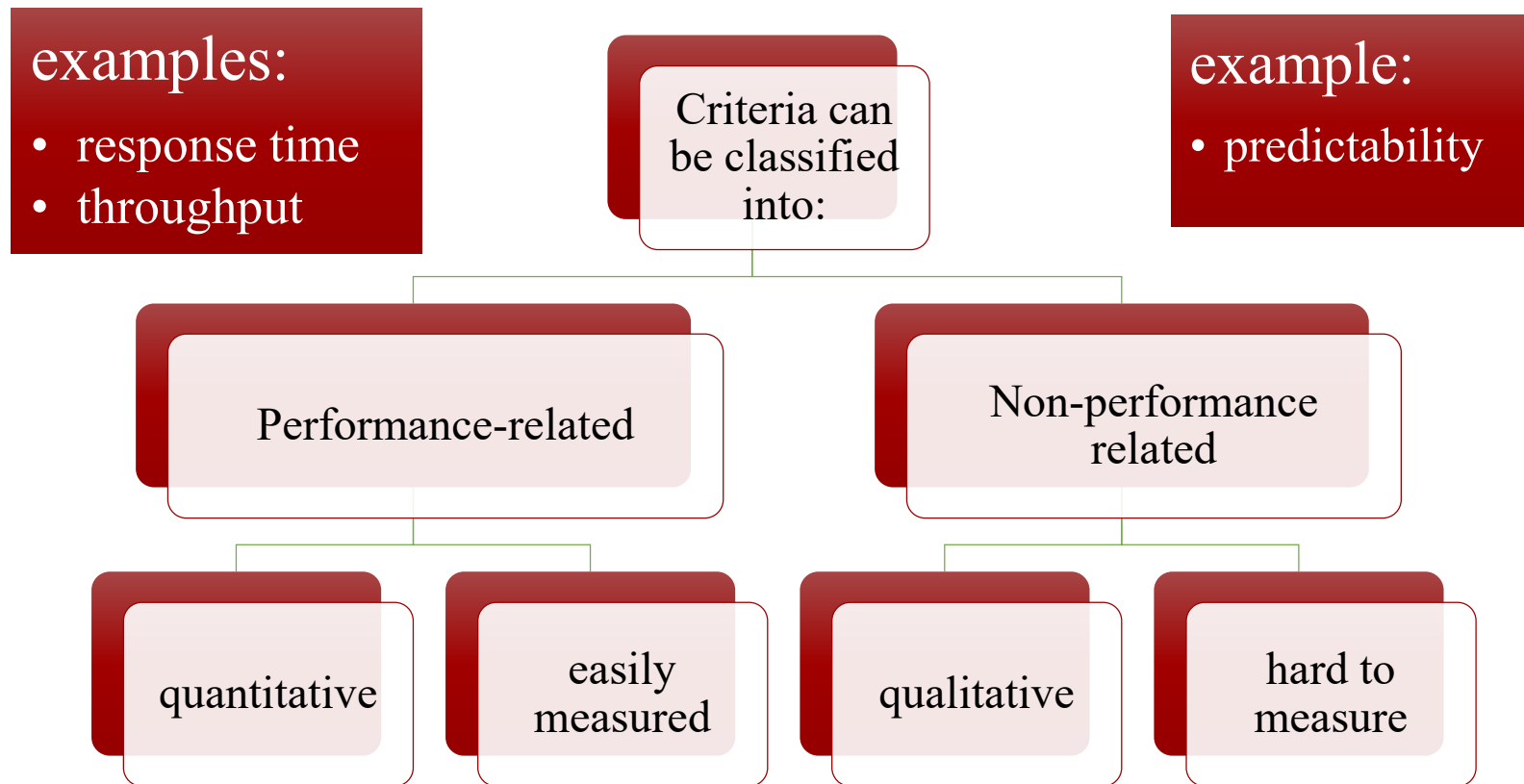
# Short Term Scheduling Criteria

---

- **Main objective is to allocate processor time to optimize certain aspects of system behavior**
- **A set of criteria is needed to evaluate the scheduling policy**
- **User-oriented criteria**
  - relate to the behavior of the system as perceived by the individual user or process (such as response time in an interactive system)
  - important on virtually all systems
- **System-oriented criteria**
  - focus in on effective and efficient utilization of the processor (rate at which processes are completed)
  - generally of minor importance on single-user systems

# Short-Term Scheduling Criteria: Performance

---



# Scheduling Criteria

---

## ■ User Oriented

### ➤ Turnaround time

- This is the interval of time between the submission of a process and its completion. Includes actual execution time plus time spent waiting for resources, including the processor. This is an appropriate measure for a batch job

### ➤ Deadlines

- When process completion deadlines can be specified, the scheduling discipline should subordinate other goals to that of maximizing the percentage of deadlines met

### ➤ Response time

- For an interactive process, this is the time from the submission of a request until the response begins to be received.
- The scheduling discipline should attempt to achieve low response time and to maximize the number of interactive users receiving acceptable response time.

### ➤ Predictability

- A given job should run in about the same amount of time and at about the same cost regardless of the load on the system. A wide variation in response time or turnaround time is distracting to users

# Scheduling Criteria

---

## ■ System Oriented

### ➤ Throughput

- The scheduling policy should attempt to maximize the number of processes completed per unit of time. This is a measure of how much work is being performed.

