

## Chap17 Queueing Theory

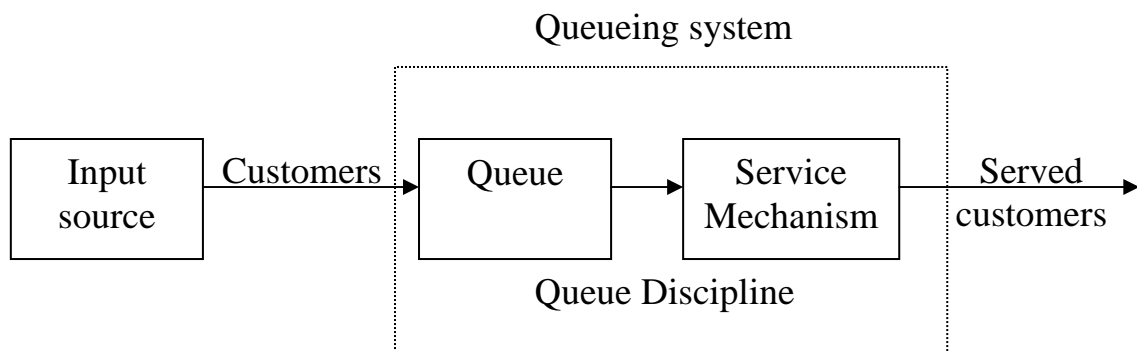
### □ Introduction

- ✓ Queues (waiting line) are a part of everyday life.
- ✓ Providing too much service involves excessive costs. And not providing enough service capacity causes the waiting line to become excessively long.
- ✓ The ultimate goal is to achieve an economic balance between the cost of service and the cost associated with the waiting for that service.
- ✓ Queueing theory is the study of waiting in all these various guises.

### □ Prototype Example—Doctor Requirement in a Emergency Room

- ✓ Consider assigning an extra doctor to the emergency room, which has one doctor already.
- ✓ How much can we reduce the average waiting time for patients if the extra doctor is hired?

### □ Basic Structure of Queueing Models

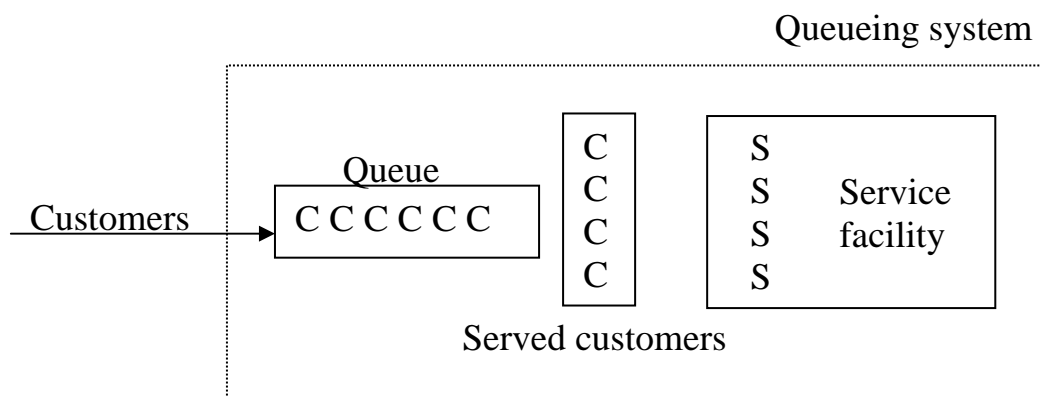


- ✓ Input Source (Calling Population)
  - One characteristic of the input source is its size. The size is the total number of customers. The size may be infinite (default one) or finite.
  - When will each one arrive? Associate with a distribution—usually, Poisson distribution (the number of customers generated until any specific time) or Exponential distribution (**interarrival time**).
  - A customer may be balking, who refuses to enter the system and is lost if the queue is too long.

- ✓ Queue
  - The queue is where customers wait before being served.
  - A queue is characterized by the maximum permissible number of customers that it can contain. Queue may be infinite (default one) or finite.
- ✓ Queue Discipline
  - Refers to the order in which members of the queue are selected for service.
  - First-come-first-serve is normally used.
- ✓ Service Mechanism
  - Consists of one or more service facilities, each of which contains one or more parallel service channels, called servers.
  - At a given facility, the customer enters one of the parallel service channels and is served by that server.
  - Most elementary models assume one service facility with either one or a finite number of servers.
  - Service time is usually defined by a probability distribution.

#### □ An Elementary Queueing Process

- ✓ A single waiting line forms in the front of a single service facility, within which are stationed one or more servers. Each customer is serviced by one of the servers, perhaps after some waiting in the queue.



- ✓ The prototype example is of this type.
- ✓ We usually label a queueing model as  $M/M/1$ 
  - The first spot is for distribution of interarrival times. The second spot is for distribution of service times. The third one is for number of servers.
  - $M$  = exponential distribution (Markovian), which is the most widely used.
  - $D$  = degenerate distribution (constant time).
  - $E_k$  = Erlang distribution.

- $G$  = general distribution (any arbitrary distribution allowed)

## □ Terminology and Notation

- ✓ State of system = number of customers in queueing system.
- ✓ Queue length = number of customers waiting for service to begin = state of system minus number of customers being served.
- ✓  $N(t)$  = number of customers in queueing system at time  $t$ .
- ✓  $P_n(t)$  = probability of exactly  $n$  customers in queueing system at time  $t$ .
- ✓  $s$  = number of servers (parallel service channels) in queueing system.
- ✓  $\lambda_n$  = mean arrival rate (expected number of arrival per unit time) of new customers when  $n$  customers are in system.
  - When  $\lambda_n$  is a constant for all  $n$ , this constant is denoted by  $\lambda$ .
  - $1/\lambda$  is the expected interarrival time.
- ✓  $\mu_n$  = mean service rate for overall system (expected number of customers completing service per unit time) when  $n$  customers are in system.
  - $\mu_n$  represents combined rate at which all busy servers achieve service completions.
  - When the mean service rate *per busy server* is a constant for all  $n \geq 1$ , this constant is denoted by  $\mu$ .
  - $\mu_n = s\mu$  when  $n \geq s$  (all servers are busy).
  - $1/\mu$  is the expected service time.
- ✓  $\rho = \lambda/(s\mu)$  is the **utilization factor** for the service facility, i.e., the expected fraction of time the individual servers are busy.
- ✓ **Transient condition**—when a queueing system has recently begun, the state of the system will be greatly affected by the initial state and by the time that has since elapsed.
- ✓ **Steady-state condition**—after sufficient time has elapsed, the state of the system becomes essentially independent of the initial state and the elapsed time.
  - Queueing theory has tended to focus largely on the steady-state condition.
- ✓ More notations are defined in a steady-state condition.
- ✓  $P_n$  = probability of exact  $n$  customers in queueing system.
- ✓  $L$  = expected number of customers in queueing system =  $\sum_{n=0}^{\infty} nP_n$ .

- ✓  $L_q = \text{expected queue length (excludes customers being served)} = \sum_{n=s}^{\infty} (n-s)P_n$ .
- ✓  $W = \text{waiting time in system (includes service time) for each customer.}$
- ✓  $W = E(W)$ .
- ✓  $W_q = \text{waiting time in queue (exclude service time) for each customer.}$
- ✓  $W_q = E(W_q)$ .

#### □ Relationships between $L$ , $W$ , $L_q$ , and $W_q$

- ✓ Assume that  $\lambda_n$  is a constant for all  $n$ .
- ✓ In a steady-state queueing process,  $L = \lambda W$  (Little's formula) and  $L_q = \lambda W_q$ .
- ✓ If the  $\lambda_n$  are not equal, then  $\lambda$  can be replaced in these equation by  $\bar{\lambda}$ , the average arrival rate over the long time.
- ✓ Assume that the mean service time ( $1/\mu$ ) is a constant. Thus,  $W = W_q + \frac{1}{\mu}$ .
- ✓ These four fundamental quantities ( $L$ ,  $W$ ,  $L_q$ , and  $W_q$ ) could be immediately determined as soon as one is found analytically.

#### □ Examples—commercial service system, transportation service system, internal service system, and social service system.

#### □ The Role of the Exponential Distribution

- ✓ The mostly commonly used distribution for interarrival and service time is exponential distribution.
- ✓ A random variable (interarrival or service times),  $T$ , is said to have an **exponential distribution** with parameter  $\alpha$  if its probability density function is

$$f_T(t) = \begin{cases} \alpha e^{-\alpha t} & \text{for } t \geq 0 \\ 0 & \text{for } t < 0 \end{cases}$$

- ✓ The cumulative probabilities are  $P\{T \leq t\} = 1 - e^{-\alpha t}$ ,  $P\{T > t\} = e^{-\alpha t}$  for  $t \geq 0$ .
- ✓  $E(T) = \frac{1}{\alpha}$ ,  $\text{var}(T) = \frac{1}{\alpha^2}$ .

✓ Property 1:  $f_T(t)$  is a strictly decreasing function of  $t$ .

➤  $P\{0 \leq T \leq \Delta t\} > P\{t \leq T \leq t + \Delta t\}$  for any strictly positive of  $t$  and  $\Delta t$ .

➤ The value  $T$  takes on is more likely to be “small” [less than half of  $E(T)$ ] than “near” its expected value. Is this property real?

$$P\left\{0 \leq T \leq \frac{1}{2} \frac{1}{\alpha}\right\} = 0.393, \quad P\left\{\frac{1}{2} \frac{1}{\alpha} \leq T \leq \frac{3}{2} \frac{1}{\alpha}\right\} = 0.383.$$

➤ Not real when the service required is essentially identical for each customer, with the server always performing the same sequence of service operations.

➤ It is suitable for the situations where the specific tasks required of the server differ among customers (hospital or banking cases).

✓ Property 2: Lack of memory:  $P\{T > t + \Delta t \mid T > \Delta t\} = P\{T > t\}$  for any positive of  $t$  and  $\Delta t$ .

➤ The probability distribution of the remaining time until the event occurs always is the same, regardless of how much time already has passed.

➤ The process “forgets” its history.

➤ The phenomenon occurs with the exponential distribution.

➤ For interarrival time, the time until next arrival is completely uninfluenced by when the last arrival occurred.

✓ Property 3: The minimum of several independent exponential random variables has an exponential distribution.

➤ Let  $T_1, T_2, \dots, T_n$  be independent exponential random variables with parameters  $\alpha_1, \alpha_2, \dots, \alpha_n$ . Also, let  $U$  be the random variable that takes on the value equal to the minimum of the values of  $T_1, T_2, \dots, T_n$ .

- If  $T_i$  represents the time until a particular event occurs, then  $U$  represents the time until the first of the  $n$  different events occurs.
  - $U$  indeed has an exponential distribution with parameter  $\alpha = \sum_{i=1}^n \alpha_i$ .
  - If there are  $n$  different types of customers (interarrival time is exponential with parameter  $\alpha_i$ ), the interarrival time for the queueing system as a whole, has an exponential distribution with parameter  $\alpha = \sum_{i=1}^n \alpha_i$ .
  - Suppose all  $n$  servers have the same exponential service-time distribution with parameter  $\mu$ . The time until the next service completion from any server has an exponential distribution with parameter  $\alpha = n\mu$ .
- ✓ Property 4: Relationship to the Poisson distribution.
- Suppose that the time between consecutive occurrences of some particular kind of event has an exponential distribution with parameter  $\alpha$ .
  - Then, the number of occurrence by time  $t$  ( $X(t)$ ) has a Poisson distribution with parameter  $\alpha t$ .
- $$P\{X(t) = n\} = \frac{(\alpha t)^n e^{-\alpha t}}{n!}$$
- With  $n = 0$ ,  $P\{X(t) = 0\} = e^{-\alpha t}$ , which is just the probability from the exponential distribution that the first event occurs after time  $t$ .
  - The mean of this Poisson distribution is  $E\{X(t)\} = \alpha t$ , so that the expected number of events per unit time is  $\alpha$ .  $\alpha$  is said to be the mean rate at which the events occurs.
  - When the events are counted on a continuing basis, the counting process  $\{X(t); t \geq 0\}$  is said to be a **Poisson process** with parameter  $\alpha$ .
  - Define  $X(t)$  as the number of service completions achieved by a continuously busy server in elapsed time  $t$ , where  $\alpha = \mu$ . For multiple-server queueing models,  $X(t)$  can also be defined as the number of service completions achieved by  $n$  servers in elapsed time  $t$ , where  $\alpha = n\mu$ .
  - Suppose the interarrival times have an exponential distribution with parameter  $\lambda$ . In this case,  $X(t)$  is the number of arrivals in elapsed time  $t$ , where  $\alpha = \lambda$  is the mean arrival rate. Therefore, arrivals occur according to a **Poisson input process** with parameter  $\lambda$ .