### ➤ Processor utilization

- This is the percentage of time that the processor is busy. For an expensive shared system, this is a significant criterion. In single-user systems and in some other systems, such as real-time systems, this criterion is less important than some of the others.

### ➤ Fairness

- In the absence of guidance from the user or other system-supplied guidance, processes should be treated the same, and no process should suffer starvation.

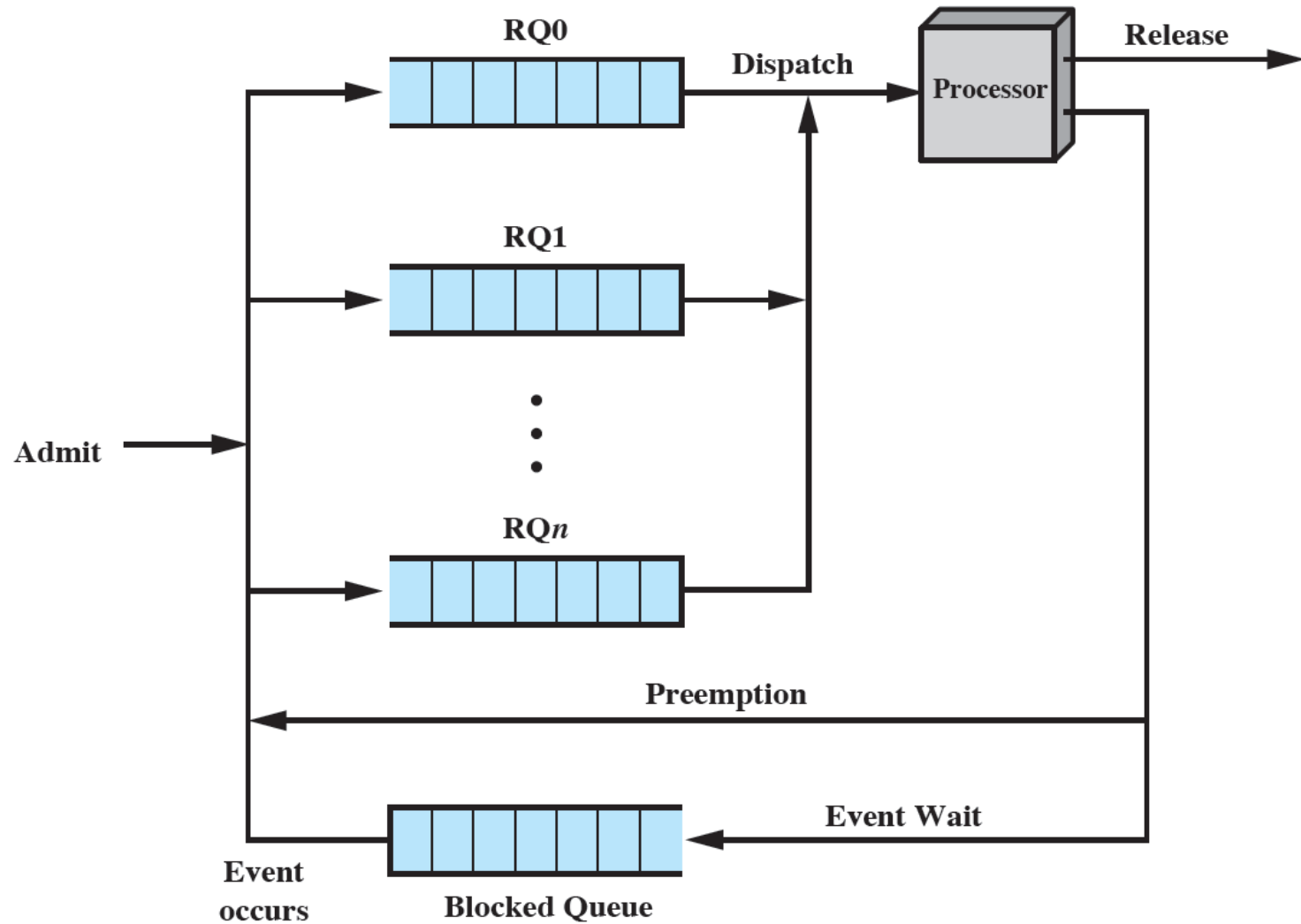
### ➤ Enforcing priorities

- When processes are assigned priorities, the scheduling policy should favor higher-priority processes

### ➤ Balancing resources

- The scheduling policy should keep the resources of the system busy. Processes that will underutilize stressed resources should be favored. This criterion also involves medium-term and long-term scheduling

# Priority Queueing





# Characteristics of Various Scheduling Policies

	FCFS	Round robin	SPN	SRT	HRRN	Feedback
<b>Selection function</b>	$\max[w]$	constant	$\min[s]$	$\min[s - e]$	$\max[(w+s)/s]$	(see text)
<b>Decision mode</b>	Non-preemptive	Preemptive (at time quantum)	Non-preemptive	Preemptive (at arrival)	Non-preemptive	Preemptive (at time quantum)
<b>Throughput</b>	Not emphasized	May be low if quantum is too small	High	High	High	Not emphasized
<b>Response time</b>	May be high, especially if there is a large variance in process execution times	Provides good response time for short processes	Provides good response time for short processes	Provides good response time	Provides good response time	Not emphasized
<b>Overhead</b>	Minimum	Minimum	Can be high	Can be high	Can be high	Can be high
<b>Effect on processes</b>	Penalizes short processes; penalizes I/O bound processes	Fair treatment	Penalizes long processes	Penalizes long processes	Good balance	May favor I/O bound processes
<b>Starvation</b>	No	No	Possible	Possible	No	Possible