- ✓ Property 5: For all positive values of  $t$ ,  $P\{T \leq t + \Delta t \mid T > t\} \approx \alpha \Delta t$ , for small  $\Delta t$ .
  - The series expansion of  $e^x$  for any exponent  $x$  is  $e^x = 1 + x + \sum_{n=2}^{\infty} \frac{x^n}{n!}$ .
  - $P\{T \leq t + \Delta t \mid T > t\} = P\{T \leq \Delta t\} = 1 - e^{-\alpha \Delta t} = 1 - 1 + \alpha \Delta t - \sum_{n=2}^{\infty} \frac{(-\alpha \Delta t)^n}{n!} \approx \alpha \Delta t$ .
- ✓ Property 6: Unaffected by aggregation or disaggregation.
  - If there are  $n$  different types of customers (each is a Poisson input process with parameter  $\lambda_i$ ), the aggregated input is a Poisson with  $\lambda = \lambda_1 + \lambda_2 + \dots + \lambda_n$ .
  - Assuming that each arriving customer has a fixed probability  $p_i$  of being of type  $i$ , with  $\lambda_i = p_i \lambda$  and  $\sum_{i=1}^n p_i = 1$ , the property says that the input process for customers of type  $i$  also must be Poisson with parameter  $\lambda_i$ .

### □ The Birth-and-Death Process

- ✓ **birth** = arrival; **death** = departure; state,  $N(t)$ , is the number of customers in the queueing system at time  $t$ .
  - ✓ The birth-and-death process describes probabilistically how  $N(t)$  changes as  $t$  increases.
  - ✓ Assumption 1: Given  $N(t) = n$ , the current probability distribution of the remaining time until the next birth (arrival) is exponential with parameter  $\lambda_n$ .
  - ✓ Assumption 2: Given  $N(t) = n$ , the current probability distribution of the remaining time until the next death (service completion) is exponential with parameter  $\mu_n$ .
  - ✓ Assumption 3: The random variable of assumption 1 and the random variable of assumption 2 are mutually independent. The next transition in the state of the process is either  $n \rightarrow n+1$  or  $n \rightarrow n-1$  depending on whether the former or latter random variable is smaller.
  - ✓ That is, the birth-and-death process can be illustrated by the rate diagram.
- 
- ✓ The following analysis only focuses on the steady state condition.
  - ✓  $E_n(t)$  = number of times that process enters state  $n$  by time  $t$ .
  - ✓  $L_n(t)$  = number of times that process leaves state  $n$  by time  $t$ .

- ✓  $|E_n(t) - L_n(t)| \leq 1$ . Dividing both sides by  $t$  and letting  $t \rightarrow \infty$ .

$$\left| \frac{E_n(t)}{t} - \frac{L_n(t)}{t} \right| \leq \frac{1}{t}, \text{ so } \lim_{t \rightarrow \infty} \left| \frac{E_n(t)}{t} - \frac{L_n(t)}{t} \right| = 0.$$

- ✓  $\lim_{n \rightarrow \infty} \frac{E_n(t)}{t} =$  mean rate at which process enters state  $n$ .

- ✓  $\lim_{n \rightarrow \infty} \frac{L_n(t)}{t} =$  mean rate at which process leaves state  $n$ .

- ✓ **Rate In = Rate Out Principle:** for any state of the system  $n$ , mean entering rate = mean leaving rate. The equation expressing this principle is called the **balance equation** for state  $n$ .

- ✓  $P_i$  is the steady-state probability of being in state  $i$ .

- ✓ Consider state 0: the mean entering rate of state 0 is  $\mu_1 P_1$ ; the mean leaving rate of state is  $\lambda_0 P_0$ . The balance equation for state 0 is  $\mu_1 P_1 = \lambda_0 P_0$

- ✓ For every other state there are two possible transitions both into and out of the state. Therefore, each side of the balance equations for these states represents the sum of the mean rates for the two transitions involved.

- ✓ We can write the balance equations for the other states.

- ✓ Notice that there is always one “extra” variable in these equations. Solve  $P_1, P_2 \dots$  in term of  $P_0$ .



- ✓ Thus, the steady-state probabilities are

$$P_n = C_n P_0, \text{ where } C_n = \frac{\lambda_{n-1} \lambda_{n-2} \dots \lambda_0}{\mu_n \mu_{n-1} \dots \mu_1} \text{ for } n = 1, 2, \dots; \text{ when } n = 0, C_n = 1.$$

- ✓ Then, use the sum of all probabilities equal 1 to solve for  $P_0$ .

$$\sum_{n=0}^{\infty} P_n = 1 \text{ implies that } \left( \sum_{n=0}^{\infty} C_n \right) P_0 = 1, \text{ so that } P_0 = \left( \sum_{n=0}^{\infty} C_n \right)^{-1}.$$

- ✓ The key measures of performance for the queueing system ( $L$ ,  $L_q$ ,  $W$ , and  $W_q$ ) can be obtained immediately after calculating the  $P_n$ .

$$L = \sum_{n=0}^{\infty} n P_n, \quad L_q = \sum_{n=s}^{\infty} (n-s) P_n,$$

$$W = \frac{L}{\lambda}, \quad W_q = \frac{L_q}{\lambda},$$

where  $\bar{\lambda}$  is the average arrival rate over the long run and  $\bar{\lambda} = \sum_{n=0}^{\infty} \lambda_n P_n$ .