# Selection Function

---

- **Determines which process, among ready processes, is selected next for execution**
- **May be based on priority, resource requirements, or the execution characteristics of the process**
- **If based on execution characteristics then important quantities are:**
  - $w$  = time spent in system so far, waiting
  - $e$  = time spent in execution so far
  - $s$  = total service time required by the process, including  $e$ ; generally, this quantity must be estimated or supplied by the user

# Decision Mode

---

- Specifies the instants in time at which the selection function is exercised
- Two categories:
  - Nonpreemptive
    - once a process is in the running state, it will continue until it terminates or blocks itself for I/O
  - Preemptive
    - currently running process may be interrupted and moved to ready state by the OS
    - preemption may occur when new process arrives, on an interrupt, or periodically

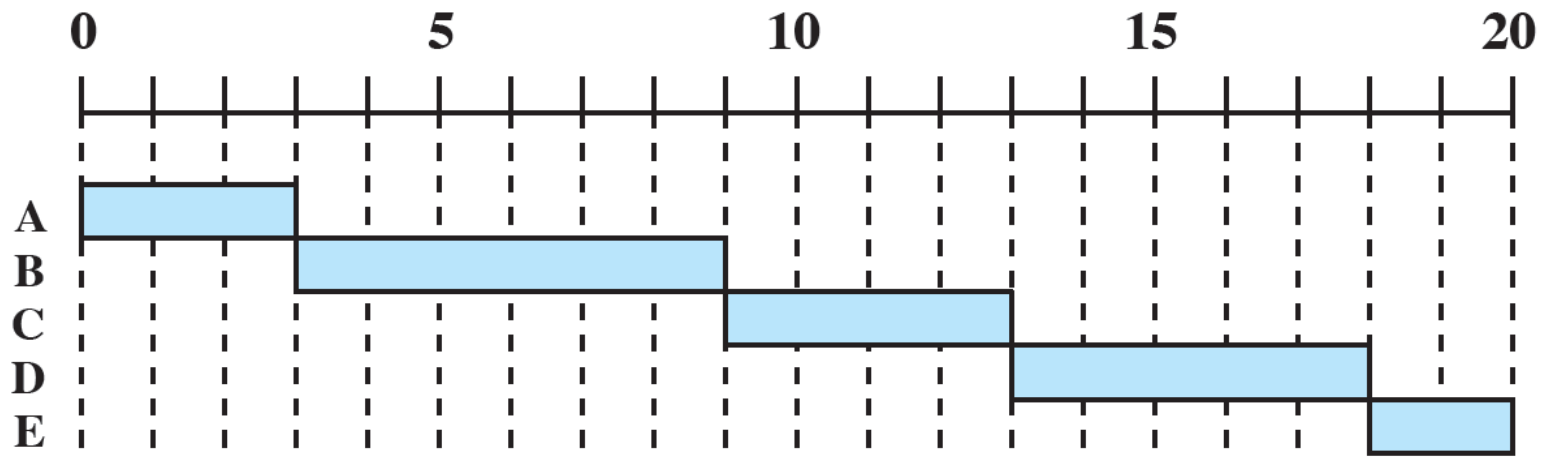
# Process Scheduling Example

---

Process	Arrival Time	Service Time
A	0	3
B	2	6
C	4	4
D	6	5
E	8	2

# First-Come-First-Served (FCFS)

- Also known as first-in-first-out (FIFO) or a strict queuing scheme
- Performs much better for long processes than short ones
- When the current process ceases to execute, the longest process in the Ready queue is selected
- Tends to favor processor-bound processes over I/O-bound processes