- ✓ The above calculations are based on the steady-state condition. Steady-state condition holds if  $\lambda_n = 0$  or  $\rho = \lambda / (s\mu) < 1$ .

#### □ The $M/M/s$ model

- ✓ All interarrival and service times are independently and identically distributed according to an exponential distribution. The number of servers is  $s$ .
- ✓ This model is a special case of the birth-and-death process where the queueing system's mean arrival rate and mean service rate per busy server are constant.
- ✓ When the system has just a single server ( $s = 1$ ), the parameters for the birth-and-death process are  $\lambda_n = \lambda$ , and  $\mu_n = \mu$ . Rate diagram is as follow.

- ✓ When the system has multiple server ( $s > 1$ ),
  - $\mu_n$  represents the mean service rate for the overall queueing system when there are  $n$  customers in the system.

- The service rate per busy server is  $\mu$ , the overall mean service rate for  $n$  busy servers must be  $n\mu$ .
- Therefore,  $\mu_n = n\mu$  when  $n \leq s$ , and  $\mu_n = s\mu$  when  $n \geq s$ .

- ✓ When the maximum mean service rate  $s\mu$  exceeds the mean arrival rate  $\lambda$ , that is, when  $\rho = \frac{\lambda}{s\mu} < 1$ , a queueing system will eventually reach a steady-state condition. We can use the results derived in the birth-and-death model.

#### □ Results for the Single-Server Case (M/M/1)

- ✓ The  $C_n$  factors reduce to  $C_n = \left(\frac{\lambda}{\mu}\right)^n = \rho^n$ , for  $n = 0, 1, 2, \dots$
- ✓ Therefore,  $P_n = \rho^n P_0$ , for  $n = 0, 1, 2, \dots$ , where  $P_0 = \left(\sum_{n=0}^{\infty} \rho^n\right)^{-1} = \left(\frac{1}{1-\rho}\right)^{-1} = 1 - \rho$ .
- ✓ Thus,  $P_n = (1 - \rho)\rho^n$ , for  $n = 0, 1, 2, \dots$
- ✓ Consequently,  $L = \sum_{n=0}^{\infty} n(1 - \rho)\rho^n = \frac{\lambda}{\mu - \lambda}$

$$\checkmark L_q = \sum_{n=1}^{\infty} (n-1)P_n = \frac{\lambda^2}{\mu(\mu - \lambda)}$$

- ✓ When  $\lambda \geq \mu$  the mean arrival rate exceeds the mean service rate, the preceding solution “blows up”.

- ✓ Consider the case when  $\lambda < \mu$  and the queue discipline is first-come-first-served. We can derive the probability distribution of the waiting time in the system  $W$  for a random arrival.
  - ✓ If this arrival finds  $n$  customers already in the system, then the arrival will have to wait through  $n+1$  exponential service times, including his/her own.
  - ✓ Let  $T_1, T_2, \dots$  be independent service-time random variables having an exponential distribution with parameter  $\mu$ , and let  $S_{n+1} = T_1 + T_2 + \dots + T_{n+1}$ .
  - ✓  $P\{W > t\} = \sum_{n=0}^{\infty} P_n P\{S_{n+1} > t\} = e^{-\mu(1-\rho)t}$ . That is,  $W$  has an exponential distribution with parameter  $\mu(1-\rho)$ .
  - ✓ Therefore,  $W = E(W) = \frac{1}{\mu(1-\rho)} = \frac{1}{\mu - \lambda}$ .
  - ✓ Sometimes, we are concern about  $W_q$ , the waiting time in the queue.
  - ✓ If this arrival finds no customers already in the system, there is no waiting time in queue.  $P\{W_q = 0\} = P_0 = 1 - \rho$ .
  - ✓ If this arrival finds  $n > 0$  customers already in the system, then the arrival has to wait through  $n$  exponential service times until his/her own service begins.
  - ✓  $P\{W_q > t\} = \sum_{n=1}^{\infty} P_n P\{S_n > t\} = \rho e^{-\mu(1-\rho)t}$ .
- 
- ✓  $W_q$  does not quite have an exponential distribution, because  $P\{W_q = 0\} > 0$ .
  - ✓ The conditional distribution of  $W_q$ , given that  $W_q > 0$ , does have an exponential distribution with parameter  $\mu(1-\rho)$ , because
 
$$P\{W_q > t \mid W_q > 0\} = \frac{P\{W_q > t\}}{P\{W_q > 0\}} = e^{-\mu(1-\rho)t}.$$
  - ✓ By deriving the mean of the (unconditional) distribution of  $W_q$  (or applying either  $L_q = \lambda W_q$  or  $W_q = W - \frac{1}{\mu}$ ),  $W_q = E(W_q) = \frac{\lambda}{\mu(\mu - \lambda)}$ .

□ **Results for the Multiple-Server Case (M/M/s) —  $s > 1$**

✓ For  $n = 1, 2, \dots, s$ ,  $C_n = \frac{(\lambda/\mu)^n}{n!}$ . For  $n = s, s+1, \dots$ ,  $C_n = \frac{(\lambda/\mu)^n}{s!s^{n-s}}$ .

✓  $P_0 =$

$$= 1 / \left[ \sum_{n=0}^{s-1} \frac{(\lambda/\mu)^n}{n!} + \frac{(\lambda/\mu)^s}{s!} \frac{1}{1 - \lambda/(s\mu)} \right].$$

✓ For  $0 \leq n \leq s$ ,  $P_n = \frac{(\lambda/\mu)^n}{n!} P_0$ . For  $n \geq s$ ,  $P_n = \frac{(\lambda/\mu)^n}{s!s^{n-s}} P_0$ .

✓  $L_q = \frac{P_0(\lambda/\mu)^s \rho}{s!(1-\rho)^2}$

✓  $W_q = L_q / \lambda$ ,  $W = W_q + \frac{1}{\mu}$ ,  $L = \lambda(W_q + \frac{1}{\mu}) = L_q + \frac{\lambda}{\mu}$

$$\checkmark P\{W > t\} = e^{-\mu} \left[ 1 + \frac{P_0 (\lambda / \mu)^s}{s!(1-\rho)} \left( \frac{1 - e^{-\mu(s-1-\lambda/\mu)}}{s-1-\lambda/\mu} \right) \right]$$

$$\checkmark P\{W_q > t\} = (1 - P\{W_q = 0\})e^{-s\mu(1-\rho)t}, \text{ where } P\{W_q = 0\} = \sum_{n=0}^{s-1} P_n.$$

### □ The County Hospital Example with the $M/M/s$ Model

✓ Estimate that patients will arrive at an average rate of 1 every 1/2 hour. A doctor requires an average of 20 minutes to treat each patient. That is,  $\lambda = 2$  customers per hour,  $\mu = 3$  customers per hour.

#### ✓ Steady-State Results

	$s = 1$	$s = 2$
$\rho$	2/3	1/3
$P_0$	1/3	1/2
$P_1$	2/9	1/3
$P_n$ (for $n \geq 2$ )	$(1/3)(2/3)^n$	$(1/3)^n$
$L_q$	4/3	1/12
$L$	2	3/4
$W_q$	2/3	1/24 (hour)
$W$	1	3/8 (hour)
$P\{W_q > 0\}$	0.667	0.167
$P\{W_q > 1/2\}$	0.404	0.022
$P\{W_q > 1\}$	0.245	0.003
$P\{W_q > t\}$	$2/3e^{-t}$	$1/6e^{-4t}$
$P\{W > t\}$	$e^{-t}$	$1/2e^{-3t}(3-e^{-t})$

### □ The Finite Queue Variation of the $M/M/s$ Model (Called the $M/M/s/K$ Model)

✓ The queueing systems sometimes have a finite queue. That is, the number of customers in the system is not permitted to exceed some specified number (denoted by  $K$ ). Thus, the queue capacity is  $K - s$ .

✓ Any customer that arrives while the queue is “full” is refused entry into the system. That is, the mean input rate becomes zero at these times.

➤ For  $n = 0, 1, 2, \dots, K-1$ ,  $\lambda_n = \lambda$ ;

For  $n \geq K$ ,  $\lambda_n = 0$ .

✓ A queueing system that fits this model will eventually reach a steady-state condition.

### □ Results for the Single-Server Case ( $M/M/1/K$ )

✓ For  $n = 0, 1, 2, \dots, K$ ,  $C_n = \left(\frac{\lambda}{\mu}\right)^n = \rho^n$ ; for  $n > K$ ,  $C_n = 0$ .

✓ Therefore, for  $\rho \neq 1$ ,  $P_0 = \frac{1-\rho}{1-\rho^{K+1}}$ , so that  $P_n = \frac{1-\rho}{1-\rho^{K+1}} \rho^n$ , for  $n = 0, 1, 2, \dots, K$ .

✓ Thence,  $L = \frac{\rho}{1-\rho} - \frac{(K+1)\rho^{K+1}}{1-\rho^{K+1}}$ ,  $L_q = L - (1 - P_0)$ .

✓ Notice that the preceding results do not require that  $\lambda < \mu$ .

✓ When  $\rho < 1$ , the second term in the expression for  $L$  converges to 0 as  $K \rightarrow \infty$ , so that all the preceding results converge to the corresponding results given for the  $M/M/1$  model.

✓  $W = \frac{L}{\lambda}$ ,  $W_q = \frac{L_q}{\lambda}$ , where  $\bar{\lambda} = \sum_{n=0}^{\infty} \lambda_n P_n = \lambda(1 - P_K)$

## □ Results for the Multiple-Server Case ( $M/M/s/K$ ) – $s > 1$

$$\checkmark \text{ For } n = 1, 2, \dots, s, C_n = \frac{(\lambda / \mu)^n}{n!};$$

$$\text{For } n = s, s+1, \dots, K, C_n = \frac{(\lambda / \mu)^n}{s!s^{n-s}};$$

$$\text{For } n > K, C_n = 0.$$

$$\checkmark \text{ Hence, For } n = 1, 2, \dots, s, P_n = \frac{(\lambda / \mu)^n}{n!} P_0;$$

$$\text{For } n = s, s+1, \dots, K, P_n = \frac{(\lambda / \mu)^n}{s!s^{n-s}} P_0;$$

$$\text{For } n > K, P_n = 0.$$

$$\text{where } P_0 = 1 / \left[ \sum_{n=0}^s \frac{(\lambda / \mu)^n}{n!} + \frac{(\lambda / \mu)^s}{s!} \sum_{n=s+1}^K \left( \frac{\lambda}{s\mu} \right)^{n-s} \right].$$

$$\checkmark L_q = \frac{P_0 (\lambda / \mu)^s \rho}{s!(1-\rho)^2} [1 - \rho^{K-s} - (K-s)\rho^{K-s}(1-\rho)], \text{ where } \rho = \lambda / (s\mu).$$

$$\checkmark L = \sum_{n=0}^{s-1} nP_n + L_q + s(1 - \sum_{n=0}^{s-1} P_n).$$

$\checkmark$   $W$  and  $W_q$  are obtained from these quantities as shown for the single-server case.

### □ The Finite Calling Population Variation of the $M/M/s$ Model

- $\checkmark$  The input source is limited; i.e., the size of the calling population is finite.
- $\checkmark$  Let  $N$  denote the size of the calling population.
- $\checkmark$  When the number of customers in the queueing system is  $n$ , there are only  $N - n$  potential customers remaining in the input source.
- $\checkmark$  The most important application of this model has been to the machine repair problem.
  - $\triangleright$  There are  $N$  machines. The maintenance people are servers.
- $\checkmark$  Each member of the calling population alternates between inside and outside the queueing system.
- $\checkmark$  Assumes that each member's outside time has an exponential distribution with parameter  $\lambda$ .