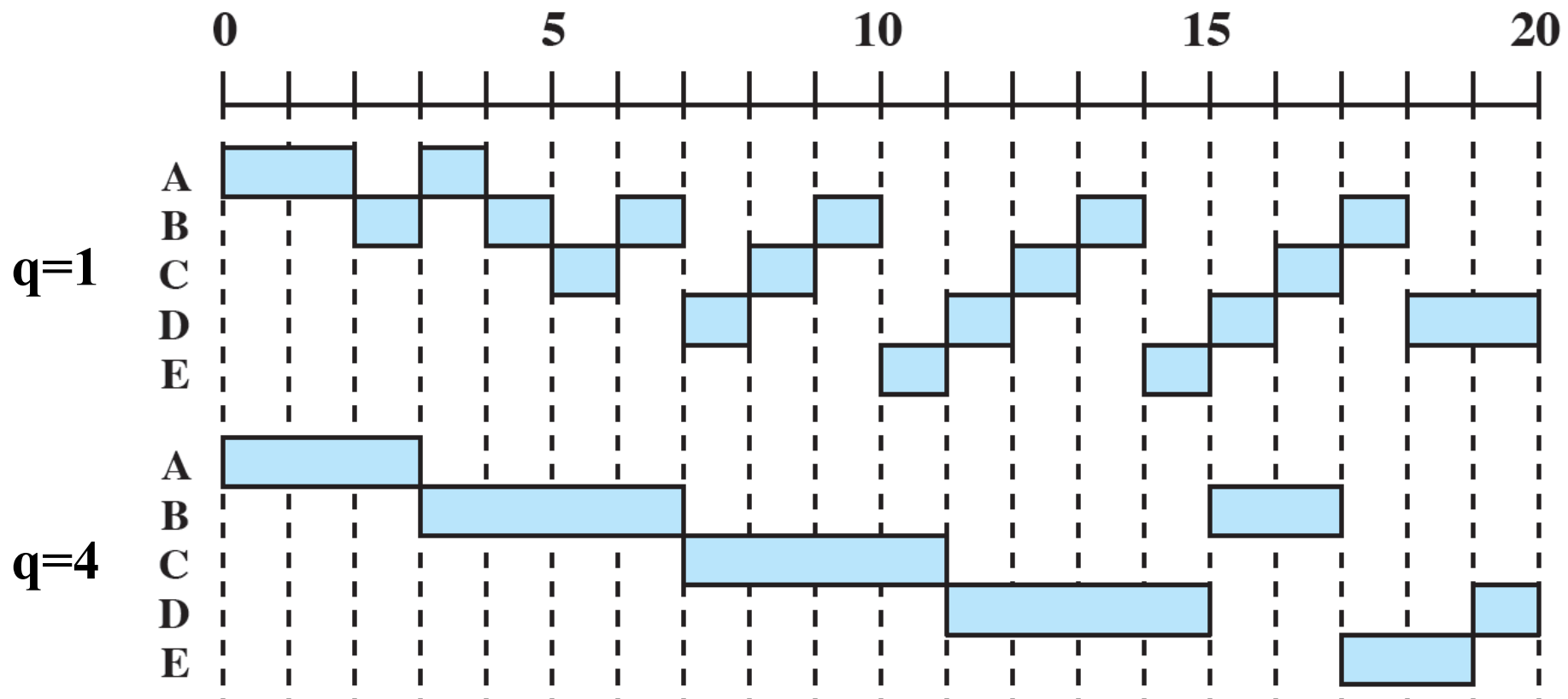


# Round Robin

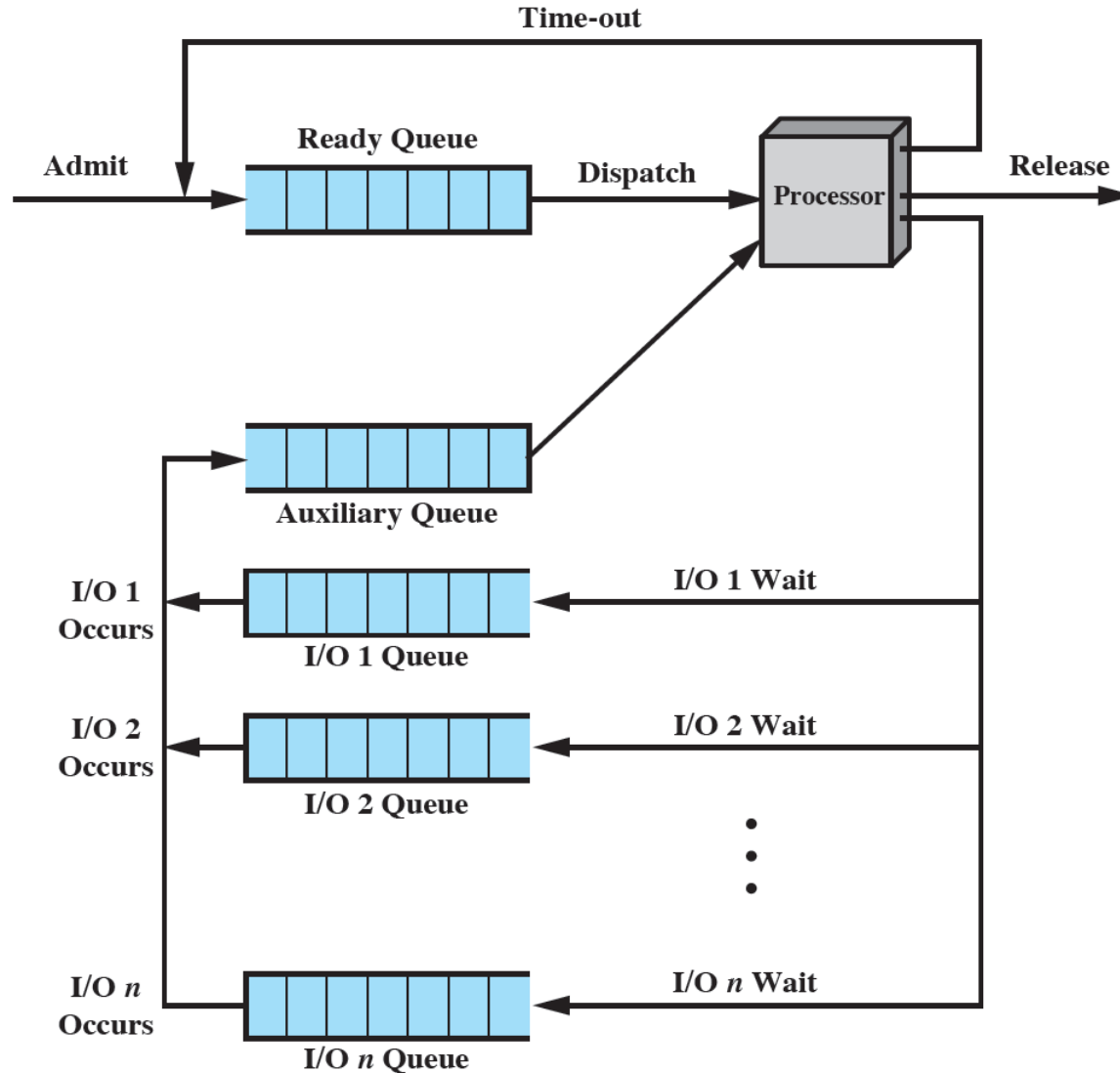
---

- Uses preemption based on a clock
- Also known as time slicing because each process is given a slice of time before being preempted
- Principal design issue is the length of the time quantum, or slice, to be used
- Particularly effective in a general-purpose time-sharing system or transaction processing system
- One drawback is its relative treatment of processor-bound and I/O-bound processes