- ✓ When  $n$  members are inside, and so  $N - n$  members are outside, the probability distribution of the remaining time until the next arrival is the distribution of the minimum of the remaining outside times for the latter  $N - n$  members.
  - This distribution is exponential with parameter  $\lambda_n = (N - n)\lambda$ .
  - The rate diagram is as followings.

### □ Results for the Single-Server Case ( $s=1$ ) for the Finite Call Population

✓ For  $n \leq N$ , 
$$C_n = N(N-1)\dots(N-n+1)\left(\frac{\lambda}{\mu}\right)^n = \frac{N!}{(N-n)!}\left(\frac{\lambda}{\mu}\right)^n$$

For  $n > N$ , 
$$C_n = 0$$

✓ Therefore, 
$$P_0 = 1 / \sum_{n=0}^N \left[ \frac{N!}{(N-n)!} \left(\frac{\lambda}{\mu}\right)^n \right]$$

$$P_n = \frac{N!}{(N-n)!} \left(\frac{\lambda}{\mu}\right)^n P_0, \text{ if } n = 1, 2, \dots, N$$

$$L_q = \sum_{n=1}^N (n-1)P_n = N - \frac{\lambda + \mu}{\lambda} (1 - P_0)$$

$$L = \sum_{n=0}^N nP_n = L_q + 1 - P_0 = N - \frac{\mu}{\lambda} (1 - P_0)$$

$$W = \frac{L}{\lambda} \text{ and } W_q = \frac{L_q}{\lambda}$$



$$\text{where } \bar{\lambda} = \sum_{n=0}^{\infty} \lambda_n P_n = \sum_{n=0}^N (N-n)\lambda P_n = \lambda(N-L)$$

### □ Results for the Multiple-Server Case ( $s > 1$ ) for the Finite Call Population

$$\checkmark \text{ For } n = 0, 1, 2, \dots, s, \quad C_n = \frac{N!}{(N-n)!n!} \left(\frac{\lambda}{\mu}\right)^n.$$

$$\text{For } n = s, s+1, \dots, N, \quad C_n = \frac{N!}{(N-n)!s!s^{n-s}} \left(\frac{\lambda}{\mu}\right)^n.$$

$$\text{For } n > N, \quad C_n = 0.$$

$$\checkmark \text{ If } 0 \leq n \leq s, \quad P_n = \frac{N!}{(N-n)!n!} \left(\frac{\lambda}{\mu}\right)^n P_0$$

$$\text{If } s \leq n \leq N, \quad P_n = \frac{N!}{(N-n)!s!s^{n-s}} \left(\frac{\lambda}{\mu}\right)^n P_0$$

$$\text{If } n > N, \quad P_n = 0$$

$$\checkmark P_0 = 1 / \left[ \sum_{n=0}^{s-1} \frac{N!}{(N-n)!n!} \left(\frac{\lambda}{\mu}\right)^n + \sum_{n=s}^N \frac{N!}{(N-n)!s!s^{n-s}} \left(\frac{\lambda}{\mu}\right)^n \right]$$

$$\checkmark L_q = \sum_{n=s}^N (n-s)P_n$$

$$\checkmark L = \sum_{n=0}^{s-1} nP_n + L_q + s \left( 1 - \sum_{n=0}^{s-1} P_n \right)$$

### □ Queueing Models Involving Nonexponential Distributions – M/G/1

- ✓ No restrictions are imposed on what the service-time distribution can be.
- ✓ It is only necessary to know the mean  $1/\mu$  and variance  $\sigma^2$  of this distribution.
- ✓ Thanks for the Pollaczek-Khintchine formula, we have (when  $\rho < 1$ )

$$P_0 = 1 - \rho$$

$$L_q = \frac{\lambda^2 \sigma^2 + \rho^2}{2(1 - \rho)}$$

$$L = \rho + L_q$$

$$W_q = \frac{L_q}{\lambda}$$

$$W = W_q + \frac{1}{\mu}$$

- ✓ For any fixed expected service time  $1/\mu$ , notice that  $L_q$ ,  $L$ ,  $W_q$ , and  $W$  all increase as  $\sigma^2$  is increased.
- ✓ When the service-time distribution is exponential,  $\sigma^2 = 1/\mu^2$ , and the preceding results will reduce to the corresponding results for the  $M/M/1$  model.

#### □ Queueing Models Involving Nonexponential Distributions – $M/D/s$

- ✓ When the service consists of essentially the same routine task for all customers, there tends to be little variation in the service time required.
- ✓ It assumes that all service times equal some fixed constant (the degenerate service-time distribution).
- ✓ When there is just a single server, the  $M/D/1$  model is just the special case of the  $M/G/1$  model where  $\sigma^2 = 0$ .

➤ The Pollaczek-Khintchine formula reduces to  $L_q = \frac{\rho^2}{2(1 - \rho)}$ .

#### □ Queueing Models Involving Nonexponential Distributions – $M/E_k/s$

- ✓  $M/D/s$  model assumes zero variation and the exponential distribution assumes a very large variation ( $\sigma = 1/\mu$ ).
- ✓ Between these two rather extreme cases lies another distribution – the **Erlang Distribution**.
- ✓ The PDF is  $f(t) = \frac{(\mu k)^k}{(k-1)!} t^{k-1} e^{-k\mu t}$ , where  $\mu$  and  $k$  are strictly positive parameters and  $k$  needs to be integer. Mean =  $\frac{1}{\mu}$ , and standard deviation =  $\frac{1}{\sqrt{k}} \frac{1}{\mu}$ .
- ✓ Suppose that  $T_1, T_2, \dots, T_k$  are  $k$  independent random variables with an identical exponential distribution whose mean is  $1/(k\mu)$ . Then their sum  $T = T_1 + T_2 + \dots + T_k$  has an Erlang distribution with parameters  $\mu$  and  $k$ .

- ✓ The Erlang distribution is a large (two-parameter) family of distributions permitting only nonnegative values.
  - Both the exponential and the degenerate distribution are special cases of the Erlang distribution with  $k = 1$  and  $k = \infty$ , respectively.
- ✓ Consider  $M/Ek/1$  model, which is just the special case of the  $M/G/1$  model. Applying Pollaczek-Khintchine formula with  $\sigma^2 = 1/(k\mu^2)$ , we have

$$L_q = \frac{\lambda^2 / (k\mu^2) + \rho^2}{2(1-\rho)} = \frac{1+k}{2k} \frac{\lambda^2}{\mu(\mu-\lambda)}$$

$$W_q = \frac{1+k}{2k} \frac{\lambda}{\mu(\mu-\lambda)}$$

$$W = W_q + \frac{1}{\mu}$$

$$L = \lambda W$$

### □ Priority-Discipline Queueing Models

- ✓ The queue discipline is based on a priority system.
- ✓ With **nonpreemptive priorities**, a customer being served cannot be ejected back into the queue (preempted) if a higher-priority customer enters.
- ✓ With **preemptive priorities**, the lowest-priority customer being served is preempted whenever a higher-priority customer enters.

### □ Results for the Nonpreemptive Priorities Model

- ✓ Let  $W_k$  be the steady-state expected waiting time in the system (including service time) for a member of priority class  $k$ .

$$W_k = \frac{1}{AB_{k-1}B_k} + \frac{1}{\mu}, \text{ for } k = 1, 2, \dots, N,$$

$$\text{where } A = s! \frac{s\mu - \lambda}{r^s} \sum_{j=0}^{s-1} \frac{r^j}{j!} + s\mu,$$

$$B_0 = 1, \quad B_k = 1 - \frac{\sum_{i=1}^k \lambda_i}{s\mu},$$

$s$  = number of servers

$\mu$  = mean service rate per busy server,

$\lambda_i$  = mean arrival rate for priority class  $i$ ,

$$\lambda = \sum_{i=0}^N \lambda_i, \quad r = \frac{\lambda}{\mu}.$$

- ✓ Little's formula still applies to individual priority class, so  $L_k = \lambda_k W_k$ .

### □ A Single-Server Variation of the Nonpreemptive Priority Model

- ✓ Different priority classes have different expected service time.
- ✓ Let  $1/\mu_k$  denote the mean of the exponential service-time distribution for priority class  $k$ , so  $\mu_k =$  mean service rate for priority class  $k$ .

$$\checkmark \quad W_k = \frac{a_k}{b_{k-1}b_k} + \frac{1}{\mu_k}, \text{ where } a_k = \sum_{i=1}^k \frac{\lambda_i}{\mu_i^2}, \quad b_0 = 1, \quad b_k = 1 - \sum_{i=1}^k \frac{\lambda_i}{\mu_i}.$$

### □ Results for the Preemptive Priority Model

- ✓ Assume the expected service time is the same for all priority classes.
- ✓  $W_k = \frac{1/\mu}{B_{k-1}B_k}$ , for the single-server cases.
- ✓ When  $s > 1$ ,  $W_k$  can be calculated by an iterative procedure that will soon be illustrated by an example.
- ✓  $L_k = \lambda_k W_k$ .

### □ The County Hospital Example with Priorities

- ✓ There are three patient categories: (1) critical (10%), (2) serious (30%), and (3) stable (60%).
- ✓ Give  $\lambda = 2$  and  $\mu = 3$ , we have  $\lambda_1 = 0.2$ ,  $\lambda_2 = 0.6$ , and  $\lambda_3 = 1.2$ .

	Preemptive Priorities		Nonpreemptive Priorities	
	$s = 1$	$s = 2$	$s = 1$	$s = 2$
$A$	--	--	4.5	36
$B_1$	0.933	--	0.933	0.967
$B_2$	0.733	--	0.733	0.867
$B_3$	0.333	--	0.333	0.667
$W_1 - 1/\mu$	0.024 hr	0.00037 hr	0.238 hr	0.029 hr
$W_2 - 1/\mu$	0.154 hr	0.00793 hr	0.325 hr	0.033 hr
$W_3 - 1/\mu$	1.033 hr	0.06542 hr	0.889 hr	0.048 hr

- ✓ The waiting time for priority class 1 customers are unaffected by the presence of lower-priority classes customers. Thus,  $W_1$  must equal  $W$  for the corresponding one-class  $M/M/1$  model with  $s = 2$ ,  $\mu = 3$ , and  $\lambda = 0.2$ , which yield  $W_1 = 0.3337$  hour.
  - Waiting time in the queue for class 1 customers is
 
$$W_1 - 1/\mu = 0.33370 - 0.33333 = 0.00037$$
- ✓ Now consider the first two priority classes. Let  $\bar{W}_{1-2}$  be the expected waiting time in the system of a random arrival in either of these two classes.
  - The probability is  $\lambda_1 / (\lambda_1 + \lambda_2) = 1/4$  that this arrival is in class 1 and  $\lambda_2 / (\lambda_1 + \lambda_2) = 3/4$  that it is in class 2.
 
$$\bar{W}_{1-2} = \frac{1}{4}W_1 + \frac{3}{4}W_2.$$
  - $\bar{W}_{1-2}$  must equal  $W$  for the  $M/M/s$  model with  $s = 2$ ,  $\mu = 3$ , and  $\lambda = 0.8$ , which yields  $\bar{W}_{1-2} = W = 0.33937$ .
  - We already know the value of  $W_1$ , so  $W_2 = 0.34126$ .  $W_2 - 1/\mu = 0.00793$ .
- ✓ Similarly,  $\bar{W}_{1-3} = 0.1W_1 + 0.3W_2 + 0.6W_3$ .  $\bar{W}_{1-3}$  equals  $W$  for the  $M/M/s$  model with  $s = 2$ ,  $\mu = 3$ , and  $\lambda = 2$ , which yields  $\bar{W}_{1-3} = W = 0.375$ . So,  $W_3 = 0.39875$ .  $W_3 - 1/\mu = 0.06542$ .