# Round Robin



# Virtual Round Robin (VRR)





# Shortest Process Next (SPN)

- **Nonpreemptive policy**

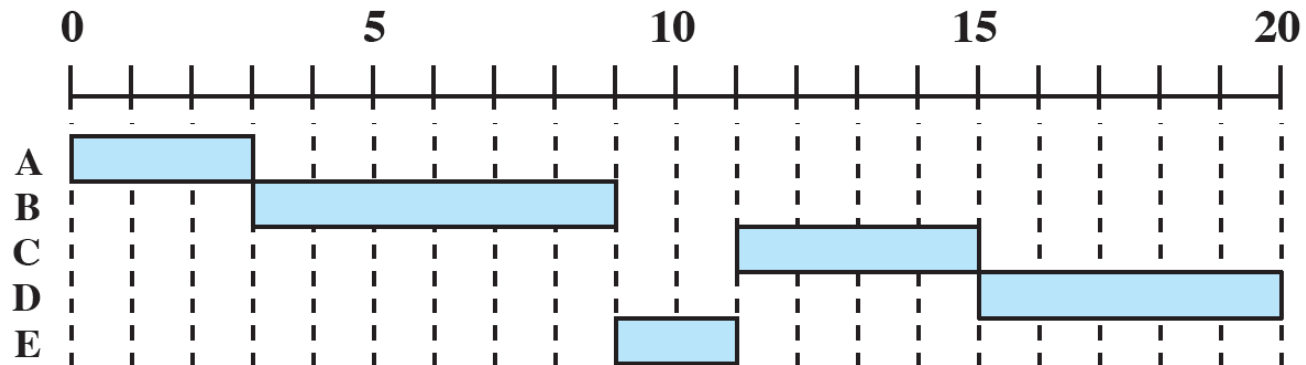
- the process with the shortest expected processing time is selected next

- **A short process will jump to the head of the queue**

- **Possibility of starvation for longer processes**

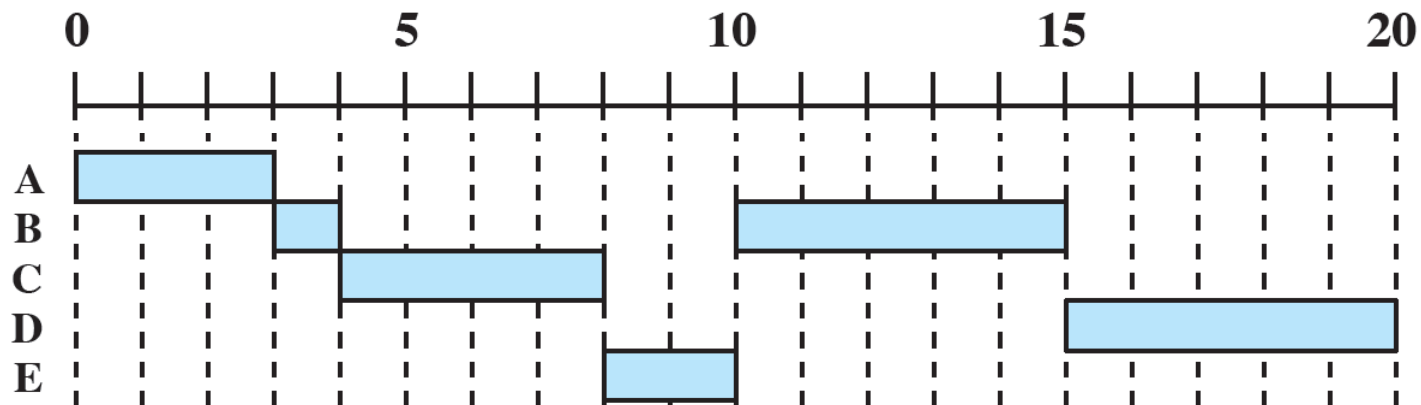
- **One difficulty is the need to know, or at least estimate, the required processing time of each process**

- If the programmer's estimate is substantially under the actual running time, the system may abort the job



# Shortest Remaining Time (SRT)

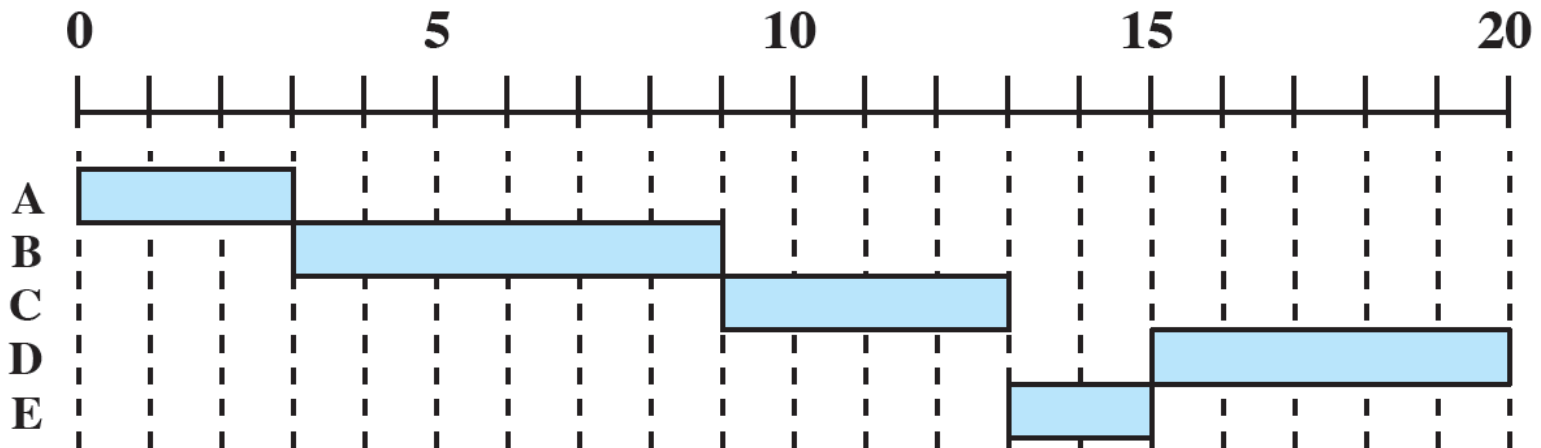
- Preemptive version of SPN
- Scheduler always chooses the process that has the shortest expected remaining processing time
- Risk of starvation of longer processes
- Should give superior turnaround time performance to SPN because a short job is given immediate preference to a running longer job



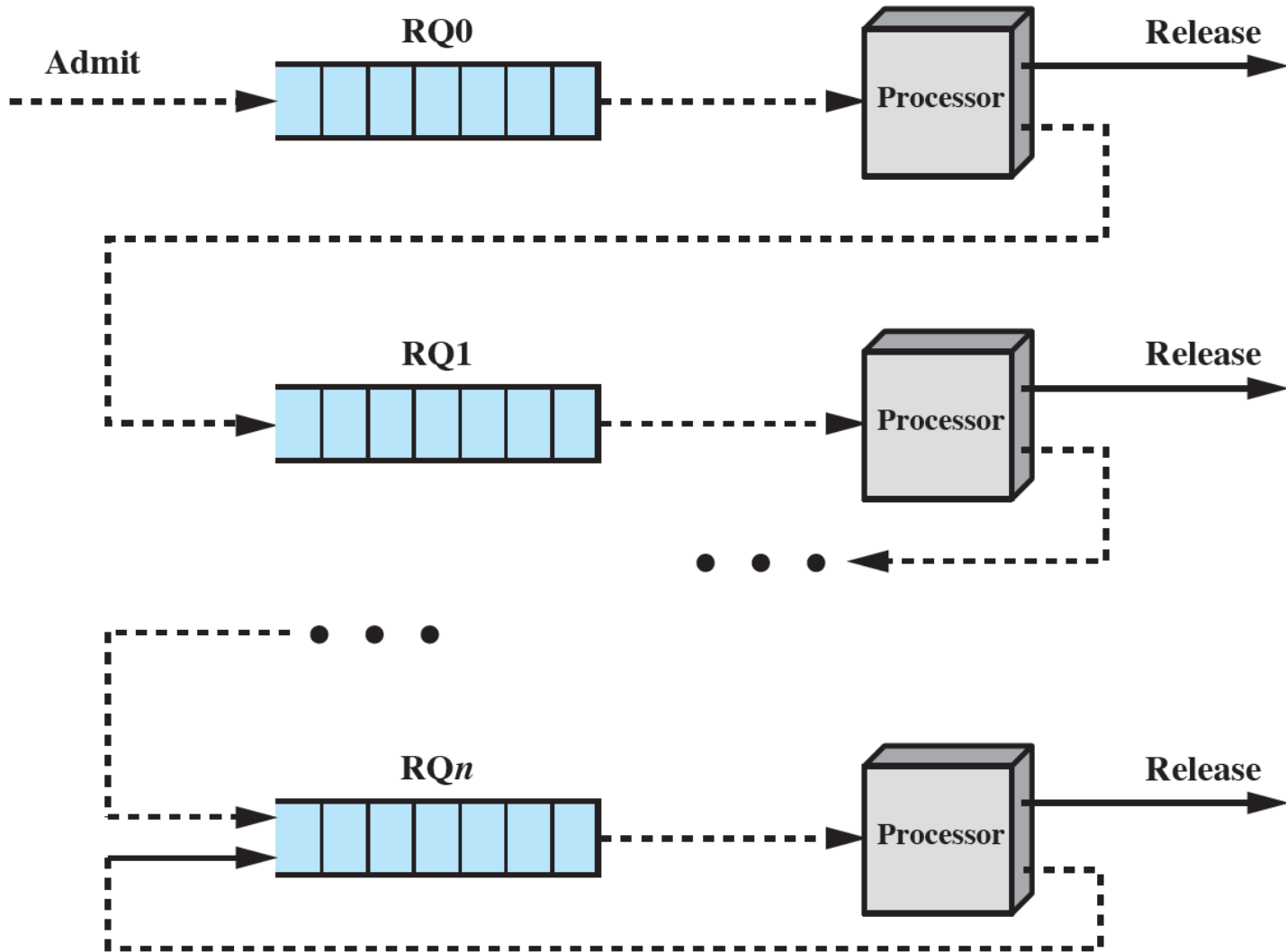
# Highest Response Ratio Next (HRRN)