## □ Queueing Network

- ✓ Thus far we have considered only queueing systems that have a single service facility with one or more servers.
- ✓ Queueing systems are sometimes actually queueing networks, i.e., networks of service facilities where customers must receive service at some of or all these facilities.
  - For example, orders being processed through a job shop must be routed through a sequence of machine groups.

## □ Equivalence Property

- ✓ Assume that a service facility with  $s$  servers and an infinite queue has a Poisson input with parameter  $\lambda$  and the same exponential service-time distribution with parameter  $\mu$  for each server (the  $M/M/s$  model), where  $s\mu > \lambda$ . Then the steady-state output of this service facility is also a Poisson process with parameter  $\lambda$ .

## □ Jackson Networks (with $m$ facilities, $i = 1, 2, \dots, m$ )

- ✓ An infinite queue.
- ✓ Customers arriving from outside the system according to a Poisson input process with parameter  $a_i$ .

- ✓  $s_i$  servers with an exponential service-time distribution with parameter  $\mu_i$ .
- ✓ A customer leaving facility  $i$  is routed next to facility  $j$  ( $j = 1, 2, \dots, m$ ) with probability  $p_{ij}$  or departs the system with probability  $q_i = 1 - \sum_{j=1}^m p_{ij}$ .

### □ Key properties of Jackson Network

- ✓ Under steady-state conditions, each facility  $j$  ( $j = 1, 2, \dots, m$ ) in a Jackson network behaves as if it were an independent  $M/M/s$  system with arrival rate  $\lambda_j = a_j + \sum_{i=1}^m \lambda_i p_{ij}$ , where  $s_j \mu_j > \lambda_j$ .
- ✓ Consider a Jackson network with three service facilities that have the parameter shown below.

Facility $j$	$s_j$	$\mu_j$	$a_j$	$P_{ij}$		
				$i = 1$	$i = 2$	$i = 3$
$j = 1$	1	10	1	0	0.1	0.4
$j = 2$	2	10	4	0.6	0	0.4
$j = 3$	1	10	3	0.3	0.3	0

We have

$$\lambda_1 = 1 + 0.1 \lambda_2 + 0.4 \lambda_3$$

$$\lambda_2 = 4 + 0.6 \lambda_1 + 0.4 \lambda_3$$

$$\lambda_3 = 3 + 0.3 \lambda_1 + 0.3 \lambda_2$$

The simultaneous solution for this system is  $\lambda_1 = 5$ ,  $\lambda_2 = 10$ ,  $\lambda_3 = 7.5$ .

- ✓ Each of the three facilities now can be analyzed independently by using the formulas for the  $M/M/s$  model given before.

$$\rho_i = \frac{\lambda_i}{s_i \mu_i} \quad (\text{That is, } \rho_1 = 1/2, \rho_2 = 1/2, \rho_3 = 3/4).$$

$$P_{n1} =$$

$$P_{n2} =$$

$$P_{n3} =$$

The join probability of  $(n_1, n_2, n_3) = P_{n1}P_{n2}P_{n3}$

$$L_1 =$$

$$L_2 =$$

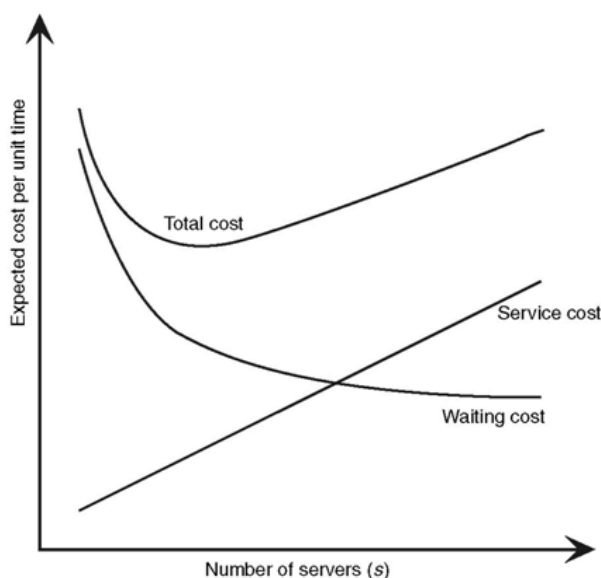
$$L_3 =$$

$$L =$$

$$W = L/\lambda = L/(a_1 + a_2 + a_3) =$$

### □ The Application of Queueing Theory – How Many Servers Should be Provided?

- ✓ The two primary considerations in making these decisions are (1) the cost of the service capacity and (2) the waiting cost of customers.
- ✓  $E(\text{TC})$  = expected total cost per unit time.  
 $E(\text{SC})$  = expected service cost per unit time.  
 $E(\text{WC})$  = expected waiting cost per unit time.
- ✓ When each server costs the same, the service cost is  $E(\text{SC}) = C_s s$ , where  $C_s$  is the marginal cost of a server per unit time.
- ✓ To evaluate WC for any value of  $s$ ,  $E(\text{WC}) = C_w L$ , where  $C_w$  is the waiting cost per unit time for each customer.
- ✓ Therefore, after estimating the constants,  $C_s$  and  $C_w$ , the goal is to choose the value of  $s$  so as to Minimize  $E(\text{TC}) = C_s s + C_w L$ .



### □ An Example for Determining the Number of Servers

- ✓ For a  $M/M/s$  model with  $\lambda = 120$  customers per hour and  $\mu = 80$  customers per hour.
- ✓ We need at least two servers to reach the steady-state, since  $\frac{120}{2(80)} < 1$ .
- ✓ Each server costs \$20 per hour and waiting cost is \$48 per hour.
- ✓ We want to minimize  $20s + 48L$ .
- ✓ Use the results obtained before, we have

$s$	$L$	$E(SC) = C_s s$	$E(WC) = C_w L$	$E(TC) = E(SC) + E(WC)$
2	3.43	\$40	\$164.57	\$204.57
3	1.74	\$60	\$83.37	\$143.37
4	1.54	\$80	\$74.15	\$154.15
5	1.51	\$100	\$72.41	\$172.41