- Chooses next process with the greatest ratio
- Attractive because it accounts for the age of the process
- While shorter jobs are favored, aging without service increases the ratio so that a longer process will eventually get past competing shorter jobs

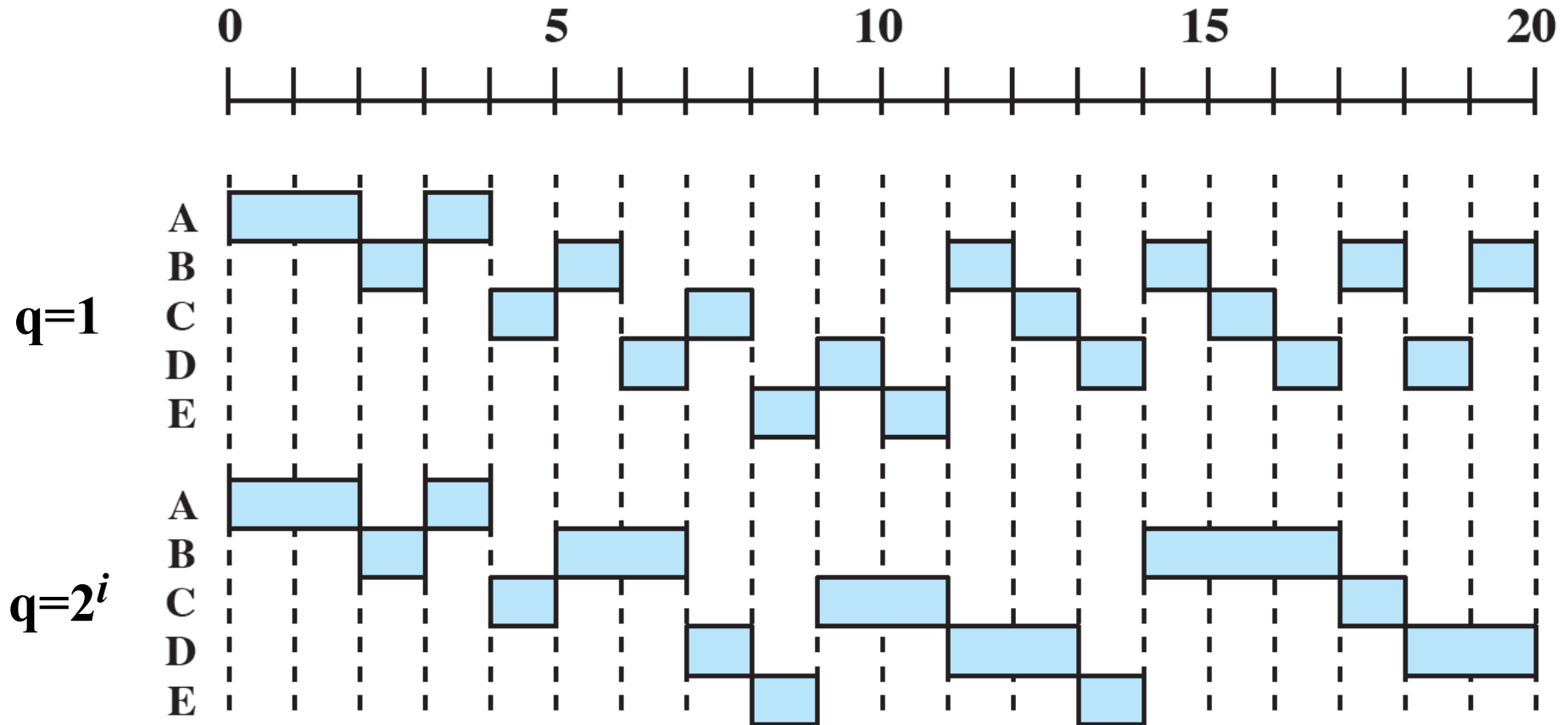
$$\text{Ratio} = \frac{\text{time spent waiting} + \text{expected service time}}{\text{expected service time}}$$



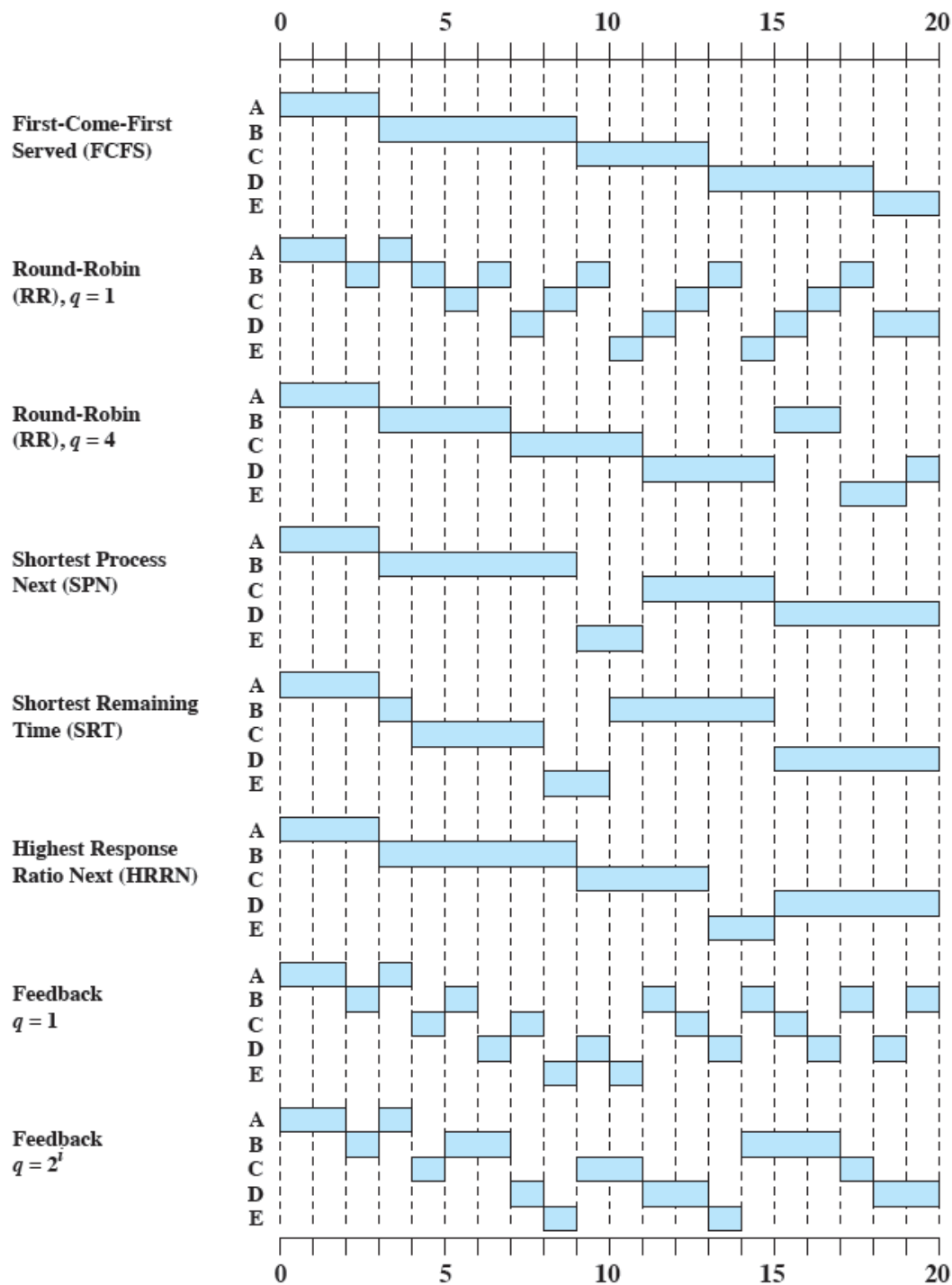
# Feedback Scheduling



# Feedback Performance



# A Comparison of Scheduling Policies



Process	Arrival Time	Service Time
A	0	3
B	2	6
C	4	4
D	6	5
E	8	2

# A Comparison of Scheduling Policies

Process	A	B	C	D	E	
Arrival Time	0	2	4	6	8	
Service Time (Ts)	3	6	4	5	2	Mean
FCFS						
Finish Time	3	9	13	18	20	
Turnaround Time (Tr)	3	7	9	12	12	8.60
Tr/Ts	1.00	1.17	2.25	2.40	6.00	2.56
RR q = 1						
Finish Time	4	18	17	20	15	
Turnaround Time (Tr)	4	16	13	14	7	10.80
Tr/Ts	1.33	2.67	3.25	2.80	3.50	2.71
RR q = 4						
Finish Time	3	17	11	20	19	
Turnaround Time (Tr)	3	15	7	14	11	10.00
Tr/Ts	1.00	2.5	1.75	2.80	5.50	2.71
SPN						
Finish Time	3	9	15	20	11	
Turnaround Time (Tr)	3	7	11	14	3	7.60
Tr/Ts	1.00	1.17	2.75	2.80	1.50	1.84
SRT						
Finish Time	3	15	8	20	10	
Turnaround Time (Tr)	3	13	4	14	2	7.20
Tr/Ts	1.00	2.17	1.00	2.80	1.00	1.59
HRRN						
Finish Time	3	9	13	20	15	
Turnaround Time (Tr)	3	7	9	14	7	8.00
Tr/Ts	1.00	1.17	2.25	2.80	3.5	2.14
FB q = 1						
Finish Time	4	20	16	19	11	
Turnaround Time (Tr)	4	18	12	13	3	10.00
Tr/Ts	1.33	3.00	3.00	2.60	1.5	2.29

# Summary

---

## ■ The OS must make three types of scheduling decisions with respect to the execution of processes

- Long-term – determines when new processes are admitted to the system
- Medium-term – part of the swapping function and determines when a program is brought into main memory so that it may be executed
- Short-term – determines which ready process will be executed next by the processor

## ■ From a user's point of view

- response time is generally the most important characteristic of a system

## ■ From a system point of view

- throughput or processor utilization is important

## ■ Algorithms:

- FCFS
- Round Robin
- SPN
- SRT
- HRRN
- Feedback